
Pimlico Documentation

Release 1.0rc

Mark Granroth-Wilding

Oct 25, 2018

Contents

1 Contents	3
Python Module Index	207

The **Pimlico Processing Toolkit** is a toolkit for building pipelines of tasks for **processing large datasets** (corpora). It is especially focussed on processing linguistic corpora and provides wrappers around many existing, widely used **NLP** (Natural Language Processing) tools.

Note: These are the docs for the **release candidate for v1.0**.

This brings with it a big project to change how datatypes work internally (previously in branch `datatypes`) and requires all datatypes and modules to be updated to the new system. [More info...](#)

Modules marked with `!!` in the docs are waiting to be updated and don't work. Other known outstanding tasks are marked with todos: [full todo list](#).

These issues will be resolved before v1.0 is released.

It makes it easy to write large, potentially complex pipelines with the following key goals:

- to provide **clear documentation** of what has been done;
- to make it easy to **incorporate standard NLP tasks**,
- and to extend the code with **non-standard tasks, specific to a pipeline**;
- to support simple **distribution of code** for reproduction, for example, on other datasets.

The toolkit takes care of managing data between the steps of a pipeline and checking that everything's executed in the right order.

The core toolkit is written in Python. Pimlico is open source, released under the GPLv3 license. It is available from [its Github repository](#). To get started with a Pimlico project, follow the [getting-started guide](#).

Pimlico is short for *Pipelined Modular Linguistic COrpus processing*.

More NLP tools will gradually be added. See [my wishlist](#) for current plans.

1.1 Pimlico guides

Step-by-step guides through common tasks while using Pimlico.

1.1.1 Super-quick Pimlico setup

This is a very quick walk-through of the process of starting a new project using Pimlico. For more details, explanations, etc see *the longer getting-started guide*.

First, make sure Python is installed.

System-wide configuration

Choose a location on your file system where Pimlico will store all the output from pipeline modules. For example, `/home/me/.pimlico_store/`.

Create a file in your home directory called `.pimlico` that looks like this:

```
store=/home/me/.pimlico_store
```

This is not specific to a pipeline: separate pipelines use separate subdirectories.

Set up new project

Create a new, empty directory to put your project in. E.g.:

```
cd ~
mkdir myproject
```

Download `newproject.py` into this directory and run it:

```
wget https://raw.githubusercontent.com/markgw/pimlico/master/admin/newproject.py
python newproject.py myproject
```

This fetches the latest Pimlico codebase (in `pimlico/`) and creates a template pipeline (`myproject.conf`).

Customizing the pipeline

You've got a basic pipeline config file now (`myproject.conf`).

Add sections to it to configure modules that make up your pipeline.

For guides to doing that, see the *the longer setup guide* and individual module documentation.

Running Pimlico

Check the pipeline can be loaded and take a look at the list of modules you've configured:

```
./pimlico.sh myproject.conf status
```

Tell the modules to fetch all the dependencies you need:

```
./pimlico.sh myproject.conf install all
```

If there's anything that can't be installed automatically, this should output instructions for manual installation.

Check the pipeline's ready to run a module that you want to run:

```
./pimlico.sh myproject.conf run MODULE --dry-run
```

To run the next unexecuted module in the list, use:

```
./pimlico.sh myproject.conf run
```

1.1.2 Setting up a new project using Pimlico

Todo: Setup guide has a lot that needs to be updated for the new datatypes system. I've updated up to **Getting input**.

You've decided to use Pimlico to implement a data processing pipeline. So, where do you start?

This guide steps through the basic setup of your project. You don't have to do everything exactly as suggested here, but it's a good starting point and follows Pimlico's recommended procedures. It steps through the setup for a very basic pipeline.

A shorter version of this guide that zooms through the essential setup steps is also available.

System-wide configuration

Pimlico needs you to specify certain parameters regarding your local system. Typically this is just a file in your home directory called `.pimlico`. *More details*.

It needs to know where to put output files as it executes. These settings apply to all Pimlico pipelines you run. Pimlico will make sure that different pipelines don't interfere with each other's output (provided you give them different names).

Most of the time, you only need to specify one storage location, using the `store` parameter in your local config file. (You can specify multiple: [more details](#)).

Create a file `~/ .pimlico` that looks like this:

```
store=/path/to/storage/directory
```

All pipelines will use different subdirectories of this one.

Getting started with Pimlico

The procedure for starting a new Pimlico project, using the latest release, is very simple.

Create a new, empty directory to put your project in. Download [newproject.py](#) into the project directory.

Choose a name for your project (e.g. `myproject`) and run:

```
python newproject.py myproject
```

This fetches the latest version of Pimlico (now in the `pimlico/` subdirectory) and creates a basic config file, which will define your pipeline.

It also retrieves libraries that Pimlico needs to run. Other libraries required by specific pipeline modules will be installed as necessary when you use the modules.

Building the pipeline

You've now got a config file in `myproject.conf`. This already includes a pipeline section, which gives the basic pipeline setup. It will look something like this:

```
[pipeline]
name=myproject
release=<release number>
python_path=%(project_root)s/src/python
```

The name needs to be distinct from any other pipelines that you run – it's what distinguishes the storage locations.

`release` is the release of Pimlico that you're using: it's automatically set to the latest one, which has been downloaded.

If you later try running the same pipeline with an updated version of Pimlico, it will work fine as long as it's the same major version (the first digit). Otherwise, there may be backwards incompatible changes, so you'd need to update your config file, ensuring it plays nicely with the later Pimlico version.

Getting input

Now we add our first module to the pipeline. This reads input from a collection of text files. We use a small subset of the [Europarl corpus](#) as an example here. This can be simply adapted to reading the real Europarl corpus or any other corpus stored in this straightforward way.

[Download and extract the small corpus from here](#)

In the example below, we have extracted the files to a directory `data/europarl_demo` in the home directory.

```
[input-text]
type=pimlico.modules.input.text.raw_text_files
files=%(home)s/data/europarl_demo/*
```

Todo: Continue writing from here

Doing something: tokenization

Now, some actual linguistic processing, albeit somewhat uninteresting. Many NLP tools assume that their input has been divided into sentences and tokenized. The OpenNLP-based tokenization module does both of these things at once, calling OpenNLP tools.

Notice that the output from the previous module feeds into the input for this one, which we specify simply by naming the module.

```
[tokenize]
type=pimlico.modules.opennlp.tokenize
input=tar-grouper
```

Doing something more interesting: POS tagging

Many NLP tools rely on part-of-speech (POS) tagging. Again, we use OpenNLP, and a standard Pimlico module wraps the OpenNLP tool.

```
[pos-tag]
type=pimlico.modules.opennlp.pos
input=tokenize
```

Running Pimlico

Now we've got our basic config file ready to go. It's a simple linear pipeline that goes like this:

```
read input docs -> group into batches -> tokenize -> POS tag
```

Before we can run it, there's one thing missing: three of these modules have their own dependencies, so we need to get hold of the libraries they use. The input reader uses the Beautiful Soup python library and the tokenization and POS tagging modules use OpenNLP.

Checking everything's dandy

Now you can run the `status` command to check that the pipeline can be loaded and see the list of modules.

```
./pimlico.sh myproject.conf status
```

To check that specific modules are ready to run, with all software dependencies installed, use the `run` command with `--dry-run` (or `--dry`) switch:

```
./pimlico.sh myproject.conf run tokenize --dry
```

With any luck, all the checks will be successful. There might be some missing software dependencies.

Fetching dependencies

All the standard modules provide easy ways to get hold of their dependencies automatically, or as close as possible. Most of the time, all you need to do is tell Pimlico to install them.

Use the `run` command, with a module name and `--dry-run`, to check whether a module is ready to run.

```
./pimlico.sh myproject.conf run tokenize --dry
```

In this case, it will tell you that some libraries are missing, but they can be installed automatically. Simply issue the `install` command for the module.

```
./pimlico.sh myproject.conf install tokenize
```

Simple as that.

There's one more thing to do: the tools we're using require statistical models. We can simply download the pre-trained English models from the OpenNLP website.

At present, Pimlico doesn't yet provide a built-in way for the modules to do this, as it does with software libraries, but it does include a GNU Makefile to make it easy to do:

```
cd ~/myproject/pimlico/models
make opennlp
```

Note that the modules we're using default to these standard, pre-trained models, which you're now in a position to use. However, if you want to use different models, e.g. for other languages or domains, you can specify them using extra options in the module definition in your config file.

If there are any other library problems shown up by the dry run, you'll need to address them before going any further.

Running the pipeline

What modules to run?

Pimlico suggests an order in which to run your modules. In our case, this is pretty obvious, seeing as our pipeline is entirely linear – it's clear which ones need to be run before others.

```
./pimlico.sh myproject.conf status
```

The output also tells you the current status of each module. At the moment, all the modules are `UNEXECUTED`.

You'll notice that the `tar-grouper` module doesn't feature in the list. This is because it's a filter – it's run on the fly while reading output from the previous module (i.e. the input), so doesn't have anything to run itself.

You might be surprised to see that `input-text` *does* feature in the list. This is because, although it just reads the data out of a corpus on disk, there's not quite enough information in the corpus, so we need to run the module to collect a little bit of metadata from an initial pass over the corpus. Some input types need this, others not. In this case, all we're lacking is a count of the total number of documents in the corpus.

Note: To make running your pipeline even simpler, you can abbreviate the command by using a **shebang** in the config file. Add a line at the top of `myproject.conf` like this:

```
#!/pimlico.sh
```

Then make the conf file executable by running (on Linux):

```
chmod ug+x myproject.conf
```

Now you can run Pimlico for your pipeline by using the config file as an executable command:

```
./myproject.conf status
```

Running the modules

The modules can be run using the `run` command and specifying the module by name. We do this manually for each module.

```
./pimlico.sh myproject.conf run input-text
./pimlico.sh myproject.conf run tokenize
./pimlico.sh myproject.conf run pos-tag
```

Adding custom modules

Most likely, for your project you need to do some processing not covered by the built-in Pimlico modules. At this point, you can start implementing your own modules, which you can distribute along with the config file so that people can replicate what you did.

The `newproject.py` script has already created a directory where our custom source code will live: `src/python`, with some subdirectories according to the standard code layout, with module types and datatypes in separate packages.

The template pipeline also already has an option `python_path` pointing to this directory, so that Pimlico knows where to find your code. Note that the code's in a subdirectory of that containing the pipeline config and we specify the custom code path relative to the config file, so it's easy to distribute the two together.

Now you can create Python modules or packages in `src/python`, following the same conventions as the built-in modules and overriding the standard base classes, as they do. The following articles tell you more about how to do this:

- [Writing Pimlico modules](#)
- [Writing document map modules](#)
- [Pimlico module structure](#)

Your custom modules and datatypes can then simply be used in the config file as module types.

1.1.3 Running a pipeline

This guide takes you through what to do if you have received someone else's code for a Pimlico project and would like to run it.

This guide is written for Unix/Mac users. You'll need to make some adjustments if using another OS.

What you've got

Hopefully got at least a pipeline config file. This will have the extension `.conf`. In the examples below, we'll use the name `myproject.conf`.

You've probably got a whole directory, with some subdirectories, containing this config file (or even several) together with other related files – datasets, code, etc. This top-level directory is what we'll refer to as the *project root*.

The project may include some code, probably defining some custom Pimlico module types and datatypes. If all is well, you won't need to delve into this, as its location will be given in the config file and Pimlico will take care of the rest.

Getting Pimlico

You hopefully didn't receive the whole Pimlico codebase together with the pipeline and code. It's recommended not to distribute Pimlico, as it can be fetched automatically for a given pipeline.

You'll need Python installed.

Download the [Pimlico bootstrap script](#) from here and put it in the project root.

Now run it:

```
python bootstrap.py myproject.conf
```

The bootstrap script will look in the config file to work out what version of Pimlico to use and then download it.

If this works, you should now be able to run Pimlico.

Using the bleeding edge code

By default, the bootstrap script will fetch a release of Pimlico that the config file declares as being that which it was built with.

If you want the very latest version of Pimlico, with all the dangers that entails and with the caveat that it might not work with the pipeline you're trying to run, you can tell the bootstrap script to checkout Pimlico from its Git repository.

```
python bootstrap.py --git myproject.conf
```

Running Pimlico

Perhaps the project root contains a (link to a) script called `pimlico.sh`.

If not, create one like this:

```
ln -s pimlico/bin/pimlico.sh .
```

Now run `pimlico.sh` with the config file as an argument, issuing the `command_status` command to see the contents of the pipeline:

```
./pimlico.sh myproject.conf status
```

Pimlico will now run and set itself up, before proceeding with your command and showing the pipeline status. This might take a bit of time. It will install a Python virtual environment and some basic packages needed for it to run.

1.1.4 Writing Pimlico modules

Pimlico comes with a fairly large number of *module types* that you can use to run many standard NLP, data processing and ML tools over your datasets.

For some projects, this is all you need to do. However, often you'll want to mix standard tools with your own code, for example, using the output from the tools. And, of course, there are many more tools you might want to run that aren't built into Pimlico: you can still benefit from Pimlico's framework for data handling, config files and so on.

For a detailed description of the structure of a Pimlico module, see *Pimlico module structure*. This guide takes you through building a simple module.

Note: In any case where a module will process a corpus one document at a time, you should write a *document map module*, which takes care of a lot of things for you, so you only need to say what to do with each document.

Todo: Module writing guide needs to be updated for new datatypes.

In particular, the executor example and datatypes in the module definition need to be updated.

Code layout

If you've followed the *basic project setup guide*, you'll have a project with a directory structure like this:

```
myproject/
  pipeline.conf
  pimlico/
    bin/
    lib/
    src/
    ...
  src/
    python/
```

If you've not already created the `src/python` directory, do that now.

This is where your custom Python code will live. You can put all of your custom module types and datatypes in there and use them in the same way as you use the Pimlico core modules and datatypes.

Add this option to the `[pipeline]` section of your config file, so Pimlico knows where to find your code:

```
python_path=src/python
```

To follow the conventions used in Pimlico's codebase, we'll create the following package structure in `src/python`:

```
src/python/myproject/
  __init__.py
  modules/
    __init__.py
  datatypes/
    __init__.py
```

Write a module

A Pimlico module consists of a Python package with a special layout. Every module has a file `info.py`. This contains the definition of the module's metadata: its inputs, outputs, options, etc.

Most modules also have a file `execute.py`, which defines the routine that's called when it's run. You should take care when writing `info.py` not to import any non-standard Python libraries or have any time-consuming operations that get run when it gets imported.

`execute.py`, on the other hand, will only get imported when the module is to be run, after dependency checks.

For the example below, let's assume we're writing a module called `nmf` and create the following directory structure for it:

```
src/python/myproject/modules/
  __init__.py
  nmf/
    __init__.py
    info.py
    execute.py
```

Easy start

To help you get started, Pimlico provides a wizard in the `newmodule` command.

This will ask you a series of questions, guiding you through the most common tasks in creating a new module. At the end, it will generate a template to get you started with your module's code. You then just need to fill in the gaps and write the code for what the module actually does.

Read on to learn more about the structure of modules, including things not covered by the wizard.

Metadata

Module metadata (everything apart from what happens when it's actually run) is defined in `info.py` as a class called `ModuleInfo`.

Here's a sample basic `ModuleInfo`, which we'll step through. (It's based on the Scikit-learn `matrix_factorization` module.)

```
from pimlico.core.dependencies.python import PythonPackageOnPip
from pimlico.core.modules.base import BaseModuleInfo
from pimlico.datatypes.arrays import ScipySparseMatrix, NumpyArray

class ModuleInfo(BaseModuleInfo):
    module_type_name = "nmf"
    module_readable_name = "Sklearn non-negative matrix factorization"
    module_inputs = [("matrix", ScipySparseMatrix)]
    module_outputs = [("w", NumpyArray), ("h", NumpyArray)]
    module_options = {
        "components": {
            "help": "Number of components to use for hidden representation",
            "type": int,
            "default": 200,
        },
    }

    def get_software_dependencies(self):
        return super(ModuleInfo, self).get_software_dependencies() + \
            [PythonPackageOnPip("sklearn", "Scikit-learn")]
```

The `ModuleInfo` should always be a subclass of `BaseModuleInfo`. There are some subclasses that you might want to use instead (e.g., see [Writing document map modules](#)), but here we just use the basic one.

Certain class-level attributes should pretty much always be overridden:

- `module_type_name`: A name used to identify the module internally

- `module_readable_name`: A human-readable short description of the module
- `module_inputs`: Most modules need to take input from another module (though not all)
- `module_outputs`: Describes the outputs that the module will produce, which may then be used as inputs to another module

Inputs are given as pairs `(name, type)`, where `name` is a short name to identify the input and `type` is the datatype that the input is expected to have. Here, and most commonly, this is a subclass of `PimlicoDatatype` and Pimlico will check that a dataset supplied for this input is either of this type, or has a type that is a subclass of this.

Here we take just a single input: a sparse matrix.

Outputs are given in a similar way. It is up to the module's executor (see below) to ensure that these outputs get written, but here we describe the datatypes that will be produced, so that we can use them as input to other modules.

Here we produce two Numpy arrays, the factorization of the input matrix.

Dependencies: Since we require Scikit-learn to execute this module, we override `get_software_dependencies()` to specify this. As Scikit-learn is available through Pip, this is very easy: all we need to do is specify the Pip package name. Pimlico will check that Scikit-learn is installed before executing the module and, if not, allow it to be installed automatically.

Finally, we also define some **options**. The values for these can be specified in the pipeline config file. When the `ModuleInfo` is instantiated, the processed options will be available in its `options` attribute. So, for example, we can get the number of components (specified in the config file, or the default of 200) using `info.options["components"]`.

Executor

Here is a sample executor for the module info given above, placed in the file `execute.py`.

```
from pimlico.core.modules.base import BaseModuleExecutor
from pimlico.datatypes.arrays import NumpyArrayWriter
from sklearn.decomposition import NMF

class ModuleExecutor(BaseModuleExecutor):
    def execute(self):
        input_matrix = self.info.get_input("matrix").array
        self.log.info("Loaded input matrix: %s" % str(input_matrix.shape))

        # Convert input matrix to CSR
        input_matrix = input_matrix.tocsr()
        # Initialize the transformation
        components = self.info.options["components"]
        self.log.info("Initializing NMF with %d components" % components)
        nmf = NMF(components)

        # Apply transformation to the matrix
        self.log.info("Fitting NMF transformation on input matrix" % transform_type)
        transformed_matrix = transformer.fit_transform(input_matrix)

        self.log.info("Fitting complete: storing H and W matrices")
        # Use built-in Numpy array writers to output results in an appropriate format
        with NumpyArrayWriter(self.info.get_absolute_output_dir("w")) as w_writer:
            w_writer.set_array(transformed_matrix)
        with NumpyArrayWriter(self.info.get_absolute_output_dir("h")) as h_writer:
            h_writer.set_array(transformer.components_)
```

The executor is always defined as a class in `execute.py` called `ModuleExecutor`. It should always be a subclass of `BaseModuleExecutor` (though, again, note that there are more specific subclasses and class factories that we might want to use in other circumstances).

The `execute()` method defines what happens when the module is executed.

The instance of the module's `ModuleInfo`, complete with **options** from the pipeline config, is available as `self.info`. A standard Python **logger** is also available, as `self.log`, and should be used to keep the user updated on what's going on.

Getting hold of the **input data** is done through the module info's `get_input()` method. In the case of a Scipy matrix, here, it just provides us with the matrix as an attribute.

Then we do whatever our module is designed to do. At the end, we write the output data to the appropriate output directory. This should always be obtained using the `get_absolute_output_dir()` method of the module info, since Pimlico takes care of the exact location for you.

Most Pimlico datatypes provide a corresponding **writer**, ensuring that the output is written in the correct format for it to be read by the datatype's reader. When we leave the `with` block, in which we give the writer the data it needs, this output is written to disk.

Pipeline config

Our module is now ready to use and we can refer to it in a pipeline config file. We'll assume we've prepared a suitable Scipy sparse matrix earlier in the pipeline, available as the default output of a module called `matrix`. Then we can add section like this to use our new module:

```
[matrix]
...(Produces sparse matrix output)...

[factorize]
type=myproject.modules.nmf
components=300
input=matrix
```

Note that, since there's only one input, we don't need to give its name. If we had defined multiple inputs, we'd need to specify this one as `input_matrix=matrix`.

You can now run the module as part of your pipeline in the usual ways.

Skeleton new module

To make developing a new module a little quicker, here's a skeleton module info and executor.

```
from pimlico.core.modules.base import BaseModuleInfo

class ModuleInfo(BaseModuleInfo):
    module_type_name = "NAME"
    module_readable_name = "READABLE NAME"
    module_inputs = [("NAME", REQUIRED_TYPE)]
    module_outputs = [("NAME", PRODUCED_TYPE)]
    # Delete module_options if you don't need any
    module_options = {
        "OPTION_NAME": {
            "help": "DESCRIPTION",
            "type": TYPE,
            "default": VALUE,
```

(continues on next page)

(continued from previous page)

```

    },
}

def get_software_dependencies(self):
    return super(ModuleInfo, self).get_software_dependencies() + [
        # Add your own dependencies to this list
        # Remove this method if you don't need to add any
    ]

```

```

from pimlico.core.modules.base import BaseModuleExecutor

class ModuleExecutor(BaseModuleExecutor):
    def execute(self):
        input_data = self.info.get_input("NAME")
        self.log.info("MESSAGES")

        # DO STUFF

        with SOME_WRITER(self.info.get_absolute_output_dir("NAME")) as writer:
            # Do what the writer requires

```

1.1.5 Writing document map modules

Todo: Write a guide to building document map modules.

For now, the skeletons below are a useful starting point, but there should be a more fulsome explanation here of what document map modules are all about and how to use them.

Todo: Document map module guides needs to be updated for new datatypes.

Skeleton new module

To make developing a new module a little quicker, here's a skeleton module info and executor for a document map module. It follows the most common method for defining the executor, which is to use the multiprocessing-based executor factory.

```

from pimlico.core.modules.map import DocumentMapModuleInfo
from pimlico.datatypes.tar import TarredCorpusType

class ModuleInfo(DocumentMapModuleInfo):
    module_type_name = "NAME"
    module_readable_name = "READABLE NAME"
    module_inputs = [("NAME", TarredCorpusType(DOCUMENT_TYPE))]
    module_outputs = [("NAME", PRODUCED_TYPE)]
    module_options = {
        "OPTION_NAME": {
            "help": "DESCRIPTION",
            "type": TYPE,
            "default": VALUE,

```

(continues on next page)

(continued from previous page)

```

    },
}

def get_software_dependencies(self):
    return super(ModuleInfo, self).get_software_dependencies() + [
        # Add your own dependencies to this list
    ]

def get_writer(self, output_name, output_dir, append=False):
    if output_name == "NAME":
        # Instantiate a writer for this output, using the given output dir
        # and passing append in as a kwarg
        return WRITER_CLASS(output_dir, append=append)

```

A bare-bones executor:

```

from pimlico.core.modules.map.multiproc import multiprocessing_executor_factory

def process_document(worker, archive_name, doc_name, *data):
    # Do something to process the document...

    # Return an object to send to the writer
    return output

ModuleExecutor = multiprocessing_executor_factory(process_document)

```

Or getting slightly more sophisticated:

```

from pimlico.core.modules.map.multiproc import multiprocessing_executor_factory

def process_document(worker, archive_name, doc_name, *data):
    # Do something to process the document

    # Return a tuple of objects to send to each writer
    # If you only defined a single output, you can just return a single object
    return output1, output2, ...

# You don't have to, but you can also define pre- and postprocessing
# both at the executor level and worker level

def preprocess(executor):
    pass

def postprocess(executor, error=None):
    pass

def set_up_worker(worker):
    pass

def tear_down_worker(worker, error=None):

```

(continues on next page)

```
pass
```

```
ModuleExecutor = multiprocessing_executor_factory(  
    process_document,  
    preprocess_fn=preprocess, postprocess_fn=postprocess,  
    worker_set_up_fn=set_up_worker, worker_tear_down_fn=tear_down_worker,  
)
```

1.1.6 Filter modules

Filter modules appear in pipeline config, but never get executed directly, instead producing their output on the fly when it is needed.

There are two types of filter modules in Pimlico:

- All *document map modules* can be used as filters.
- Other modules may be defined in such a way that they always function as filters.

Using document map modules as filters

See *this guide* for how to create document map modules, which process each document in an input iterable corpus, producing one document in the output corpus for each. Many of the core Pimlico modules are document map modules.

Any document map module can be used as a filter simply by specifying `filter=True` in its options. It will then not appear in the module execution schedule (output by the `status` command), but will get executed on the fly by any module that uses its output. It will be initialized when the downstream module starts accessing the output, and then the single-document processing routine will be run on each document to produce the corresponding output document as the downstream module iterates over the corpus.

It is possible to chain together filter modules in sequence.

Other filter modules

Todo: Filter module guide needs to be updated for new datatypes. This section is currently completely wrong – **ignore it!** This is quite a substantial change.

The difficulty of describing what you need to do here suggests we might want to provide some utilities to make this easier!

A module can be defined so that it always functions as a filter by setting `module_executable=False` on its module-info class. Pimlico will assume that its outputs are ready as soon as its inputs are ready and will not try to execute it. The module developer must ensure that the outputs get produced when necessary.

This form of filter is typically appropriate for very simple transformations of data. For example, it might perform a simple conversion of one datatype into another to allow the output of a module to be used as if it had a different datatype. However, it is possible to do more sophisticated processing in a filter module, though the implementation is a little more tricky (`tar_filter` is an example of this).

Defining

Define a filter module something like this:

```
class ModuleInfo(BaseModuleInfo):
    module_type_name = "my_module_name"
    module_executable = False # This is the crucial instruction to treat this as a
    ↪filter
    module_inputs = [] # Define inputs
    module_outputs = [] # Define at least one output, which we'll produce as
    ↪needed
    module_options = {} # Any options you need

    def instantiate_output_datatype(self, output_name, output_datatype, **kwargs):
        # Here we produce the desired output datatype,
        # using the inputs acquired from self.get_input(name)
        return MyOutputDatatype()
```

You don't need to create an `execute.py`, since it's not executable, so Pimlico will not try to load a module executor. Any processing you need to do should be put inside the datatype, so that it's performed when the datatype is used (e.g. when iterating over it), but not when `instantiate_output_datatype()` is called or when the datatype is instantiated, as these happen every time the pipeline is loaded.

A trick that can be useful to wrap up functionality in a filter datatype is to define a new datatype that does the necessary processing on the fly and to set its class attribute `emulated_datatype` to point to a datatype class that should be used instead for the purposes of type checking. The built-in `tar_filter` module uses this trick.

Either way, you should **take care with imports**. Remember that the `execute.py` of executable modules is only imported when a module is to be run, meaning that we can load the pipeline config without importing any dependencies needed to run the module. If you put processing in a specially defined datatype class that has dependencies, make sure that they're not imported at the top of `info.py`, but only when the datatype is used.

1.1.7 Multistage modules

Multistage modules are used to encapsulate a module that is executed in several consecutive runs. You can think of each stage as being its own module, but where the whole sequence of modules is always executed together. The multistage module simply chains together these individual modules so that you only include a single module instance in your pipeline definition.

One common example of a use case for multistage modules is where some fairly time-consuming preprocessing needs to be done on an input dataset. If you put all of the processing into a single module, you can end up in an irritating situation where the lengthy data preprocessing succeeds, but something goes wrong in the main execution code. You then fix the problem and have to run all the preprocessing again.

Most obvious solution to this is to separate the preprocessing and main execution into two separate modules. But then, if you want to reuse your module sometime in the future, you have to remember to always put the preprocessing module before the main one in your pipeline (or infer this from the datatypes!). And if you have more than these two modules (say, a sequence of several, or preprocessing of several inputs) this starts to make pipeline development frustrating.

A multistage module groups these internal modules into one logical unit, allowing them to be used together by including a single module instance and also to share parameters.

Defining a multistage module

Component stages

The first step in defining a multistage module is to define its individual stages. These are actually defined in exactly the same way as normal modules. (This means that they can also be used separately.)

If you're writing these modules specifically to provide the stages of your multistage module (rather than tying together already existing modules for convenience), you probably want to put them all in subpackages.

For an ordinary module, *we used the directory structure*:

```
src/python/myproject/modules/  
  __init__.py  
  mymodule/  
    __init__.py  
    info.py  
    execute.py
```

Now, we'll use something like this:

```
src/python/myproject/modules/  
  __init__.py  
  my_ms_module/  
    __init__.py  
    info.py  
    module1/  
      __init__.py  
      info.py  
      execute.py  
    module2/  
      __init__.py  
      info.py  
      execute.py
```

Note that `module1` and `module2` both have the typical structure of a module definition: an `info.py` to define the module-info, and an `execute.py` to define the executor. At the top level, we've just got an `info.py`. It's in here that we'll define the multistage module. We don't need an `execute.py` for that, since it just ties together the other modules, using their executors at execution time.

Multistage module-info

With our component modules that constitute the stages defined, we now just need to tie them together. We do this by defining a module-info for the multistage module in its `info.py`. Instead of subclassing `BaseModuleInfo`, as usual, we create the `ModuleInfo` class using the factory function `multistage_module()`.

```
ModuleInfo = multistage_module("module_name",  
  [  
    # Stages to be defined here...  
  ]  
)
```

In other respects, this module-info works in the same way as usual: it's a class (return by the factory) called `ModuleInfo` in the `info.py`.

`multistage_module()` takes two arguments: a module name (equivalent to the `module_name` attribute of a normal module-info) and a list of instances of `ModuleStage`.

Connecting inputs and outputs

Connections between the outputs and inputs of the stages work in a very similar way to connections between module instances in a pipeline. The same type checking system is employed and data is passed between the stages (i.e. between consecutive executions) as if the stages were separate modules.

Each stage is defined as an instance of *ModuleStage*:

```
[
  ModuleStage("stage_name", TheModuleInfoClass, connections=[...], output_
↪connections=[...])
]
```

The parameter `connections` defines how the stage's inputs are connected up to either the outputs of previous stages or inputs to the multistage module. Just like in pipeline config files, if no explicit input connections are given, the default input to a stage is connected to the default output from the previous one in the list.

There are two classes you can use to define input connections.

InternalModuleConnection This makes an explicit connection to the output of another stage.

You must specify the name of the input (to this stage) that you're connecting. You may specify the name of the output to connect it to (defaults to the default output). You may also give the name of the stage that the output comes from (defaults to the previous one).

```
[
  ModuleStage("stage1", FirstInfo,
    # FirstInfo has an output called "corpus", which we connect explicitly to the_
↪next stage
    # We could leave out the "corpus" here, if it's the default output from_
↪FirstInfo
    ModuleStage("stage2", SecondInfo, connections=[InternalModuleConnection("data"
↪", "corpus")]),
    # We connect the same output from stage1 to stage3
    ModuleStage("stage3", ThirdInfo, connections=[InternalModuleConnection("data",
↪ "corpus", "stage1")]),
]
```

ModuleInputConnection: This makes a connection to an input to the whole multistage module.

Note that you don't have to explicitly define the multistage module's inputs anywhere: you just mark certain inputs to certain stages as coming from outside the multistage module, using this class.

```
[
  ModuleStage("stage1", FirstInfo, [ModuleInputConnection("raw_data")]),
  ModuleStage("stage2", SecondInfo, [InternalModuleConnection("data", "corpus"
↪)]),
  ModuleStage("stage3", ThirdInfo, [InternalModuleConnection("data", "corpus",
↪"stage1")]),
]
```

Here, the module type `FirstInfo` has an input called `raw_data`. We've specified that this needs to come in directly as an input to the multistage module – when we use the multistage module in a pipeline, it must be connected up with some earlier module.

The multistage module's input created by doing this will also have the name `raw_data` (specified using a parameter `input_raw_data` in the config file). You can override this, if you want to use a different name:

```
[
  ModuleStage("stage1", FirstInfo, [ModuleInputConnection("raw_data", "data
↪")]),
  ModuleStage("stage2", SecondInfo, [InternalModuleConnection("data", "corpus
↪")]),
  ModuleStage("stage3", ThirdInfo, [InternalModuleConnection("data", "corpus",
↪"stage1")]),
]
```

This would be necessary if two stages both had inputs called `raw_data`, which you want to come from different data sources. You would then simply connect them to different inputs to the multistage module:

```
[
  ModuleStage("stage1", FirstInfo, [ModuleInputConnection("raw_data", "first_
↪data")]),
  ModuleStage("stage2", SecondInfo, [ModuleInputConnection("raw_data", "second_
↪data")]),
  ModuleStage("stage3", ThirdInfo, [InternalModuleConnection("data", "corpus",
↪"stage1")]),
]
```

Conversely, you might deliberately connect the inputs from two stages to the same input to the multistage module, by using the same multistage input name twice. (Of course, the two stages are not required to have overlapping input names for this to work.) This will result in the multistage just requiring one input, which get used by both stages.

```
[
  ModuleStage("stage1", FirstInfo,
    [ModuleInputConnection("raw_data", "first_data"), ↪
↪ModuleInputConnection("dict", "vocab")]),
  ModuleStage("stage2", SecondInfo,
    [ModuleInputConnection("raw_data", "second_data"), ↪
↪ModuleInputConnection("vocabulary", "vocab")]),
  ModuleStage("stage3", ThirdInfo, [InternalModuleConnection("data", "corpus",
↪"stage1")]),
]
```

By default, the multistage module has just a single output: the default output of the last stage in the list. You can specify any of the outputs of any of the stages to be provided as an output to the multistage module. Use the `output_connections` parameter when defining the stage.

This parameter should be a list of instances of *ModuleOutputConnection*. Just like with input connections, if you don't specify otherwise, the multistage module's output will have the same name as the output from the stage module. But you can override this when giving the output connection.

```
[
  ModuleStage("stage1", FirstInfo, [ModuleInputConnection("raw_data", "first_data
↪")]),
  ModuleStage("stage2", SecondInfo, [ModuleInputConnection("raw_data", "second_data
↪")]),
    output_connections=[ModuleOutputConnection("model")]), # This ↪
↪output will just be called "model"
  ModuleStage("stage3", ThirdInfo, [InternalModuleConnection("data", "corpus",
↪"stage1"),
    output_connections=[ModuleOutputConnection("model", "stage3_model")]),
]
```

Module options

The parameters of the multistage module that can be specified when it is used in a pipeline config (those usually defined in the `module_options` attribute) include all of the options to all of the stages. The option names are simply `<stage_name>_<option_name>`.

So, in the above example, if `FirstInfo` has an option called `threshold`, the multistage module will have an option `stage1_threshold`, which gets passed through to `stage1` when it is run.

Often you might wish to specify one parameter to the multistage module that gets used by several stages. Say `stage2` had a `cutoff` parameter and we always wanted to use the same value as the `threshold` for `stage1`. Instead of having to specify `stage1_threshold` and `stage2_cutoff` every time in your config file, you can assign a single name to an option (say `threshold`) for the multistage module, whose value gets passed through to the appropriate options of the stages.

Do this by specifying a dictionary as the `option_connections` parameter to `ModuleStage`, whose keys are names of the stage module type's options and whose values are the new option names for the multistage module that you want to map to those stage options. You can use the same multistage module option name multiple times, which will cause only a single option to be added to the multistage module (using the definition from the first stage), which gets mapped to multiple stage options.

To implement that above example, you would give:

```
[
  ModuleStage("stage1", FirstInfo, [ModuleInputConnection("raw_data", "first_data
↪"),
                                option_connections={"threshold": "threshold"}),
  ModuleStage("stage2", SecondInfo, [ModuleInputConnection("raw_data", "second_data
↪"),
                                     [ModuleOutputConnection("model")],
                                     option_connections={"cutoff": "threshold"}),
  ModuleStage("stage3", ThirdInfo, [InternalModuleConnection("data", "corpus",
↪"stage1"),
                                    [ModuleOutputConnection("model", "stage3_model")]],
]
```

If you know that the different stages have distinct option name, or that they should always tie their values together where their option names overlap, you can set `use_stage_option_names=True` on the stages. This will cause the stage-name prefix not to be added to the option name when connecting it to the multistage module's option.

You can also force this behaviour for all stages by setting `use_stage_option_names=True` when you call `multistage_module()`. Any explicit option name mappings you provide via `option_connections` will override this.

Running

To run a multistage module once you've used it in your pipeline config, you run one stage at a time, as if they were separate module instances.

Say we've used the above multistage module in a pipeline like so:

```
[model_train]
type=myproject.modules.my_ms_module
stage1_threshold=10
stage2_cutoff=10
```

The normal way to run this module would be to use the `run` command with the module name:

```
./pimlico.sh mypipeline.conf run model_train
```

If we do this, Pimlico will choose the next unexecuted stage that's ready to run (presumably `stage1` at this point). Once that's done, you can run the same command again to execute `stage2`.

You can also select a specific stage to execute by using the module name `<ms_module_name>:<stage_name>`, e.g. `model_train:stage2`. (Note that `stage2` doesn't actually depend on `stage1`, so it's perfectly plausible that we might want to execute them in a different order.)

If you want to execute multiple stages at once, just use this scheme to specify each of them as a module name for the run command. Remember, Pimlico can take any number of modules and execute them in sequence:

```
./pimlico.sh mypipeline.conf run model_train:stage1 model_train:stage2
```

Or, if you want to execute all of them, you can use the stage name `*` or `all` as a shorthand:

```
./pimlico.sh mypipeline.conf run model_train:all
```

Finally, if you're not sure what stages a multistage module has, use the module name `<ms_module_name>:?.` The run command will then just output a list of stages and exit.

1.1.8 Running one pipeline on multiple computers

Multiple servers

In most of the examples, we've been setting up a pipeline, with a config file, some source code and some data, all on one machine. Then we run each module in turn, checking that it has all the software and data that it needs to run.

But it's not unusual to find yourself needing to process a dataset across different computers. For example, you have access to a server with lots of CPUs and one module in your pipeline would benefit greatly from parallelizing lots of little tasks over them. However, you don't have permission to install software on that server that you need for another module.

This is not a problem: you can simply put your config file and code on both machines. After running one module on one machine, you copy over its output to the place on the other machine where Pimlico expects to find it. Then you're ready to run the next module on the second machine.

Pimlico is designed to handle this situation nicely.

- **It doesn't expect software requirements for all modules to be satisfied before you can run any of them.** Software dependencies are checked only for modules about to be run and the code used to execute a module is not even loaded until you actually run the module.
- **It doesn't require you to execute your pipeline in order.** If the output from a module is available where it's expected to be, you can happily run any modules that take that data as input, even if the pipeline up to that point doesn't appear to have been executed (e.g. if it's been run on another machine).
- **It provides you with tools to make it easier to copy data between machines.** You can easily copy the output data from one module to the appropriate location on another server, so it's ready to be used as input to another module there.

Copying data between computers

Let's assume you've got your pipeline set up, with identical config files, on two computers: `server_a` and `server_b`. You've run the first module in your pipeline, `module1`, on `server_a` and want to run the next, `module2`, which takes input from `module1`, on `server_b`.

The procedure is as follows:

- **Dump** the data from the pipeline on `server_a`. This packages up the output data for a module in a single file.
- **Copy** the dumped file from `server_a` to `server_b`, in whatever way is most convenient, e.g., using `scp`.
- **Load** the dumped file into the pipeline on `server_b`. This unpacks the data directory for the file and puts it in Pimlico's data directory for the module.

For example, on `server_a`:

```
$ ./pimlico.sh pipeline.conf dump module1
$ scp ~/module1.tar.gz server_b:~/
```

Note that the `dump` command created a `.tar.gz` file in your home directory. If you want to put it somewhere else, use the `--output` option to specify a directory. The file is named after the module that you're dumping.

Now, log into `server_b` and load the data.

```
$ ./pimlico.sh pipeline.conf load ~/module1.tar.gz
```

Now `module1`'s output data is in the right place and ready for use by `module2`.

The `dump` and `load` commands can also process data for multiple modules at once. For example:

```
$ mkdir ~/modules
$ ./pimlico.sh pipeline.conf dump module1 ... module10 --output ~/modules
$ scp -r ~/modules server_b:~/
```

Then on `server_b`:

```
$ ./pimlico.sh pipeline.conf load ~/modules/*
```

Other issues

Aside from getting data between the servers, there are certain issues that often arise when running a pipeline across multiple servers.

- **Shared Pimlico codebase.** If you share the directory that contains Pimlico's code across servers (e.g. NFS or `rsync`), you can have problems resulting from sharing the libraries it installs. See *instructions for using multiple virtualenvs* for the solution.
- **Shared home directory.** If you share your home directory across servers, using the same `.pimlico` local config file might be a problem. See *Local configuration* for various possible solutions.

1.1.9 Documenting your own code

Pimlico's documentation is produced using [Sphinx](#). The Pimlico codebase includes a tool for generating documentation of Pimlico's built-in modules, including things like a table of the module's available config options and its input and outputs.

You can also use this tool yourself to generate documentation of your own code that uses Pimlico. Typically, you will use in your own project some of Pimlico's built-in modules and some of your own.

Refer to Sphinx's documentation for how to build normal Sphinx documentation – writing your own ReST documents and using the `apidoc` tool to generate API docs. Here we describe how to create a basic Sphinx setup that will generate a reference for your custom Pimlico modules.

It is assumed that you've got a working Pimlico setup and have already successfully written some modules.

Basic doc setup

Create a `docs` directory in your project root (the directory in which you have `pimlico/` and your own `src/`, etc).

Put a Sphinx `conf.py` in there. You can start from the very basic skeleton [here](#).

You'll also want a `Makefile` to build your docs with. You can use the basic Sphinx one as a starting point. Here's a version of that that already includes an extra target for building your module docs.

Finally, create a root document for your documentation, `index.rst`. This should include a table of contents which includes the generated module docs. You can use [this one](#) as a template.

Building the module docs

Take a look in the `Makefile` (if you've used our one as a starting point) and set the variables at the top to point to the Python package that contains the Pimlico modules you want to document.

The make target there runs the tool `modulegen` in the Pimlico codebase. Just run, in the `docs/`:

```
make modules
```

You can also do this manually:

```
python -m pimlico.utils.docs.modulegen --path python.path.to.modules modules/
```

(The Pimlico codebase must, of course, be importable. The simplest way to ensure this is to use Pimlico's `python` alias in its `bin/` directory.)

There is now a set of `.rst` files in the `modules/` output directory, which can be built using Sphinx by running `make html`.

Your beautiful docs are now in the `_build/` directory!

1.2 Core docs

A set of articles on the core aspects and features of Pimlico.

1.2.1 Downloading Pimlico

To start a new project using Pimlico, download the `newproject.py` script. It will create a template pipeline config file to get you started and download the latest version of Pimlico to accompany it.

See *Setting up a new project using Pimlico* for more detail.

Pimlico's source code is available on [Github](#).

Manual setup

If for some reason you don't want to use the `newproject.py` script, you can set up a project yourself. Download Pimlico [from Github](#).

Simply download the whole source code as a `.zip` or `.tar.gz` file and uncompress it. This will produce a directory called `pimlico`, followed by a long incomprehensible string, which you can rename simply `pimlico`.

Pimlico has a few basic dependencies, but these will be automatically downloaded the first time you load it.

1.2.2 Pipeline config

A Pimlico pipeline, as read from a config file (`pimlico.core.config.PipelineConfig`) contains all the information about the pipeline being processed and provides access to specific modules in it. A config file looks much like a standard `.ini` file, with sections headed by `[section_name]` headings, containing key-value parameters of the form `key=value`.

Each section, except for `vars` and `pipeline`, defines a module instance in the pipeline. Some of these can be executed, others act as filters on the outputs of other modules, or input readers.

Each section that defines a module has a `type` parameter. Usually, this is a fully-qualified Python package name that leads to the module type's Python code (that package containing the `info` Python module). A special type is `alias`. This simply defines a module alias – an alternative name for an already defined module. It should have exactly one other parameter, `input`, specifying the name of the module we're aliasing.

Special sections

- **vars:** May contain any variable definitions, to be used later on in the pipeline. Further down, expressions like `%(varname)s` will be expanded into the value assigned to `varname` in the `vars` section.
- **pipeline:** Main pipeline-wide configuration. The following options are required for every pipeline:
 - `name`: a single-word name for the pipeline, used to determine where files are stored
 - `release`: the release of Pimlico for which the config file was written. It is considered compatible with later minor versions of the same major release, but not with later major releases. Typically, a user receiving the pipeline config will get hold of an appropriate version of the Pimlico codebase to run it with.

Other optional settings:

- `python_path`: a path or paths, relative to the directory containing the config file, in which Python modules/packages used by the pipeline can be found. Typically, a config file is distributed with a directory of Python code providing extra modules, datatypes, etc. Multiple paths are separated by colons (`:`).

Special variable substitutions

Certain variable substitutions are always available, in addition to those defined in `vars` sections. Use them anywhere in your config file with an expression like `%(varname)s` (note the `s` at the end).

- **pimlico_root:** Root directory of Pimlico, usually the directory `pimlico/` within the project directory.
- **project_root:** Root directory of the whole project. Current assumed to always be the parent directory of `pimlico_root`.
- **output_dir:** Path to output dir (usually `output` in Pimlico root).
- **home:** Running user's home directory (on Unix and Windows, see Python's `os.path.expanduser()`).
- **test_data_dir:** Directory in Pimlico distribution where test data is stored (`test/data` in Pimlico root). Used in test pipelines, which take all their input data from this directory.

For example, to point a parameter to a file located within the project root:

```
param=%(project_root)s/data/myfile.txt
```

Directives

Certain special directives are processed when reading config files. They are lines that begin with `%%`, followed by the directive name and any arguments.

- **variant:** Allows a line to be included only when loading a particular variant of a pipeline. For more detail on pipeline variants, see *Pipeline variants*.

The variant name is specified as part of the directive in the form: `variant:variant_name`. You may include the line in more than one variant by specifying multiple names, separated by commas (and no spaces). You can use the default variant “main”, so that the line will be left out of other variants. The rest of the line, after the directive and variant name(s) is the content that will be included in those variants.

```
[my_module]
type=path.to.module
%%variant:main size=52
%%variant:smaller size=7
```

An alternative notation for the variant directive is provided to make config files more readable. Instead of `variant:variant_name`, you can write `(variant_name)`. So the above example becomes:

```
[my_module]
type=path.to.module
%%(main) size=52
%%(smaller) size=7
```

- **novariant:** A line to be included only when not loading a variant of the pipeline. Equivalent to `variant:main`.

```
[my_module]
type=path.to.module
%%novariant size=52
%%variant:smaller size=7
```

- **include:** Include the entire contents of another file. The filename, specified relative to the config file in which the directive is found, is given after a space.
- **abstract:** Marks a config file as being abstract. This means that Pimlico will not allow it to be loaded as a top-level config file, but only allow it to be included in another config file.
- **copy:** Copies all config settings from another module, whose name is given as the sole argument. May be used multiple times in the same module and later copies will override earlier. Settings given explicitly in the module’s config override any copied settings.

All parameters are copied, including things like `type`. Any parameter can be overridden in the copying module instance. Any parameter can be excluded from the copy by naming it after the module name. Separate multiple exclusions with spaces.

The directive even allows you to copy parameters from multiple modules by using the directive multiple times, though this is not very often useful. In this case, the values are copied (and overridden) in the order of the directives.

For example, to reuse all the parameters from `module1` in `module2`, only specifying them once:

```
[module1]
type=some.module.type
input=moduleA
param1=56
param2=never
```

(continues on next page)

(continued from previous page)

```
param3=0.75

[module2]
# Copy all params from module1
%%copy module1
# Override the input module
input=moduleB
```

Multiple parameter values

Sometimes you want to write a whole load of modules that are almost identical, varying in just one or two parameters. You can give a parameter multiple values by writing them separated by vertical bars (`|`). The module definition will be expanded to produce a separate module for each value, with all the other parameters being identical.

For example, this will produce three module instances, all having the same `num_lines` parameter, but each with a different `num_chars`:

```
[my_module]
type=module.type.path
num_lines=10
num_chars=3|10|20
```

You can even do this with multiple parameters of the same module and the expanded modules will cover all combinations of the parameter assignments.

For example:

```
[my_module]
type=module.type.path
num_lines=10|50|100
num_chars=3|10|20
```

Tying alternatives

You can change the behaviour of alternative values using the `tie_alts` option. `tie_alts=T` will cause parameters within the same module that have multiple alternatives to be expanded in parallel, rather than taking the product of the alternative sets. So, if `option_a` has 5 values and `option_b` has 5 values, instead of producing 25 pipeline modules, we'll only produce 5, matching up each pair of values in their alternatives.

```
[my_module]
type=module.type.path
tie_alts=T
option_a=1|2|3|4|5
option_b=one|two|three|four|five
```

If you want to tie together the alternative values on some parameters, but not others, you can specify groups of parameter names to tie using the `tie_alts` option. Each group is separated by spaces and the names of parameters to tie within a group are separated by `|` s. Any parameters that have alternative values but are not specified in one of the groups are not tied to anything else.

For example, the following module config will tie together `option_a`'s alternatives with `option_b`'s, but produce all combinations of them with `option_c`'s alternatives, resulting in $3*2=6$ versions of the module (`my_module[option_a=1~option_b=one~option_c=x]`,

`my_module[option_a=1~option_b=one~option_c=y],my_module[option_a=2~option_b=two~option_c=x]`
etc).

```
[my_module]
type=module.type.path
tie_alts=option_a|option_b
option_a=1|2|3
option_b=one|two|three
option_c=x|y
```

Using this method, you must give the parameter names in `tie_alts` exactly as you specify them in the config. For example, although for a particular module you might be able to specify a certain input (the default) using the name `input` or a specific name like `input_data`, these will not be recognised as being the same parameter in the process of expanding out the combinations of alternatives.

Naming alternatives

Each module will be given a distinct name, based on the varied parameters. If just one is varied, the names will be of the form `module_name[param_value]`. If multiple parameters are varied at once, the names will be `module_name[param_name0=param_value0~param_name1=param_value1~...]`. So, the first example above will produce: `my_module[3],my_module[10]` and `my_module[20]`. And the second will produce: `my_module[num_lines=10~num_chars=3], my_module[num_lines=10~num_chars=10]`, etc.

You can also specify your own identifier for the alternative parameter values, instead of using the values themselves (say, for example, if it's a long file path). Specify it surrounded by curly braces at the start of the value in the alternatives list. For example:

```
[my_module]
type=module.type.path
file_path={small}/home/me/data/corpus/small_version|{big}/home/me/data/corpus/big_
↪version
```

This will result in the modules `my_module[small]` and `my_module[big]`, instead of using the whole file path to distinguish them.

An alternative approach to naming the expanded alternatives can be selected using the `alt_naming` parameter. The default behaviour described above corresponds to `alt_naming=full`. If you choose `alt_naming=pos`, the alternative parameter settings (using names where available, as above) will be distinguished like positional arguments, without making explicit what parameter each value corresponds to. This can make for nice concise names in cases where it's clear what parameters the values refer to.

If you specify `alt_naming=full` explicitly, you can also give a further option `alt_naming=full(inputnames)`. This has the effect of removing the `input_` from the start of named inputs. This often makes for intuitive module names, but is not the default behaviour, since there's no guarantee that the input name (without the initial `input_`) does not clash with an option name.

Another possibility, which is occasionally appropriate, is `alt_naming=option(<name>)`, where `<name>` is the name of an option that has alternatives. In this case, the names of the alternatives for the whole module will be taken directly from the alternative names on that option only. (E.g. specified by `{name}` or inherited from a previous module, see below). You may specify multiple option names, separated by commas, and the corresponding alt names will be separated by `~`. If there's only one option with alternatives, this is equivalent to `alt_naming=pos`. If there are multiple, it might often lead to name clashes. The circumstance in which this is most commonly appropriate is where you use `tie_alts=T`, so it's sufficient to distinguish the alternatives by the name associated with just one option.

Expanding alternatives down the pipeline

If a module takes input from a module that has been expanded into multiple versions for alternative parameter values, it too will automatically get expanded, as if all the multiple versions of the previous module had been given as alternative values for the input parameter. For example, the following will result in 3 versions of `my_module` (`my_module [1]`, etc) and 3 corresponding versions of `my_next_module` (`my_next_module [1]`, etc):

```
[my_module]
type=module.type.path
option_a=1|2|3

[my_next_module]
type=another.module.type.path
input=my_module
```

Where possible, names given to the alternative parameter values in the first module will be carried through to the next.

Module variables: passing information through the pipeline

When a pipeline is read in, each module instance has a set of *module variables* associated with it. In your config file, you may specify assignments to the variables for a particular module. Each module inherits all of the variable assignments from modules that it receives its inputs from.

The main reason for having module variables is to be able to do things in later modules that depend on what path through the pipeline an input came from. Once you have defined the sequence of processing steps that pass module variables through the pipeline, apply mappings to them, etc, you can use them in the parameters passed into modules.

Basic assignment

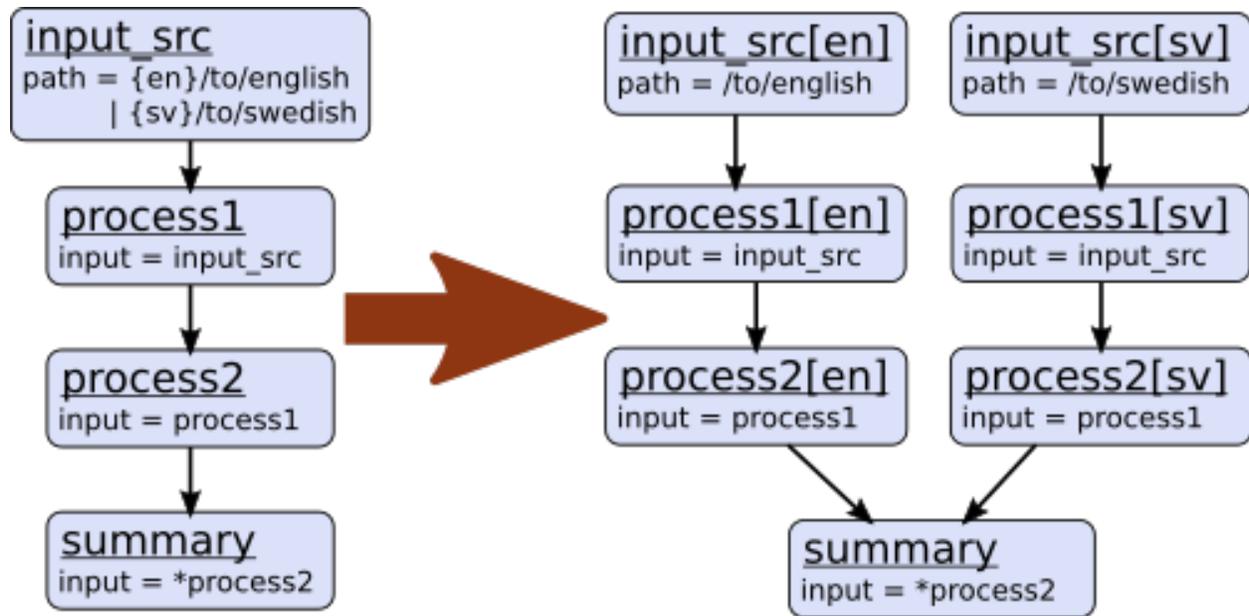
Module variables are set by including parameters in a module's config of the form `modvar_<name> = <value>`. This will assign `value` to the variable `name` for this module. The simplest form of assignment is just a string literal, enclosed in double quotes:

```
[my_module]
type=module.type.path
modvar_myvar = "Value of my variable"
```

Names of alternatives

Say we have a simple pipeline that has a single source of data, with different versions of the dataset for different languages (English and Swedish). A series of modules process each language in an identical way and, at the end, outputs from all languages are collected by a single `summary` module. This final module may need to know what language each of its incoming datasets represents, so that it can output something that we can understand.

The two languages are given as alternative values for a parameter `path`, and the whole pipeline gets automatically expanded into two paths for the two alternatives:



The `summary` module gets its two inputs for the two different languages as a multiple-input: this means we could expand this pipeline to as many languages as we want, just by adding to the `input_src` module's `path` parameter.

However, as far as `summary` is concerned, this is just a list of datasets – it doesn't know that one of them is English and one is Swedish. But let's say we want it to output a table of results. We're going to need some labels to identify the languages.

The solution is to add a module variable to the first module that takes different values when it gets expanded into two modules. For this, we can use the `altname` function in a `modvar` assignment: this assigns the name of the expanded module's alternative for a given parameter that has alternatives in the config.

```
[input_src]
path={en}/to/english | {sv}/to/swedish
modvar_lang=altname(path)
```

Now the expanded module `input_src[en]` will have the module variable `lang="en"` and the Swedish version `lang="sv"`. This value gets passed from module to module down the two paths in the pipeline.

Other assignment syntax

A further function `map` allows you to apply a mapping to a value, rather like a Python dictionary lookup. Its first argument is the value to be mapped (or anything that expands to a value, using `modvar` assignment syntax). The second is the mapping. This is simply a space-separated list of source-target mappings of the form `source -> target`. Typically both the sources and targets will be string literals.

Now we can give our languages legible names. (Here we're splitting the definition over multiple indented lines, as permitted by config file syntax, which makes the mapping easier to read.)

```
[input_src]
path={en}/to/english | {sv}/to/swedish
modvar_lang=map(
    altname(path),
    "en" -> "English"
    "sv" -> "Svenska")
```

The assignments may also reference variable names, including those previously assigned to in the same module and those received from the input modules.

```
[input_src]
path={en}/to/english | {sv}/to/swedish
modvar_lang=altname(path)
modvar_lang_name=map(
    lang,
    "en" -> "English"
    "sv" -> "Svenska")
```

If a module gets two values for the same variable from multiple inputs, the first value will simply be overridden by the second. Sometimes it's useful to map module variables from specific inputs to different modvar names. For example, if we're combining two different languages, we might need to keep track of what the two languages we combined were. We can do this using the notation `input_name.var_name`, which refers to the value of module variable `var_name` that was received from input `input_name`.

```
[input_src]
path={en}/to/english | {sv}/to/swedish
modvar_lang=altname(path)

[combiner]
type=my.language.combiner
input_lang_a=lang_data
input_lang_b=lang_data
modvar_first_lang=lang_a.lang
modvar_second_lang=lang_b.lang
```

If a module inherits multiple values for the same variable from the **same input** (i.e. a multiple-input), they are all kept and treated as a list. The most common way to then use the values is via the `join` function. Like Python's `string.join`, this turns a list into a single string by joining the values with a given separator string. Use `join(sep, list)` to join the values coming from some list `modvar list` on the separator `sep`.

You can get the number of values in a list `modvar` using `len(list)`, which works just like Python's `len()`.

Use in module parameters

To make something in a module's execution dependent on its module variables, you can insert them into module parameters.

For example, say we want one of the module's parameters to make use of the `lang` variable we defined above:

```
[input_src]
path={en}/to/english | {sv}/to/swedish
modvar_lang=altname(path)
some_param=${lang}
```

Note the difference to other variable substitutions, which use the `%(varname)s` notation. For modvars, we use the notation `$(varname)`.

We can also put the value in the middle of other text:

```
[input_src]
path={en}/to/english | {sv}/to/swedish
modvar_lang=altname(path)
some_param=myval-${lang}-continues
```

The modvar processing to compute a particular module's set of variable assignments is performed before the substitution. This means that you can do any modvar processing specific to the module instance, in the various ways defined above, and use the resulting value in other parameters. For example:

```
[input_src]
path={en}/to/english | {sv}/to/swedish
modvar_lang=altname(path)
modvar_mapped_lang=map(lang,
    "en" -> "eng"
    "sv" -> "swe"
)
some_param=$(mapped_lang)
```

You can also place in the `$(...)` construct any of the variable processing operations shown above for assignments to module variables. This is a little more concise than first assigning values to modvars, if you don't need to use the variables again anywhere else. For example:

```
[input_src]
path={en}/to/english | {sv}/to/swedish
some_param=$(map(altname(path),
    "en" -> "eng"
    "sv" -> "swe"
))
```

Usage in module code

A module's executor can also retrieve the values assigned to module variables from the `module_variables` attribute of the module-info associated with the input dataset. Sometimes this can be useful when you are writing your own module code, though the above usage to pass values from (or dependent on) module variables into module parameters is more flexible, so should generally be preferred.

```
# Code in executor
# This is a MultipleInput-type input, so we get a list of datasets
datasets = self.info.get_input()
for d in datasets:
    language = d.module.module_variables["lang"]
```

1.2.3 Pipeline variants

You can create several different versions of a pipeline, called pipeline *variants* in a single config file. The data corresponding to each will be kept completely separate. This is useful when you want multiple versions of a pipeline that are almost identical, but have some small differences.

The most common use of this, though by no means the only, is to create a variant that is faster to run than the main pipeline for the purposes of quickly testing the whole pipeline during development.

Every pipeline has by default one variant, called `main`. You define other variants simply by using special directives to mark particular lines as belonging to a particular variant. Lines with no variant marking will appear in all variants.

Loading variants

If you don't specify otherwise when loading a pipeline, the `main` variant will be loaded. Use the `--variant` parameter (or `-v`) to specify another variant by name:

```
./pimlico.sh mypipeline.conf -v smaller status
```

To see a list of all available variants of a particular pipeline, use the *variants* command:

```
./pimlico.sh mypipeline.conf variants
```

Variant directives

Directives are processed when a pipeline config file is read in, before the file is parsed to build a pipeline. They are lines that begin with `%%`, followed by the directive name and any arguments. See *Directives* for details of other directives.

- **variant:** This line will be included only when loading a particular variant of a pipeline.

The variant name is specified in the form: `variant:variant_name`. You may include the line in more than one variant by specifying multiple names, separated by commas (and no spaces). You can use the default variant “main”, so that the line will be left out of other variants. The rest of the line, after the directive and variant name(s) is the content that will be included in those variants.

```
[my_module]
type=path.to.module
%%variant:main size=52
%%variant:smaller size=7
```

An alternative notation makes config files more readable. Instead of `%%variant:variant_name`, write `%%(variant_name)`. So the above example becomes:

```
[my_module]
type=path.to.module
%%(main) size=52
%%(smaller) size=7
```

- **novariant:** A line to be included only when not loading a variant of the pipeline. Equivalent to `variant:main`.

```
[my_module]
type=path.to.module
%%novariant size=52
%%variant:smaller size=7
```

Example

The following example config file, defines one variant, `small`, aside from the default `main` variant.

```
[pipeline]
name=myvariants
release=0.8
python_path=%(project_root)s/src/python

# Load a dataset
[input_data]
type=pimlico.modules.input.text.raw_text_files
files=%(home)s/data/*
```

(continues on next page)

(continued from previous page)

```
# For the small version, we cut down the dataset to just 10 documents
# We don't need this module at all in the main variant
%(small) [small_data]
%(small) type=pimlico.modules.corpora.subset
%(small) size=10

# Tokenize the text
# Control where the input data comes from in the different variants
# The main variant simply uses the full, uncut corpus
[tokenize]
type=pimlico.modules.text.simple_tokenize
%(small) input=small_data
%(main) input=input_data
```

The main variant will be loaded if you don't specify otherwise. In this version the module `small_data` doesn't exist at all and `tokenize` takes its input from `input_data`.

```
./pimlico.sh myvariants.conf status
```

You can load the small variant by giving its name on the command line. This includes the `small_data` module and `tokenize` gets its input from there, making it much faster to test.

```
./pimlico.sh myvariants.conf -v small status
```

1.2.4 Pimlico module structure

This document describes the code structure for Pimlico module types in full.

For a basic guide to writing your own modules, see [Writing Pimlico modules](#).

Todo: Write documentation for this

1.2.5 Module dependencies

In a Pimlico pipeline, you typically use lots of different external software packages. Some are Python packages, others system tools, Java libraries, whatever. Even the core modules that core with Pimlico between them depend on a huge amount of software.

Naturally, we don't want to have to install *all* of this software before you can run even a simple Pimlico pipeline that doesn't use all (or any) of it. So, we keep the core dependencies of Pimlico to an absolute minimum, and then check whether the necessary software dependencies are installed each time a pipeline module is going to be run.

Core dependencies

Certain dependencies are required for Pimlico to run at all, or needed so often that you wouldn't get far without installing them. These are defined in `pimlico.core.dependencies.core`, and when you run the Pimlico command-line interface, it checks they're available and tries to install them if they're not.

Module dependencies

Each module type defines its own set of software dependencies, if it has any. When you try to run the module, Pimlico runs some checks to try to make sure that all of these are available.

If some of them are not, it may be possible to install them automatically, straight from Pimlico. In particular, many Python packages can be very easily installed using [Pip](#). If this is the case for one of the missing dependencies, Pimlico will tell you in the error output, and you can install them using the `install` command (with the module name/number as an argument).

Virtualenv

In order to simplify automatic installation, Pimlico is always run within a virtual environment, using [Virtualenv](#). This means that any Python packages installed by Pip will live in a local directory within the Pimlico codebase that you're running and won't interfere with anything else on your system.

When you run Pimlico for the first time, it will create a new virtualenv for this purpose. Every time you run it after that, it will use this same environment, so anything you install will continue to be available.

Custom virtualenv

Most of the time, you don't even need to be aware of the virtualenv that Python's running in¹. Under certain circumstances, you might need to use a custom virtualenv.

For example, say you're running your pipeline over different servers, but have the pipeline and Pimlico codebase on a shared network drive. Then you can find that the software installed in the virtualenv on one machine is incompatible with the system-wide software on the other.

You can specify a name for a custom virtualenv using the environment variable `PIMENV`. The first time you run Pimlico with this set, it will automatically create the new virtualenv.

```
$ PIMENV=myenv ./pimlico.sh mypipeline.conf status
```

Replace `myenv` with a name that better reflects its use (e.g. name of the server).

Every time you run Pimlico on that server, set the `PIMENV` environment variable in the same way.

In case you want to get to the virtualenv itself, you can find it in `pimlico/lib/virtualenv/myenv`.

Note: Pimlico previously used another environment variable `VIRTUALENV`, which gave a path to the virtualenv. You can still use this, but, unless you have a good reason to, it's easier to use `PIMENV`.

Defining module dependencies

Todo: Describe how module dependencies are defined for different types of deps

¹ If you're interested, it lives in `pimlico/lib/virtualenv/default`

Some examples

Todo: Include some examples from the core modules of how deps are defined and some special cases of software fetching

1.2.6 Local configuration

As well as knowing about the pipeline you're running, Pimlico also needs to know some things about the setup of the system on which you're running it. This is completely independent of the pipeline config: the same pipeline can be run on different systems with different local setups.

A couple of settings must always be provided for Pimlico: the **long-term** and **short-term stores** (see [Data stores](#) below). Other system settings may be specified as necessary. (At the time of writing, there aren't any, but they will be documented here as they arise.) See [Other Pimlico settings](#) below.

Specific modules may also have system-level settings. For example, a module that calls an external tool may need to know the location of that tool, or how much memory it can use on this system. Any that apply to the built-in Pimlico modules are listed below in [Settings for built-in modules](#).

Local config file location

Pimlico looks in various places to find the local config settings. Settings are loaded in a particular order, overriding earlier versions of the same setting as we go (see `pimlico.core.config.PipelineConfig.load_local_config()`).

Settings are specified with the following order of precedence (those later override the earlier):

```
local config file < host-specific config file < cmd-line overrides
```

Most often, you'll just specify all settings in the main local config file. This is a file in your home directory named `.pimlico`. This must exist for Pimlico to be able to run at all.

Host-specific config

If you share your home directory between different computers (e.g. a networked filesystem), the above setup could cause a problem, as you may need a different local config on the different computers. Pimlico allows you to have special config files that only get read on machines with a particular hostname.

For example, say I have two computers, `localbox` and `remotebox`, which share a home directory. I've created my `.pimlico` local config file on `localbox`, but need to specify a different storage location on `remotebox`. I simply create another config file called `.pimlico_remotebox`#[hostname]_`. Pimlico will load first the basic local config in ``.pimlico` and then override those settings with what it reads from the host-specific config file.

You can also specify a hostname prefix to match. Say I've got a whole load of computers I want to be able to run on, with hostnames `remotebox1`, `remotebox2`, etc. If I create a config file called `.pimlico_remotebox-`, it will be used on all of these hosts.

Command-line overrides

Occasionally, you might want to specify a local config setting just for one run of Pimlico. Use the `--override-local-config` (or `-l`) to specify a value for an individual setting in the form `setting=value`. For example:

```
./pimlico.sh mypipeline.conf -l somesetting=5 run mymodule
```

If you want to override multiple settings, simply use the option multiple times.

Custom location

If the above solutions don't work for you, you can also explicitly specify on the command line an alternative location from which to load the local config file that Pimlico typically expects to find in `~/pimlico`.

Use the `--local-config` parameter to give a filename to use instead of the `~/pimlico`.

For example, if your home directory is shared across servers and the above hostname-specific config solution doesn't work in your case, you can fall back to pointing Pimlico at your own host-specific config file.

Data stores

Pimlico needs to know where to put and find output files as it executes. Settings are given in the local config, since they apply to all Pimlico pipelines you run and may vary from system to system. Note that Pimlico will make sure that different pipelines don't interfere with each other's output (provided you give them different names): all pipelines store their output and look for their input within these same base locations.

See *Data storage* for an explanation of Pimlico's data store system.

At least one store must be given in the local config:

```
store=/path/to/storage/root
```

You may specify as many storage locations as you like, giving each a name:

```
store_fast=/path/to/fast/store
store_big=/path/to/big/store
```

If you specify named stores *and* an unnamed one, the unnamed one will be used as the default output store. Otherwise, the first in the file will be the default.

```
store=/path/to/a/store           # This will be the default output store
store_fast=/path/to/fast/store  # These will be additional, named stores
store_big=/path/to/big/store
```

Other Pimlico settings

In future, there will no doubt be more settings that you can specify at the system level for Pimlico. These will be documented here as they arise.

Settings for built-in modules

Specific modules may consult the local config to allow you to specify settings for them. We cannot document them here for all modules, as we don't know what modules are being developed outside the core codebase. However, we can provide a list here of the settings consulted by built-in Pimlico modules.

There aren't any yet, but they will be listed here as they arise.

Footnotes:

1.2.7 Data storage

Pimlico needs to know where to put and find output files as it executes, in order to store data and pass it between modules. On any particular system running Pimlico, multiple locations (**stores**) may be used as storage and Pimlico will check all of them when it's looking for a module's data.

Single store

Let's start with a simple setup with just one store. A setting `store` in the local config (see *Local configuration*) specifies the root directory of this store. This applies to all Pimlico pipelines you run on this system and Pimlico will make sure that different pipelines don't interfere with each other's output (provided you give them different names).

When you run a pipeline module, its output will be stored in a subdirectory specific to that pipeline and that module with the store's root directory. When Pimlico needs to use that data as input to another module, it will look in the appropriate directory within the store.

Multiple stores

For various reasons, you may wish to store Pimlico data in multiple locations.

For example, one common scenario is that you have access to a disk that is fast to write to (call it *fast-disk*), but not very big, and another disk (e.g. over a network filesystem) that has lots of space, but is slower (call it *big-disk*). You therefore want Pimlico to output its data, much of which might only be used fleetingly and then no longer needed, to *fast-disk*, so the processing runs quickly. Then, you want to move the output from certain modules over to *big-disk*, to make space on *fast-disk*.

We can define two stores for Pimlico to use and give them names. The first ("fast") will be used to output data to (just like the sole store in the previous section). The second ("big"), however, will also be checked for module data, meaning that we can move data from "fast" to "big" whenever we like.

Instead of using the `store` parameter in the local config, we use multiple `store_<name>` parameters. One of them (the first one, or the one given by `store` with no name, if you include that) will be treated as the default output store.

Specify the locations in the local config like this:

```
store_fast=/path/to/fast/store
store_big=/path/to/big/store
```

Remember, these paths are not specific to a pipeline: all pipelines will use different subdirectories of these ones.

To check what stores you've got in your current configuration, use the `stores` command.

Moving data between stores

Say you've got a two-store setup like in the previous example. You've now run a module that produces a lot of output and want to move it to your big disk and have Pimlico read it from there.

You don't need to replicate the directory structure yourself and move module output between stores. Pimlico has a command *movestores* to do this for you. Specify the name of the store you want to move data to (*big* in this case) and the names or numbers of the modules whose data you want to move.

Once you've done that, Pimlico should continue to behave as it did before, just as if the data was still in its original location.

Updating from the old storage system

Prior to v0.8, Pimlico used a different system of storage locations. If you have a local config file (*~/ .pimlico*) from an earlier version you will see deprecation warnings.

Change something like this:

```
long_term_store=/path/to/long/store
short_term_store=/path/to/short/store
```

to something like this:

```
store_long=/path/to/long/store
store_short=/path/to/short/store
```

Or, if you only ever needed one storage location, simply this:

```
store=/path/to/store
```

1.2.8 Python scripts

All the heavy work of your data-processing is implemented in Pimlico modules, either by loading core Pimlico modules from your pipeline config file or by writing your own modules. Sometimes, however, it can be handy to write a quick Python script to get hold of the output of one of your pipeline's modules and inspect it or do something with it.

This can be easily done writing a Python script and using the *python* shell command to run it. This command loads your pipeline config (just like all others) and then either runs a script you've specified on the command line, or enters an interactive Python shell. The advantages of this over just running the normal *python* command on the command line are that the script is run in the same execution context used for your pipeline (e.g. using the Pimlico instance's *virtualenv*) and that the loaded pipeline is available to you, so you can easily can hold of its data locations, datatypes, etc.

Accessing the pipeline

At the top of your Python script, you can get hold of the loaded pipeline config instance like this:

```
from pimlico.cli.pyshell import get_pipeline

pipeline = get_pipeline()
```

Now you can use this to get to, among other things, the pipeline's modules and their input and output datasets. A module called *module1* can be accessed by treating the pipeline like a dict:

```
module = pipeline["module1"]
```

This gives you the `ModuleInfo` instance for that module, giving access to its inputs, outputs, options, etc:

```
data = module.get_output("output_name")
```

Writing and running scripts

All of the above code to access a pipeline can be put in a Python script somewhere in your codebase and run from the command line. Let's say I create a script `src/python/scripts/myscript.py` containing:

```
from pimlico.cli.pyshell import get_pipeline

pipeline = get_pipeline()
module = pipeline["module1"]
data = module.get_output("output_name")
# Here we can start probing the data using whatever interface the datatype provides
print data
```

Now I can run this from the root directory of my project as follows:

```
./pimlico.sh mypipeline.conf python src/python/scripts/myscript.py
```

1.3 Core Pimlico modules

Pimlico comes with a substantial collection of module types that provide wrappers around existing NLP and machine learning tools, as well as a number of general tools for processing datasets that are useful for many applications.

1.3.1 !! C&C parser

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.candc
Executable	yes

Wrapper around the original C&C parser.

Takes tokenized input and parses it with C&C. The output is written exactly as it comes out from C&C. It contains both GRs and supertags, plus POS-tags, etc.

The wrapper uses C&C's SOAP server. It sets the SOAP server running in the background and then calls C&C's SOAP client for each document. If parallelizing, multiple SOAP servers are set going and each one is kept constantly fed with documents.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
documents	invalid input type specification

Outputs

Name	Type(s)
parsed	invalid output type specification

Options

Name	Description	Type
model	Absolute path to models directory or name of model set. If not an absolute path, assumed to be a subdirectory of the candc models dir (see instructions in models/candc/README on how to fetch pre-trained models)	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_candc_module]
type=pimlico.modules.candc
input_documents=module_a.some_output
```

This example usage includes more options.

```
[my_candc_module]
type=pimlico.modules.candc
input_documents=module_a.some_output
model=ccgbank
```

1.3.2 !! Stanford CoreNLP

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.corenlp
Executable	yes

Process documents one at a time with the [Stanford CoreNLP toolkit](#). CoreNLP provides a large number of NLP tools, including a POS-tagger, various parsers, named-entity recognition and coreference resolution. Most of these tools can be run using this module.

The module uses the CoreNLP server to accept many inputs without the overhead of loading models. If parallelizing, only a single CoreNLP server is run, since this is designed to set multiple Java threads running if it receives multiple queries at the same time. Multiple Python processes send queries to the server and process the output.

The module has no non-optional outputs, since what sort of output is available depends on the options you pass in: that is, on which tools are run. Use the annotations option to choose which word annotations are added. Otherwise, simply select the outputs that you want and the necessary tools will be run in the CoreNLP pipeline to produce those outputs.

Currently, the module only accepts tokenized input. If pre-POS-tagged input is given, for example, the POS tags won't be handed into CoreNLP. In the future, this will be implemented.

We also don't currently provide a way of choosing models other than the standard, pre-trained English models. This is a small addition that will be implemented in the future.

Todo: Update to new datatypes system and add test pipelines

Inputs

Name	Type(s)
documents	invalid input type specification

Outputs

No non-optional outputs

Optional

Name	Type(s)
annotations	invalid output type specification
tokenized	invalid output type specification
parse	invalid output type specification
parse-deps	invalid output type specification
dep-parse	invalid output type specification
raw	invalid output type specification
coref	invalid output type specification

Options

Name	Description	Type
gzip	If True, each output, except annotations, for each document is gzipped. This can help reduce the storage occupied by e.g. parser or coref output. Default: False	bool
time-out	Timeout for the CoreNLP server, which is applied to every job (document). Number of seconds. By default, we use the server's default timeout (15 secs), but you may want to increase this for more intensive tasks, like coref	float
readable	If True, JSON outputs are formatted in a readable fashion, pretty printed. Otherwise, they're as compact as possible. Default: False	bool
annotators	Comma-separated list of word annotations to add, from CoreNLP's annotators. Choose from: word, pos, lemma, ner	string
dep_type	Type of dependency parse to output, when outputting dependency parses, either from a constituency parse or direct dependency parse. Choose from the three types allowed by CoreNLP: 'basic', 'collapsed' or 'collapsed-ccprocessed'	'basic', 'collapsed' or 'collapsed-ccprocessed'

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_corenlp_module]
type=pimlico.modules.corenlp
input_documents=module_a.some_output
```

This example usage includes more options.

```
[my_corenlp_module]
type=pimlico.modules.corenlp
input_documents=module_a.some_output
gzip=T
timeout=0.1
readable=T
annotators=
dep_type=collapsed-ccprocessed
```

1.3.3 Corpus manipulation

Core modules for generic manipulation of mainly iterable corpora.

Corpus concatenation

Path	pimlico.modules.corpora.concat
Executable	no

Concatenate two (or more) corpora to produce a bigger corpus.

They must have the same data point type, or one must be a subtype of the other.

This is a filter module. It is not executable, so won't appear in a pipeline's list of modules that can be run. It produces its output for the next module on the fly when the next module needs it.

Inputs

Name	Type(s)
corpora	<i>list of iterable_corpus</i>

Outputs

Name	Type(s)
corpus	<i>corpus with data-point from input</i>

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_concat_module]
type=pimlico.modules.corpora.concat
input_corpora=module_a.some_output
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *concat*

Corpus statistics

Path	pimlico.modules.corpora.corpus_stats
Executable	yes

Some basic statistics about tokenized corpora

Counts the number of tokens, sentences and distinct tokens in a corpus.

Inputs

Name	Type(s)
corpus	<i>grouped_corpus</i> < <i>TokenizedDocumentType</i> >

Outputs

Name	Type(s)
stats	<i>named_file</i>

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_corpus_stats_module]
type=pimlico.modules.corpora.corpus_stats
input_corpus=module_a.some_output
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *stats*

Human-readable formatting

Path	pimlico.modules.corpora.format
Executable	yes

Corpus formatter

Pimlico provides a data browser to make it easy to view documents in a tarred document corpus. Some datatypes provide a way to format the data for display in the browser, whilst others provide multiple formatters that display the data in different ways.

This module allows you to use this formatting functionality to output the formatted data as a corpus. Since the formatting operations are designed for display, this is generally only useful to output the data for human consumption.

Inputs

Name	Type(s)
corpus	<i>grouped_corpus</i>

Outputs

Name	Type(s)
formatted	<i>grouped_corpus</i> <RawTextDocumentType>

Options

Name	Description	Type
for- mat- ter	Fully qualified class name of a formatter to use to format the data. If not specified, the default formatter is used, which uses the datatype's <code>browser_display</code> attribute if available, or falls back to just converting documents to unicode	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_format_module]
type=pimlico.modules.corpora.format
input_corpus=module_a.some_output
```

This example usage includes more options.

```
[my_format_module]
type=pimlico.modules.corpora.format
input_corpus=module_a.some_output
formatter=path.to.formatter.FormatterClass
```

Archive grouper (filter)

Path	pimlico.modules.corpora.group
Executable	no

Group the data points (documents) of an iterable corpus into fixed-size archives. This is a standard thing to do at the start of the pipeline, since it's a handy way to store many (potentially small) files without running into filesystem problems.

The documents are simply grouped linearly into a series of groups (archives) such that each (apart from the last) contains the given number of documents.

After grouping documents in this way, document map modules can be called on the corpus and the grouping will be preserved as the corpus passes through the pipeline.

Note: This module used to be called `tar_filter`, but has been renamed in keeping with other changes in the new datatype system.

There also used to be a `tar` module that wrote the grouped corpus to disk. This has now been removed, since most of the time it's fine to use this filter module instead. If you really want to store the grouped corpus, you can use the `store` module.

This is a filter module. It is not executable, so won't appear in a pipeline's list of modules that can be run. It produces its output for the next module on the fly when the next module needs it.

Inputs

Name	Type(s)
documents	<i>iterable_corpus</i>

Outputs

Name	Type(s)
documents	<i>grouped corpus with input doc type</i>

Options

Name	Description	Type
archive_size	Number of documents to include in each archive (default: 1k)	int
archive_basename	Base name to use for archive tar files. The archive number is appended to this. (Default: 'archive')	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_group_module]
type=pimlico.modules.corpora.group
input_documents=module_a.some_output
```

This example usage includes more options.

```
[my_group_module]
type=pimlico.modules.corpora.group
input_documents=module_a.some_output
archive_size=1000
archive_basename=archive
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *store*
- *group*

Interleaved corpora

Path	pimlico.modules.corpora.interleave
Executable	no

Interleave data points from two (or more) corpora to produce a bigger corpus.

Similar to *concat*, but interleaves the documents when iterating. Preserves the order of documents within corpora and takes documents two each corpus in inverse proportion to its length, i.e. spreads out a smaller corpus so we don't finish iterating over it earlier than the longer one.

They must have the same data point type, or one must be a subtype of the other.

In theory, we could find the most specific common ancestor and use that as the output type, but this is not currently implemented and may not be worth the trouble. Perhaps we will add this in future.

This is a filter module. It is not executable, so won't appear in a pipeline's list of modules that can be run. It produces its output for the next module on the fly when the next module needs it.

Inputs

Name	Type(s)
corpora	<i>list of grouped_corpus</i>

Outputs

Name	Type(s)
corpus	<i>grouped corpus with input doc type</i>

Options

Name	Description	Type
archive_size	Documents are regrouped into new archives. Number of documents to include in each archive (default: 1k)	string
archive_basename	Documents are regrouped into new archives. Base name to use for archive tar files. The archive number is appended to this. (Default: 'archive')	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_interleave_module]
type=pimlico.modules.corpora.interleave
input_corpora=module_a.some_output
```

This example usage includes more options.

```
[my_interleave_module]
type=pimlico.modules.corpora.interleave
input_corpora=module_a.some_output
archive_size=1000
archive_basename=archive
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *interleave*

Corpus document list filter

Path	pimlico.modules.corpora.list_filter
Executable	yes

Similar to *split*, but instead of taking a random split of the dataset, splits it according to a given list of documents, putting those in the list in one set and the rest in another.

Inputs

Name	Type(s)
corpus	<i>grouped_corpus</i>
list	<i>string_list</i>

Outputs

Name	Type(s)
set1	<i>grouped corpus with input doc type</i>
set2	<i>grouped corpus with input doc type</i>

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_list_filter_module]
type=pimlico.modules.corpora.list_filter
input_corpus=module_a.some_output
input_list=module_a.some_output
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *list_filter*

Corpus split

Path	pimlico.modules.corpora.split
Executable	yes

Split a tarred corpus into two subsets. Useful for dividing a dataset into training and test subsets. The output datasets have the same type as the input. The documents to put in each set are selected randomly. Running the module multiple times will give different splits.

Note that you can use this multiple times successively to split more than two ways. For example, say you wanted a training set with 80% of your data, a dev set with 10% and a test set with 10%, split it first into training and non-training 80-20, then split the non-training 50-50 into dev and test.

The module also outputs a list of the document names that were included in the first set. Optionally, it outputs the same thing for the second input too. Note that you might prefer to only store this list for the smaller set: e.g. in a training-test split, store only the test document list, as the training list will be much larger. In such a case, just put the smaller set first and don't request the optional output *doc_list2*.

Inputs

Name	Type(s)
corpus	<i>grouped_corpus</i>

Outputs

Name	Type(s)
set1	<i>grouped_corpus with input doc type</i>
set2	<i>grouped_corpus with input doc type</i>
doc_list1	<i>string_list</i>

Optional

Name	Type(s)
doc_list2	<i>string_list</i>

Options

Name	Description	Type
set1_size	Proportion of the corpus to put in the first set, float between 0.0 and 1.0. If an integer >1 is given, this is treated as the absolute number of documents to put in the first set, rather than a proportion. Default: 0.2 (20%)	float

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_split_module]
type=pimlico.modules.corpora.split
input_corpus=module_a.some_output
```

This example usage includes more options.

```
[my_split_module]
type=pimlico.modules.corpora.split
input_corpus=module_a.some_output
set1_size=0.20
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *split*

Store a corpus

Path	pimlico.modules.corpora.store
Executable	yes

Store a corpus

Take documents from a corpus and write them to disk using the standard writer for the corpus' data point type. This is useful where documents are produced on the fly, for example from some filter module or from an input reader, but where it is desirable to store the produced corpus for further use, rather than always running the filters/readers each time the corpus' documents are needed.

Inputs

Name	Type(s)
corpus	<i>grouped_corpus</i>

Outputs

Name	Type(s)
corpus	<i>grouped corpus with input doc type</i>

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_store_module]
type=pimlico.modules.corpora.store
input_corpus=module_a.some_output
```

Corpus subset

Path	pimlico.modules.corpora.subset
Executable	no

Simple filter to truncate a dataset after a given number of documents, potentially offsetting by a number of documents. Mainly useful for creating small subsets of a corpus for testing a pipeline before running on the full corpus.

Can be run on an iterable corpus or a tarred corpus. If the input is a tarred corpus, the filter will emulate a tarred corpus with the appropriate datatype, passing through the archive names from the input.

When a number of valid documents is required (calculating corpus length when skipping invalid docs), if one is stored in the metadata as `valid_documents`, that count is used instead of iterating over the data to count them up.

This is a filter module. It is not executable, so won't appear in a pipeline's list of modules that can be run. It produces its output for the next module on the fly when the next module needs it.

Inputs

Name	Type(s)
corpus	<i>iterable_corpus</i>

Outputs

Name	Type(s)
corpus	<i>corpus with data-point from input</i>

Options

Name	Description	Type
off-set	Number of documents to skip at the beginning of the corpus (default: 0, start at beginning)	int
skip_invalid	Skip over any invalid documents so that the output subset contains the chosen number of (valid) documents (or as many as possible) and no invalid ones. By default, invalid documents are passed through and counted towards the subset size	bool
size	(required) Number of documents to include	int

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_subset_module]
type=pimlico.modules.corpora.subset
input_corpus=module_a.some_output
size=100
```

This example usage includes more options.

```
[my_subset_module]
type=pimlico.modules.corpora.subset
input_corpus=module_a.some_output
offset=0
skip_invalid=T
size=100
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *subset*

Corpus vocab builder

Path	pimlico.modules.corpora.vocab_builder
Executable	yes

Builds a dictionary (or vocabulary) for a tokenized corpus. This is a data structure that assigns an integer ID to every distinct word seen in the corpus, optionally applying thresholds so that some words are left out.

Similar to *pimlico.modules.features.vocab_builder*, which builds two vocabs, one for terms and one for features.

Inputs

Name	Type(s)
text	<i>grouped_corpus</i> < <i>TokenizedDocumentType</i> >

Outputs

Name	Type(s)
vocab	<i>dictionary</i>

Options

Name	Description	Type
prune_size	Prune the dictionary if it reaches this size. Setting a lower value avoids getting stuck with too big a dictionary to be able to prune and slowing things down, but means that the final pruning will less accurately reflect the true corpus stats. Should be considerably higher than limit (if used). Set to 0 to disable. Default: 2M (Gensim's default)	int
max_include	Include terms that occur in max this proportion of documents	float
oov	Use the final index to represent chars that will be out of vocabulary after applying threshold/limit filters. Applied even if the count is 0. Represent OOVs using the given string in the vocabulary	string
limit	Limit vocab size to this number of most common entries (after other filters)	int
threshold	Minimum number of occurrences required of a term to be included	int
include	Ensure that certain words are always included in the vocabulary, even if they don't make it past the various filters, or are never seen in the corpus. Give as a comma-separated list	comma-separated list of strings

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_vocab_builder_module]
type=pimlico.modules.corpora.vocab_builder
input_text=module_a.some_output
```

This example usage includes more options.

```
[my_vocab_builder_module]
type=pimlico.modules.corpora.vocab_builder
input_text=module_a.some_output
prune_at=2000000
oov=text
limit=10k
threshold=100
include=word1,word2,...
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *vocab_builder*

!! Token frequency counter

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.corpora.vocab_counter
Executable	yes

Count the frequency of each token of a vocabulary in a given corpus (most often the corpus on which the vocabulary was built).

Note that this distribution is not otherwise available along with the vocabulary. It stores the document frequency counts - how many documents each token appears in - which may sometimes be a close enough approximation to the actual frequencies. But, for example, when working with character-level tokens, this estimate will be very poor.

The output will be a 1D array whose size is the length of the vocabulary, or the length plus one, if `oov_excluded=T` (used if the corpus has been mapped so that OOVs are represented by the ID `vocab_size+1`, instead of having a special token).

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
corpus	invalid input type specification
vocab	invalid input type specification

Outputs

Name	Type(s)
distribution	invalid output type specification

Options

Name	Description	Type
oov_excluded	Indicates that the corpus has been mapped so that OOVs are represented by the ID vocab_size+1, instead of having a special token in the vocab	bool

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_vocab_counter_module]
type=pimlico.modules.corpora.vocab_counter
input_corpus=module_a.some_output
input_vocab=module_a.some_output
```

This example usage includes more options.

```
[my_vocab_counter_module]
type=pimlico.modules.corpora.vocab_counter
input_corpus=module_a.some_output
input_vocab=module_a.some_output
oov_excluded=T
```

Tokenized corpus to ID mapper

Path	pimlico.modules.corpora.vocab_mapper
Executable	yes

Todo: Write description of vocab mapper module

Todo: Add test pipeline and test

Inputs

Name	Type(s)
text	<i>grouped_corpus</i> <TokenizedDocumentType>
vocab	<i>dictionary</i>

Outputs

Name	Type(s)
ids	<i>grouped_corpus</i> <IntegerListsDocumentType>

Options

Name	Description	Type
oov	If given, special token to map all OOV tokens to. Otherwise, use vocab_size+1 as index. Special value 'skip' simply skips over OOV tokens	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_vocab_mapper_module]
type=pimlico.modules.corpora.vocab_mapper
input_text=module_a.some_output
input_vocab=module_a.some_output
```

This example usage includes more options.

```
[my_vocab_mapper_module]
type=pimlico.modules.corpora.vocab_mapper
input_text=module_a.some_output
input_vocab=module_a.some_output
oov=value
```

1.3.4 Embedding feature extractors and trainers

Modules for extracting features from which to learn word embeddings from corpora, and for training embeddings.

Some of these don't actually learn the embeddings, they just produce features which can then be fed into an embedding learning module, such as a form of matrix factorization. Note that you can train embeddings not only using the trainers here, but also using generic matrix manipulation techniques, for example the factorization methods provided by sklearn.

!! Dependency feature extractor for embeddings

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.embeddings.dependencies
Executable	yes

Todo: Document this module

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
dependencies	invalid input type specification

Outputs

Name	Type(s)
term_features	invalid output type specification

Options

Name	Description	Type
lemma	Use lemmas as terms instead of the word form. Note that if you didn't run a lemmatizer before dependency parsing the lemmas are probably actually just copies of the word forms	bool
con-dense_prep	Where a word is modified ... TODO	string
term_pos	Only extract features for terms whose POSs are in this comma-separated list. Put a * at the end to denote POS prefixes	comma-separated list of strings
skip_types	Dependency relations to skip, separated by commas	comma-separated list of strings

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_embedding_dep_features_module]
type=pimlico.modules.embeddings.dependencies
input_dependencies=module_a.some_output
```

This example usage includes more options.

```
[my_embedding_dep_features_module]
type=pimlico.modules.embeddings.dependencies
input_dependencies=module_a.some_output
lemma=T
condense_prep=value
term_pos=
skip_types=
```

Store embeddings (internal)

Path	pimlico.modules.embeddings.store_embeddings
Executable	yes

Simply stores embeddings in the Pimlico internal format.

This is not often needed, but can be useful if reading embeddings for an input reader that is slower than reading from the internal format. Then you can use this module to do the reading and store the result before passing it to other modules.

Inputs

Name	Type(s)
embeddings	<i>embeddings</i>

Outputs

Name	Type(s)
embeddings	<i>embeddings</i>

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_store_embeddings_module]
type=pimlico.modules.embeddings.store_embeddings
input_embeddings=module_a.some_output
```

!! Store in TSV format

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.embeddings.store_tsv
Executable	yes

Takes embeddings stored in the default format used within Pimlico pipelines (see *Embeddings*) and stores them as TSV files.

These are suitable as input to the [Tensorflow Projector](<https://projector.tensorflow.org/>).

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
embeddings	invalid input type specification

Outputs

Name	Type(s)
embeddings	invalid output type specification

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_store_tsv_module]
type=pimlico.modules.embeddings.store_tsv
input_embeddings=module_a.some_output
```

!! Store in word2vec format

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.embeddings.store_word2vec
Executable	yes

Takes embeddings stored in the default format used within Pimlico pipelines (see *Embeddings*) and stores them using the `word2vec` storage format.

Uses the Gensim implementation of the storage, so depends on Gensim.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
embeddings	invalid input type specification

Outputs

Name	Type(s)
embeddings	invalid output type specification

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_store_word2vec_module]
type=pimlico.modules.embeddings.store_word2vec
input_embeddings=module_a.some_output
```

!! Word2vec embedding trainer

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.embeddings.word2vec
Executable	yes

Word2vec embedding learning algorithm, using [Gensim](#)'s implementation.

Find out more about [word2vec](#).

This module is simply a wrapper to call [Gensim Python \(+C\)](#)'s implementation of word2vec on a Pimlico corpus.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
text	invalid input type specification

Outputs

Name	Type(s)
model	invalid output type specification

Options

Name	Description	Type
iters	number of iterations over the data to perform. Default: 5	int
min_count	word2vec's min_count option: prunes the dictionary of words that appear fewer than this number of times in the corpus. Default: 5	int
negative_samples	number of negative samples to include per positive. Default: 5	int
size	number of dimensions in learned vectors. Default: 200	int

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_word2vec_module]
type=pimlico.modules.embeddings.word2vec
input_text=module_a.some_output
```

This example usage includes more options.

```
[my_word2vec_module]
type=pimlico.modules.embeddings.word2vec
input_text=module_a.some_output
iters=5
min_count=5
negative_samples=5
size=200
```

1.3.5 Feature set processing

Various tools for generic processing of extracted sets of features: building vocabularies, mapping to integer indices, etc.

!! Key-value to term-feature converter

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.features.term_feature_compiler
Executable	yes

Todo: Document this module

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
key_values	invalid input type specification

Outputs

Name	Type(s)
term_features	invalid output type specification

Options

Name	Description	Type
term_keys	Name of keys (feature names in the input) which denote terms. The first one found in the keys of a particular data point will be used as the term for that data point. Any other matches will be removed before using the remaining keys as the data point's features. Default: just 'term'	comma-separated list of strings
include_feature_keys	If True, include the key together with the value from the input key-value pairs as feature names in the output. Otherwise, just use the value. E.g. for input [prop=wordy, poss=my], if True we get features [prop_wordy, poss_my] (both with count 1); if False we get just [wordy, my]. Default: False	bool

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_term_feature_list_module]
type=pimlico.modules.features.term_feature_compiler
input_key_values=module_a.some_output
```

This example usage includes more options.

```
[my_term_feature_list_module]
type=pimlico.modules.features.term_feature_compiler
input_key_values=module_a.some_output
term_keys=term
include_feature_keys=F
```

!! Term-feature matrix builder

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.features.term_feature_matrix_builder
Executable	yes

Todo: Document this module

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
data	invalid input type specification

Outputs

Name	Type(s)
matrix	invalid output type specification

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_term_feature_matrix_builder_module]
type=pimlico.modules.features.term_feature_matrix_builder
input_data=module_a.some_output
```

!! Term-feature corpus vocab builder

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.features.vocab_builder
Executable	yes

Todo: Document this module

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
term_features	invalid input type specification

Outputs

Name	Type(s)
term_vocab	invalid output type specification
feature_vocab	invalid output type specification

Options

Name	Description	Type
feature_limit	Limit vocab size to this number of most common entries (after other filters)	int
feature_max_prop	Include features that occur in max this proportion of documents	float
term_max_prop	Include terms that occur in max this proportion of documents	float
term_threshold	Minimum number of occurrences required of a term to be included	int
feature_threshold	Minimum number of occurrences required of a feature to be included	int
term_limit	Limit vocab size to this number of most common entries (after other filters)	int

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_term_feature_vocab_builder_module]
type=pimlico.modules.features.vocab_builder
input_term_features=module_a.some_output
```

This example usage includes more options.

```
[my_term_feature_vocab_builder_module]
type=pimlico.modules.features.vocab_builder
input_term_features=module_a.some_output
feature_limit=0
feature_max_prop=0.1
term_max_prop=0.1
term_threshold=0
feature_threshold=0
term_limit=0
```

!! Term-feature corpus vocab mapper

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.features.vocab_mapper
Executable	yes

Todo: Document this module

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
data	invalid input type specification
term_vocab	invalid input type specification
feature_vocab	invalid input type specification

Outputs

Name	Type(s)
data	invalid output type specification

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_term_feature_vocab_mapper_module]
type=pimlico.modules.features.vocab_mapper
input_data=module_a.some_output
input_term_vocab=module_a.some_output
input_feature_vocab=module_a.some_output
```

1.3.6 Gensim topic modelling

Modules providing access to topic model training and other routines from [Gensim](#).

LDA trainer

Path	pimlico.modules.gensim.lda
Executable	yes

Trains LDA using Gensim's basic LDA implementation, or the multicore version.

Todo: Add test pipeline and test

Inputs

Name	Type(s)
corpus	<i>grouped_corpus</i> <IntegerListsDocumentType>
vocab	<i>dictionary</i>

Outputs

Name	Type(s)
model	<i>lda_model</i>

Options

Name	Description	Type
eval_every		int
passes	Passes parameter. Default: 1	int
num_topics	Number of topics for the trained model to have. Default: 100	int
eta	Eta prior of word distribution. May be one of special values 'auto' and 'symmetric', or a float. Default: symmetric	'symmetric', 'auto' or a float
decay	Decay parameter. Default: 0.5	float
distributed	Turn on distributed computing. Default: False. Ignored by multicore implementation	bool
minimum_phi_value		float
update_every	Model's update_every parameter. Default: 1. Ignored by multicore implementation	int
tfidf	Transform word counts using TF-IDF when presenting documents to the model for training. Default: False	bool
ignore_terms	Ignore any of these terms in the bags of words when iterating over the corpus to train the model. Typically, you'll want to include an OOV term here if your corpus has one, and any other special terms that are not part of a document's content	comma-separated list of strings
multicore	Use Gensim's multicore implementation of LDA training (gensim.models.ldamulticore). Default is to use gensim.models.ldamodel. Number of cores used for training set by Pimlico's processes parameter	bool
iterations	Max number of iterations in each update. Default: 50	int
offset	Offset parameter. Default: 1.0	float
gamma_threshold		float
alpha	Alpha prior over topic distribution. May be one of special values 'symmetric', 'asymmetric' and 'auto', or a single float, or a list of floats. Default: symmetric	'symmetric', 'asymmetric', 'auto' or a float
minimum_probability		float
chunk-size	Model's chunksize parameter. Chunk size to use for distributed/multicore computing. Default: 2000	int

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_lda_trainer_module]
type=pimlico.modules.gensim.lda
input_corpus=module_a.some_output
input_vocab=module_a.some_output
```

This example usage includes more options.

```
[my_lda_trainer_module]
type=pimlico.modules.gensim.lda
input_corpus=module_a.some_output
input_vocab=module_a.some_output
eval_every=10
```

(continues on next page)

(continued from previous page)

```

passes=1
num_topics=100
eta=symmetric
decay=0.50
distributed=F
minimum_phi_value=0.01
update_every=1
tfidf=F
ignore_terms=
multicore=F
iterations=50
offset=1.00
gamma_threshold=0.00
alpha=symmetric
minimum_probability=0.01
chunksize=2000
    
```

LDA document topic analysis

Path	pimlico.modules.gensim.lda_doc_topics
Executable	yes

Takes a trained LDA model and produces the topic vector for every document in a corpus.

The corpus is given as integer lists documents, which are the integer IDs of the words in each sentence of each document. It is assumed that the corpus uses the same vocabulary to map to integer IDs as the LDA model's training corpus, so no further mapping needs to be done.

Todo: Add test pipeline and test

Inputs

Name	Type(s)
corpus	<i>grouped_corpus</i> <IntegerListsDocumentType>
model	<i>lda_model</i>

Outputs

Name	Type(s)
vectors	<i>grouped_corpus</i> <VectorDocumentType>

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_lda_doc_topics_module]
type=pimlico.modules.gensim.lda_doc_topics
input_corpus=module_a.some_output
input_model=module_a.some_output
```

1.3.7 Input readers

Various input readers for various datatypes. These are used to read in data from some external source, such as a corpus in its distributed format (e.g. XML files or a collection of text files), and present it to the Pimlico pipeline as a Pimlico dataset, which can be used as input to other modules.

They do not typically store the data as a Pimlico dataset, but produce it on the fly, although sometimes it could be appropriate to do otherwise.

Note that there can be multiple input readers for a single datatype. For example, there are many ways to read in a corpus of raw text documents, depending on the format they're stored in. They might be in one big XML file, text files collected into compressed archives, a big text file with document separators, etc. These all require their own input reader and all of them produce the same output corpus type.

Embeddings

Read vector embeddings (e.g. word embeddings) from various storage formats.

There are several formats in common usage and we provide readers for most of these here: *FastText*, *word2vec* and *GloVe*.

FastText embedding reader

Path	pimlico.modules.input.embeddings.fasttext
Executable	yes

Reads in embeddings from the *FastText* format, storing them in the format used internally in Pimlico for embeddings.

Can be used, for example, to read the [pre-trained embeddings](#) offered by Facebook AI.

Currently only reads the text format (*.vec*), not the binary format (*.bin*).

See also:

pimlico.modules.input.embeddings.fasttext_gensim: An alternative reader that uses Gensim's *FastText* format reading code and permits reading from the binary format, which contains more information.

Inputs

No inputs

Outputs

Name	Type(s)
embeddings	<i>embeddings</i>

Options

Name	Description	Type
path	(required) Path to the FastText embedding file	string
limit	Limit to the first N words. Since the files are typically ordered from most to least frequent, this limits to the N most common words	int

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_fasttext_embedding_reader_module]
type=pimlico.modules.input.embeddings.fasttext
path=value
```

This example usage includes more options.

```
[my_fasttext_embedding_reader_module]
type=pimlico.modules.input.embeddings.fasttext
path=value
limit=0
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *fasttext_input_test*

FastText embedding reader (Gensim)

Path	pimlico.modules.input.embeddings.fasttext_gensim
Executable	yes

Reads in embeddings from the [FastText](#) format, storing them in the format used internally in Pimlico for embeddings. This version uses Gensim's implementation of the format reader, so depends on Gensim.

Can be used, for example, to read the [pre-trained embeddings](#) offered by Facebook AI.

Reads only the binary format (`.bin`), not the text format (`.vec`).

See also:

pimlico.modules.input.embeddings.fasttext: An alternative reader that does not use Gensim. It permits (only) reading the text format.

Todo: Add test pipeline. This is slightly difficult, as we need a small FastText binary file, which is harder to produce, since you can't easily just truncate a big file.

Inputs

No inputs

Outputs

Name	Type(s)
embeddings	<i>embeddings</i>

Options

Name	Description	Type
path	(required) Path to the FastText embedding file (.bin)	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_fasttext_embedding_reader_gensim_module]
type=pimlico.modules.input.embeddings.fasttext_gensim
path=value
```

GloVe embedding reader (Gensim)

Path	pimlico.modules.input.embeddings.glove
Executable	yes

Reads in embeddings from the [GloVe](#) format, storing them in the format used internally in Pimlico for embeddings. We use Gensim's implementation of the format reader, so the module depends on Gensim.

Can be used, for example, to read the pre-trained embeddings [offered by Stanford](#).

Note that the format is almost identical to *word2vec*'s text format.

Inputs

No inputs

Outputs

Name	Type(s)
embeddings	<i>embeddings</i>

Options

Name	Description	Type
path	(required) Path to the GloVe embedding file	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_glove_embedding_reader_module]
type=pimlico.modules.input.embeddings.glove
path=value
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *glove_input_test*

Word2vec embedding reader (Gensim)

Path	pimlico.modules.input.embeddings.word2vec
Executable	yes

Reads in embeddings from the *word2vec* format, storing them in the format used internally in Pimlico for embeddings. We use Gensim's implementation of the format reader, so the module depends on Gensim.

Can be used, for example, to read the pre-trained embeddings *offered by Google*.

Inputs

No inputs

Outputs

Name	Type(s)
embeddings	<i>embeddings</i>

Options

Name	Description	Type
binary	Assume input is in word2vec binary format. Default: True	bool
path	(required) Path to the word2vec embedding file (.bin)	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_word2vec_embedding_reader_module]
type=pimlico.modules.input.embeddings.word2vec
path=value
```

This example usage includes more options.

```
[my_word2vec_embedding_reader_module]
type=pimlico.modules.input.embeddings.word2vec
binary=T
path=value
```

Text corpora

Raw text files

Path	pimlico.modules.input.text.raw_text_files
Executable	no

Input reader for raw text file collections. Reads in files from arbitrary locations specified by a list of globs.

The input paths must be absolute paths (or globs), but remember that you can make use of various *special substitutions in the config file* to give paths relative to your project root, or other locations.

The file paths may use `globs` to match multiple files. By default, it is assumed that every filename should exist and every glob should match at least one file. If this does not hold, the dataset is assumed to be not ready. You can override this by placing a `?` at the start of a filename/glob, indicating that it will be included if it exists, but is not depended on for considering the data ready to use.

This is an input module. It takes no pipeline inputs and is used to read in data

Inputs

No inputs

Outputs

Name	Type(s)
corpus	<i>grouped_corpus</i> <RawTextDocumentType>

Options

Name	Description	Type
files	(required) Comma-separated list of absolute paths to files to include in the collection. Paths may include globs. Place a '?' at the start of a filename to indicate that it's optional. You can specify a line range for the file by adding ':X-Y' to the end of the path, where X is the first line and Y the last to be included. Either X or Y may be left empty. (Line numbers are 1-indexed.)	comma-separated list of (line range-limited) file paths
en-cod-ing	Encoding to assume for input files. Default: utf8	string
ex-clude	A list of files to exclude. Specified in the same way as <i>files</i> (except without line ranges). This allows you to specify a glob in <i>files</i> and then exclude individual files from it (you can use globs here too)	comma-separated list of strings
archive_size	Number of documents to include in each archive (default: 1k)	int
en-cod-ing_errors	What to do in the case of invalid characters in the input while decoding (e.g. illegal utf-8 chars). Select 'strict' (default), 'ignore', 'replace'. See Python's str.decode() for details	string
archive_basename	Base name to use for archive tar files. The archive number is appended to this. (Default: 'archive')	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_raw_text_files_reader_module]
type=pimlico.modules.input.text.raw_text_files
files=path1,path2,...
```

This example usage includes more options.

```
[my_raw_text_files_reader_module]
type=pimlico.modules.input.text.raw_text_files
files=path1,path2, ...
encoding=utf8
exclude=text,text, ...
archive_size=1000
encoding_errors=strict
archive_basename=archive
```

Annotated text

Datasets that store text with accompanying annotations, like POS tags or dependency parses.

There are lots of different ways of storing this type of data in common usage. Here we currently only implement variants on one – the VRT format, used by Korp. In future, others should be added, e.g. CoNLL dependency parses.

Datatypes exist for some of these, which should be converted to input readers in due course.

1.3.8 Malt dependency parser

Wrapper around the [Malt dependency parser](#) and data format converters to support connections to other modules.

!! Annotated text to CoNLL dep parse input converter

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.malt.conll_parser_input
Executable	yes

Converts word-annotations to CoNLL format, ready for input into the Malt parser. Annotations must contain words and POS tags. If they contain lemmas, all the better; otherwise the word will be repeated as the lemma.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
annotations	invalid input type specification

Outputs

Name	Type(s)
conll_data	invalid output type specification

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_conll_parser_input_module]
type=pimlico.modules.malt.conll_parser_input
input_annotations=module_a.some_output
```

!! Malt dependency parser

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.malt.parse
Executable	yes

Todo: Document this module

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
documents	invalid input type specification

Outputs

Name	Type(s)
parsed	invalid output type specification

Options

Name	Description	Type
model	Filename of parsing model, or path to the file. If just a filename, assumed to be Malt models dir (models/malt). Default: engmalt.linear-1.7.mco, which can be acquired by 'make malt' in the models dir	string
no_gzip	By default, we gzip each document in the output data. If you don't do this, the output can get very large, since it's quite a verbose output format	bool

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_malt_module]
type=pimlico.modules.malt.parse
input_documents=module_a.some_output
```

This example usage includes more options.

```
[my_malt_module]
type=pimlico.modules.malt.parse
input_documents=module_a.some_output
model=engmalt.linear-1.7.mco
no_gzip=F
```

1.3.9 NLTK

Modules that wrap functionality in the Natural Language Toolkit (NLTK).

Currently, not much is provided here, but adding new modules is easy to do, so hopefully more modules will gradually appear.

!! OpenNLP NIST tokenizer

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.nltk.nist_tokenize
Executable	yes

Sentence splitting and tokenization using the [NLTK NIST tokenizer](#).

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
text	invalid input type specification

Outputs

Name	Type(s)
documents	invalid output type specification

Options

Name	Description	Type
lowercase	Lowercase all output. Default: False	bool
non_european	Use the tokenizer's <code>international_tokenize()</code> method instead of <code>tokenize()</code> . Default: False	bool

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_nltk_nist_tokenizer_module]
type=pimlico.modules.nltk.nist_tokenize
input_text=module_a.some_output
```

This example usage includes more options.

```
[my_nltk_nist_tokenizer_module]
type=pimlico.modules.nltk.nist_tokenize
input_text=module_a.some_output
lowercase=F
non_european=F
```

1.3.10 OpenNLP modules

A collection of module types to wrap individual OpenNLP tools.

!! OpenNLP coreference resolution

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.opennlp.coreference
Executable	yes

Todo: Document this module

Todo: Update to new datatypes system and add test pipeline

Use local config setting `opennlp_memory` to set the limit on Java heap memory for the OpenNLP processes. If parallelizing, this limit is shared between the processes. That is, each OpenNLP worker will have a memory limit of `opennlp_memory / processes`. That setting can use *g*, *G*, *m*, *M*, *k* and *K*, as in the Java setting.

Inputs

Name	Type(s)
parses	invalid input type specification

Outputs

Name	Type(s)
coref	invalid output type specification

Options

Name	Description	Type
gzip	If True, each output, except annotations, for each document is gzipped. This can help reduce the storage occupied by e.g. parser or coref output. Default: False	bool
model	Coreference resolution model, full path or directory name. If a filename is given, it is expected to be in the OpenNLP model directory (models/opennlp/). Default: "" (standard English opennlp model in models/opennlp/)	string
readable	If True, pretty-print the JSON output, so it's human-readable. Default: False	bool
timeout	Timeout in seconds for each individual coref resolution task. If this is exceeded, an InvalidDocument is returned for that document	int

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_opennlp_coref_module]
type=pimlico.modules.opennlp.coreference
input_parsers=module_a.some_output
```

This example usage includes more options.

```
[my_opennlp_coref_module]
type=pimlico.modules.opennlp.coreference
input_parsers=module_a.some_output
gzip=T
model=
readable=T
timeout=0
```

!! OpenNLP coreference resolution

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.opennlp.coreference_pipeline
Executable	yes

Runs the full coreference resolution pipeline using OpenNLP. This includes sentence splitting, tokenization, pos tagging, parsing and coreference resolution. The results of all the stages are available in the output.

Todo: Update to new datatypes system and add test pipeline

Use local config setting `opennlp_memory` to set the limit on Java heap memory for the OpenNLP processes. If parallelizing, this limit is shared between the processes. That is, each OpenNLP worker will have a memory limit of `opennlp_memory / processes`. That setting can use *g*, *G*, *m*, *M*, *k* and *K*, as in the Java setting.

Inputs

Name	Type(s)
text	invalid input type specification

Outputs

Name	Type(s)
coref	invalid output type specification

Optional

Name	Type(s)
tokenized	invalid output type specification
pos	invalid output type specification
parse	invalid output type specification

Options

Name	Description	Type
gzip	If True, each output, except annotations, for each document is gzipped. This can help reduce the storage occupied by e.g. parser or coref output. Default: False	bool
token_model	Tokenization model. Specify a full path, or just a filename. If a filename is given it is expected to be in the opennlp model directory (models/opennlp/)	string
parser_model	Parser model, full path or directory name. If a filename is given, it is expected to be in the OpenNLP model directory (models/opennlp/)	string
timeout	Timeout in seconds for each individual coref resolution task. If this is exceeded, an InvalidDocument is returned for that document	int
coref_model	Coreference resolution model, full path or directory name. If a filename is given, it is expected to be in the OpenNLP model directory (models/opennlp/). Default: "" (standard English opennlp model in models/opennlp/)	string
readable	If True, pretty-print the JSON output, so it's human-readable. Default: False	bool
pos_model	POS tagger model, full path or filename. If a filename is given, it is expected to be in the opennlp model directory (models/opennlp/)	string
sentence_model	Sentence segmentation model. Specify a full path, or just a filename. If a filename is given it is expected to be in the opennlp model directory (models/opennlp/)	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_opennlp_coref_module]
type=pimlico.modules.opennlp.coreference_pipeline
input_text=module_a.some_output
```

This example usage includes more options.

```
[my_opennlp_coref_module]
type=pimlico.modules.opennlp.coreference_pipeline
input_text=module_a.some_output
gzip=T
token_model=en-token.bin
parse_model=en-parser-chunking.bin
timeout=0
coref_model=
readable=T
pos_model=en-pos-maxent.bin
sentence_model=en-sent.bin
```

!! OpenNLP NER

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.opennlp.ner
Executable	yes

Named-entity recognition using OpenNLP's tools.

By default, uses the pre-trained English model distributed with OpenNLP. If you want to use other models (e.g. for other languages), download them from the OpenNLP website to the models dir (*models/opennlp*) and specify the model name as an option.

Note that the default model is for identifying person names only. You can identify other name types by loading other pre-trained OpenNLP NER models. Identification of multiple name types at the same time is not (yet) implemented.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
text	invalid input type specification

Outputs

Name	Type(s)
documents	invalid output type specification

Options

Name	Description	Type
model	NER model, full path or filename. If a filename is given, it is expected to be in the opennlp model directory (models/opennlp/)	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_opennlp_ner_module]
type=pimlico.modules.opennlp.ner
input_text=module_a.some_output
```

This example usage includes more options.

```
[my_opennlp_ner_module]
type=pimlico.modules.opennlp.ner
input_text=module_a.some_output
model=en-ner-person.bin
```

!! OpenNLP constituency parser

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.opennlp.parse
Executable	yes

Todo: Document this module

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
documents	invalid input type specification or invalid input type specification

Outputs

Name	Type(s)
parser	invalid output type specification

Options

Name	Description	Type
model	Parser model, full path or directory name. If a filename is given, it is expected to be in the OpenNLP model directory (models/opennlp/)	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_opennlp_parser_module]
type=pimlico.modules.opennlp.parse
input_documents=module_a.some_output
```

This example usage includes more options.

```
[my_opennlp_parser_module]
type=pimlico.modules.opennlp.parse
input_documents=module_a.some_output
model=en-parser-chunking.bin
```

!! OpenNLP POS-tagger

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.opennlp.pos
Executable	yes

Part-of-speech tagging using OpenNLP's tools.

By default, uses the pre-trained English model distributed with OpenNLP. If you want to use other models (e.g. for other languages), download them from the OpenNLP website to the models dir (*models/opennlp*) and specify the model name as an option.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
text	invalid input type specification

Outputs

Name	Type(s)
documents	invalid output type specification

Options

Name	Description	Type
model	POS tagger model, full path or filename. If a filename is given, it is expected to be in the opennlp model directory (models/opennlp/)	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_opennlp_pos_tagger_module]
type=pimlico.modules.opennlp.pos
input_text=module_a.some_output
```

This example usage includes more options.

```
[my_opennlp_pos_tagger_module]
type=pimlico.modules.opennlp.pos
input_text=module_a.some_output
model=en-pos-maxent.bin
```

OpenNLP tokenizer

Path	pimlico.modules.opennlp.tokenize
Executable	yes

Sentence splitting and tokenization using OpenNLP's tools.

Sentence splitting may be skipped by setting the option `tokenize_only=T`. The tokenizer will then assume that each line in the input file represents a sentence and tokenize within the lines.

Inputs

Name	Type(s)
text	<i>grouped_corpus</i> <TextDocumentType>

Outputs

Name	Type(s)
documents	<i>grouped_corpus</i> <TokenizedDocumentType>

Options

Name	Description	Type
to-ken_model	Tokenization model. Specify a full path, or just a filename. If a filename is given it is expected to be in the <code>opennlp</code> model directory (<code>models/opennlp/</code>)	string
tokenize_only	By default, sentence splitting is performed prior to tokenization. If <code>tokenize_only</code> is set, only the tokenization step is executed	bool
sentence_model	Sentence segmentation model. Specify a full path, or just a filename. If a filename is given it is expected to be in the <code>opennlp</code> model directory (<code>models/opennlp/</code>)	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_opennlp_tokenizer_module]
type=pimlico.modules.opennlp.tokenize
input_text=module_a.some_output
```

This example usage includes more options.

```
[my_opennlp_tokenizer_module]
type=pimlico.modules.opennlp.tokenize
input_text=module_a.some_output
token_model=en-token.bin
tokenize_only=F
sentence_model=en-sent.bin
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *opennlp_tokenize*

1.3.11 R interfaces

Modules for interfacing with the [statistical programming language R](#). Currently, we provide just a simple way to pass data from the output of another module into an R script and run it. In the future, it may be appropriate to add more sophisticated interfaces, or expose R's functionality in a more specialised way, integrating more closely with Pimlico's datatype system.

!! R script executor

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	<code>pimlico.modules.r.script</code>
Executable	yes

Simple interface to R that just involves running a given R script, first substituting in some paths from the pipeline, making it easy to pass in data from the output of other modules.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
sources	invalid input type specification

Outputs

Name	Type(s)
output	invalid output type specification

Options

Name	Description	Type
script	(required) Path to the script to be run. The script itself may include substitutions of the form ‘{{inputX}}’, which will be replaced with the absolute path to the data dir of the Xth input, and ‘{{output}}’, which will be replaced with the absolute path to the output dir. The latter allows the script to output things other than the output file, which always exists and contains the full script’s output	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_r_script_module]
type=pimlico.modules.r.script
input_sources=module_a.some_output
script=value
```

1.3.12 Regular expressions

!! Regex annotated text matcher

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.regex.annotated_text
Executable	yes

Todo: Document this module

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
documents	invalid input type specification

Outputs

Name	Type(s)
documents	invalid output type specification

Options

Name	Description	Type
expr	(required) An expression to determine what to search for in sentences. Consists of a sequence of tokens, each matching one field in the corresponding token's annotations in the data. These are specified in the form field[x], where field is the name of a field supplied by the input data and x is the value required of that field. If x ends in a *, it will match prefixes: e.g. pos[NN*]. If no field name is given, the default 'word' is used. A token of the form 'x=y' matches the expression y as above and assigns the matching word to the extracted variable x (to be output). You may also extract a different annotation field by specifying x=f:y, where f is the field name to be extracted. E.g. 'what a=lemma:pos[NN*] lemma[come] with b=pos[NN*]' matches phrases like 'what meals come with fries', producing 'a=meal' and 'b=fries'. Both pos and lemma need to be fields in the dataset'. If you give multiple whole expressions separated by s, matches will be collected from all of them	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_annotated_text_matcher_module]
type=pimlico.modules.regex.annotated_text
input_documents=module_a.some_output
expr=value
```

1.3.13 Scikit-learn tools

Scikit-learn ('sklearn') provides easy-to-use implementations of a large number of machine-learning methods, based on Numpy/Scipy.

You can build Numpy arrays from your corpus using the *feature processing tools* and then use them as input to Scikit-learn's tools using the modules in this package.

Sklearn logistic regression

Path	pimlico.modules.sklearn.logistic_regression
Executable	yes

Provides an interface to Scikit-Learn's simple logistic regression trainer.

You may also want to consider using:

- **LogisticRegressionCV**: LR with cross-validation to choose regularization strength
- **SGDClassifier**: general gradient-descent training for classifiers, which includes logistic regression. A better choice for training on a large dataset.

Inputs

Name	Type(s)
features	<i>scored_real_feature_sets</i>

Outputs

Name	Type(s)
model	<i>sklearn_model</i>

Options

Name	Description	Type
options	Options to pass into the constructor of LogisticRegression, formatted as a JSON dictionary (potentially without the {}s). E.g.: "C":1.5, "penalty":"l2"	JSON dict

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_sklearn_log_reg_module]
type=pimlico.modules.sklearn.logistic_regression
input_features=module_a.some_output
```

This example usage includes more options.

```
[my_sklearn_log_reg_module]
type=pimlico.modules.sklearn.logistic_regression
input_features=module_a.some_output
options="C":1.5, "penalty":"l2"
```

!! Sklearn matrix factorization

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.sklearn.matrix_factorization
Executable	yes

Provides a simple interface to [Scikit-Learn](#)'s various matrix factorization models.

Since they provide a consistent training interface, you can simply choose the class name of the method you want to use and specify options relevant to that method in the `options` option. For available options, take a look at the table of parameters in the [Scikit-Learn documentation](#) for each class.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
matrix	invalid input type specification

Outputs

Name	Type(s)
w	invalid output type specification
h	invalid output type specification

Options

Name	Description	Type
class	(required) Scikit-learn class to use to fit the matrix factorization. Should be the name of a class in the package <code>sklearn.decomposition</code> that has a <code>fit_transform()</code> method and a <code>components_</code> attribute. Supported classes: NMF, SparsePCA, ProjectedGradientNMF, FastICA, FactorAnalysis, PCA, RandomizedPCA, LatentDirichletAllocation, TruncatedSVD	'NMF', 'SparsePCA', 'ProjectedGradientNMF', 'FastICA', 'FactorAnalysis', 'PCA', 'RandomizedPCA', 'LatentDirichletAllocation' or 'TruncatedSVD'
options	Options to pass into the constructor of the sklearn class, formatted as a JSON dictionary (potentially without the <code>{}</code> s). E.g.: <code>'n_components=200, solver="cd", tol=0.0001, max_iter=200'</code>	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_sklearn_mat_fac_module]
type=pimlico.modules.sklearn.matrix_factorization
input_matrix=module_a.some_output
class=value
```

This example usage includes more options.

```
[my_sklearn_mat_fac_module]
type=pimlico.modules.sklearn.matrix_factorization
input_matrix=module_a.some_output
class=value
options=value
```

1.3.14 Document-level text filters

Simple text filters that are applied at the document level, i.e. each document in a `TarredCorpus` is processed one at a time. These perform relatively simple processing, not relying on external software or involving lengthy processing times. They are therefore most often used using the `filter=T` option, so that the processing is performed on the fly.

Such filters are needed sometimes just to convert before different datapoint formats.

Probably a good deal of these will be added in due course.

!! Text to character level

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.text.char_tokenize
Executable	yes

Filter to treat text data as character-level tokenized data. This makes it simple to train character-level models, since the output appears exactly like a tokenized document, where each token is a single character. You can then feed it into any module that expects tokenized text.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
corpus	invalid input type specification

Outputs

Name	Type(s)
corpus	invalid output type specification

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_char_tokenize_module]
type=pimlico.modules.text.char_tokenize
input_corpus=module_a.some_output
```

Normalize tokenized text

Path	pimlico.modules.text.normalize
Executable	yes

Perform text normalization on tokenized documents.

Currently, this includes only the following:

- case normalization (to upper or lower case)
- blank line removal
- empty sentence removal

In the future, more normalization operations may be added.

Inputs

Name	Type(s)
corpus	<i>grouped_corpus</i> <TokenizedDocumentType>

Outputs

Name	Type(s)
corpus	<i>grouped_corpus</i> <TokenizedDocumentType>

Options

Name	Description	Type
case	Transform all text to upper or lower case. Choose from 'upper' or 'lower', or leave blank to not perform transformation	'upper', 'lower' or ''
re-move_only_punct	Skip over any sentences that are empty if punctuation is ignored	bool
re-move_empty	Skip over any empty sentences (i.e. blank lines)	bool

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_normalize_module]
type=pimlico.modules.text.normalize
input_corpus=module_a.some_output
```

This example usage includes more options.

```
[my_normalize_module]
type=pimlico.modules.text.normalize
input_corpus=module_a.some_output
case=
remove_only_punct=F
remove_empty=F
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *normalize*

Simple tokenization

Path	pimlico.modules.text.simple_tokenize
Executable	yes

Tokenize raw text using simple splitting.

This is useful where either you don't mind about the quality of the tokenization and just want to test something quickly, or text is actually already tokenized, but stored as a raw text datatype.

If you want to do proper tokenization, consider either the CoreNLP or OpenNLP core modules.

Inputs

Name	Type(s)
corpus	<i>grouped_corpus</i> <TextDocumentType>

Outputs

Name	Type(s)
corpus	<i>grouped_corpus</i> <TokenizedDocumentType>

Options

Name	Description	Type
splitter	Character or string to split on. Default: space	<type 'unicode'>

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_simple_tokenize_module]
type=pimlico.modules.text.simple_tokenize
input_corpus=module_a.some_output
```

This example usage includes more options.

```
[my_simple_tokenize_module]
type=pimlico.modules.text.simple_tokenize
input_corpus=module_a.some_output
splitter=
```

Test pipelines

This module is used by the following *test pipelines*. They are a further source of examples of the module's usage.

- *simple_tokenize*

!! Normalize raw text

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.text.text_normalize
Executable	yes

Text normalization for raw text documents.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
corpus	invalid input type specification

Outputs

Name	Type(s)
corpus	invalid output type specification

Options

Name	Description	Type
case	Transform all text to upper or lower case. Choose from ‘upper’ or ‘lower’, or leave blank to not perform transformation	‘upper’, ‘lower’ or ‘’
blank_lines	Remove all blank lines (after whitespace stripping, if requested)	bool
strip	Strip whitespace from the start and end of lines	bool

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_text_normalize_module]
type=pimlico.modules.text.text_normalize
input_corpus=module_a.some_output
```

This example usage includes more options.

```
[my_text_normalize_module]
type=pimlico.modules.text.text_normalize
input_corpus=module_a.some_output
case=
blank_lines=T
strip=T
```

!! Tokenized text to text

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.text.untokenize
Executable	yes

Filter to take tokenized text and join it together to make raw text.

This module shouldn’t be necessary and will be removed later. For the time being, it’s here as a workaround for [this problem](<https://github.com/markgw/pimlico/issues/1#issuecomment-383620759>), until it’s solved in the datatype re-design.

Tokenized text is a subtype of text, so theoretically it should be acceptable to modules that expect plain text (and is considered so by typechecking). But it provides an incompatible data structure, so things go bad if you use it like that.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
corpus	invalid input type specification

Outputs

Name	Type(s)
corpus	invalid output type specification

Options

Name	Description	Type
sentence_joiner	String to join lines/sentences on. (Default: linebreak)	<type 'unicode'>
joiner	String to join words on. (Default: space)	<type 'unicode'>

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_untokenize_module]
type=pimlico.modules.text.untokenize
input_corpus=module_a.some_output
```

This example usage includes more options.

```
[my_untokenize_module]
type=pimlico.modules.text.untokenize
input_corpus=module_a.some_output
sentence_joiner=

joiner=
```

1.3.15 General utilities

General utilities for things like filesystem manipulation.

!! Module output alias

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.utility.alias
Executable	no

Alias a datatype coming from the output of another module.

Used to assign a handy identifier to the output of a module, so that we can just refer to this alias module later in the pipeline and use its default output. This can help make for a more readable pipeline config.

For example, say we use *split* to split a dataset into two random subsets. The two splits can be accessed by referring to the two outputs of that module: *split_module.set1* and *split_module.set2*. However, it's easy to lose track of what these splits are supposed to be used for, so we might want to give them names:

```
[split_module]
type=pimlico.modules.corpora.split
set1_size=0.2

[test_set]
type=pimlico.modules.utility.alias
input=split_module.set1

[training_set]
type=pimlico.modules.utility.alias
input=split_module.set2

[training_routine]
type=...
input_corpus=training_set
```

Note the difference between using this module and using the special *alias* module type. The *alias* type creates an alias for a whole module, allowing you to refer to all of its outputs, inherit its settings, and anything else you could do with the original module name. This module, however, provides an alias for exactly one output of a module and generates a module instance of its own in the pipeline (albeit a filter module).

Todo: Update to new datatypes system and add test pipeline

This is a filter module. It is not executable, so won't appear in a pipeline's list of modules that can be run. It produces its output for the next module on the fly when the next module needs it.

Inputs

Name	Type(s)
input	invalid input type specification

Outputs

Name	Type(s)
output	invalid output type specification

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_alias_module]
type=pimlico.modules.utility.alias
input_input=module_a.some_output
```

!! Collect files

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.utility.collect_files
Executable	yes

Collect files output from different modules.

A simple convenience module to make it easier to inspect output by putting it all in one place.

Files are either collected into subdirectories or renamed to avoid clashes.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
files	<i>list</i> of invalid input type specification

Outputs

Name	Type(s)
files	<i>collected_named_file_collection</i>

Options

Name	Description	Type
sub-dirs	Use subdirectories to collect the files from different sources, rather than renaming each file. By default, a prefix is added to the filenames	bool
names	List of string identifiers to use to distinguish the files from different sources, either used as subdirectory names or filename prefixes. If not given, integer ids will be used instead	comma-separated list of strings

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_collect_files_module]
```

<<<<<< HEAD

```
type=pimlico.modules.utility.collect_files
```

```
>>>>>> 238f254c3e241cae81cba8fea74a0090eeafb35d input_files=module_a.some_output
```

This example usage includes more options.

```
[my_collect_files_module]
```

<<<<<< HEAD

```
type=pimlico.modules.utility.collect_files
```

```
>>>>>> 238f254c3e241cae81cba8fea74a0090eeafb35d input_files=module_a.some_output      subdirs=T
names=text,text,...
```

!! Copy file

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.utility.copy_file
Executable	yes

Copy a file

Simple utility for copying a file (which presumably comes from the output of another module) into a particular location. Useful for collecting together final output at the end of a pipeline.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
source	invalid input type specification

Outputs

No outputs Options =====

Name	Description	Type
tar-get_name	Name to rename the target file to. If not given, it will have the same name as the source file. Ignored if there's more than one input file	string
tar-get_dir	(required) Path to directory into which the file should be copied. Will be created if it doesn't exist	string

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_copy_file_module]
type=pimlico.modules.utility.copy_file
input_source=module_a.some_output
target_dir=value
```

This example usage includes more options.

```
[my_copy_file_module]
type=pimlico.modules.utility.copy_file
input_source=module_a.some_output
target_name=value
target_dir=value
```

1.3.16 Visualization tools

Modules for plotting and suchlike

!! Bar chart plotter

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.visualization.bar_chart
Executable	yes

Simple plotting of a bar chart from numeric data using Matplotlib

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
values	invalid input type specification

Outputs

Name	Type(s)
plot	invalid output type specification

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_bar_chart_module]
type=pimlico.modules.visualization.bar_chart
input_values=module_a.some_output
```

!! Embedding space plotter

Note: This module has not yet been updated to the new datatype system, so cannot be used in the *datatypes* branch. Soon it will be updated.

Path	pimlico.modules.visualization.embeddings_plot
Executable	yes

Plot vectors from embeddings, trained by some other module, in a 2D space using a MDS reduction and Matplotlib.

They might, for example, come from `pimlico.modules.embeddings.word2vec`. The embeddings are read in using Pimlico's generic word embedding storage type.

Uses scikit-learn to perform the MDS/TSNE reduction.

Todo: Update to new datatypes system and add test pipeline

Inputs

Name	Type(s)
vectors	invalid input type specification

Outputs

Name	Type(s)
plot	invalid output type specification

Options

Name	Description	Type
skip	Number of most frequent words to skip, taking the next most frequent after these. Default: 0	int
metric	Distance metric to use. Choose from 'cosine', 'euclidean', 'manhattan'. Default: 'cosine'	'cosine', 'euclidean' or 'manhattan'
reduction	Dimensionality reduction technique to use to project to 2D. Available: mds (Multi-dimensional Scaling), tsne (t-distributed Stochastic Neighbor Embedding). Default: mds	'mds' or 'tsne'
colors	List of colours to use for different embedding sets. Should be a list of matplotlib colour strings, one for each embedding set given in input_vectors	comma-separated list of strings
cmap	Mapping from word prefixes to matplotlib plotting colours. Every word beginning with the given prefix has the prefix removed and is plotted in the corresponding colour. Specify as a JSON dictionary mapping prefix strings to colour strings	JSON string
words	Number of most frequent words to plot. Default: 50	int

Example config

This is an example of how this module can be used in a pipeline config file.

```
[my_embeddings_plot_module]
type=pimlico.modules.visualization.embeddings_plot
input_vectors=module_a.some_output
```

This example usage includes more options.

```
[my_embeddings_plot_module]
type=pimlico.modules.visualization.embeddings_plot
input_vectors=module_a.some_output
skip=0
metric=cosine
reduction=mds
colors=text,text,...
cmap={"key1":"value"}
words=50
```

1.4 Command-line interface

The main Pimlico command-line interface (usually accessed via *pimlico.sh* in your project root) provides subcommands to perform different operations. Call it like so, using one of the subcommands documented below to access particular functionality:

```
./pimlico.sh <config-file> [general options...] <subcommand> [subcommand args/options]
```

The commands you are likely to use most often are: *status*, *run*, *reset* and maybe *browse*.

For a reference for each command's options, see the command-line documentation: `./pimlico.sh --help`, for a general reference and `./pimlico.sh <config_file> <command> --help` for a specific subcommand's reference.

Below is a more detailed guide for each subcommand, including all of the documentation available via the command line.

<i>browse</i>	View the data output by a module
<i>clean</i>	Remove all module output directories that do not correspond to a module in the pipeline
<i>deps</i>	List information about software dependencies: whether they're available, versions, etc
<i>dump</i>	Dump the entire available output data from a given pipeline module to a tarball
<i>email</i>	Test email settings and try sending an email using them
<i>inputs</i>	Show the (expected) locations of the inputs of a given module
<i>install</i>	Install missing module library dependencies
<i>load</i>	Load a module's output data from a tarball previously created by the dump command
<i>movestores</i>	Move data between stores
<i>newmodule</i>	Create a new module type
<i>output</i>	Show the location where the given module's output data will be (or has been) stored
<i>python</i>	Load the pipeline config and enter a Python interpreter with access to it in the environment
<i>reset</i>	Delete any output from the given module and restore it to unexecuted state
<i>run</i>	Execute an individual pipeline module, or a sequence
<i>shell</i>	Open a shell to give access to the data output by a module
<i>status</i>	Output a module execution schedule for the pipeline and execution status for every module
<i>stores</i>	List named Pimlico stores
<i>unlock</i>	Forcibly remove an execution lock from a module
<i>variants</i>	List the available variants of a pipeline config
<i>visualize</i>	Comming soon... visualize the pipeline in a pretty way

1.4.1 status

Command-line tool subcommand

Output a module execution schedule for the pipeline and execution status for every module.

Usage:

```
pimlico.sh [...] status [module_name] [-h] [--all] [--short] [--history] [--deps-of_↵
↵DEPS_OF] [--no-color]
```

Positional arguments

Arg	Description
[module]	Optionally specify a module name (or number). More detailed status information will be outut for this module. Alternatively, use this arg to limit the modules whose status will be output to a range by specifying 'A...B', where A and B are module names or numbers

Options

Option	Description
<code>--all</code> , <code>-a</code>	Show all modules defined in the pipeline, not just those that can be executed
<code>--short</code> , <code>-s</code>	Use a brief format when showing the full pipeline's status. Only applies when module names are not specified. This is useful with very large pipelines, where you just want a compact overview of the status
<code>--history</code> , <code>-i</code>	When a module name is given, even more detailed output is given, including the full execution history of the module
<code>--deps-only</code> , <code>-d</code>	Restrict to showing only the named/numbered module and any that are (transitive) dependencies of it. That is, show the whole tree of modules that lead through the pipeline to the given module
<code>--no-color</code> , <code>--nc</code>	Don't include terminal color characters, even if the terminal appears to support them. This can be useful if the automatic detection of color terminals doesn't work and the status command displays lots of horrible escape characters

1.4.2 variants

Command-line tool subcommand

List the available variants of a pipeline config

See *Pipeline variants* for more details.

Usage:

```
pimlico.sh [...] variants [-h]
```

1.4.3 run

Command-line tool subcommand

Main command for executing Pimlico modules from the command line *run* command.

Usage:

```
pimlico.sh [...] run [modules [modules ...]] [-h] [--force-rerun] [--all-deps] [--  
↪all] [--dry-run] [--step] [--preliminary] [--exit-on-error] [--email {modend,end}]
```

Positional arguments

Arg	Description
[modules [modules ...]]	The name (or number) of the module to run. To run a stage from a multi-stage module, use 'module:stage'. Use 'status' command to see available modules. Use 'module:?' or 'module:help' to list available stages. If not given, defaults to next incomplete module that has all its inputs ready. You may give multiple modules, in which case they will be executed in the order specified

Options

Option	Description
<code>--force-run</code> <code>-f</code>	Force, running the module(s), even if it's already been run to completion
<code>--all-dep</code> <code>-a</code>	If the given module(s) has dependent modules that have not been completed, executed them first. This allows you to specify a module late in the pipeline and execute the full pipeline leading to that point
<code>--all</code>	Run all currently unexecuted modules that have their inputs ready, or will have by the time previous modules are run. (List of modules will be ignored)
<code>--dry-run</code> <code>--dry</code> <code>--check</code>	Perform all pre-execution checks, but don't actually run the module(s)
<code>--step</code>	Enabled super-verbose debugging mode, which steps through a module's processing outputting a lot of information and allowing you to control the output as it goes. Useful for working out what's going on inside a module if it's mysteriously not producing the output you expected
<code>--preliminary</code> <code>--pre</code>	Perform a preliminary run of any modules that take multiple datasets into one of their inputs. This means that we will run the module even if not all the datasets are yet available (but at least one is) and mark it as preliminarily completed
<code>--exit-on-error</code> <code>-e</code>	If an error is encountered while executing a module that causes the whole module execution to fail, output the error and exit. By default, Pimlico will send error output to a file (or print it in debug mode) and continue to execute the next module that can be executed, if any
<code>--email</code>	Send email notifications when processing is complete, including information about the outcome. Choose from: 'modend' (send notification after module execution if it fails and a summary at the end of everything), 'end' (send only the final summary). Email sending must be configured: see 'email' command to test

1.4.4 browse

Command-line tool subcommand

View the data output by a module.

Usage:

```
pimlico.sh [...] browse module_name [output_name] [-h] [--skip-invalid] [--formatter_↵
↵FORMATTER]
```

Positional arguments

Arg	Description
<code>module_name</code>	The name (or number) of the module whose output to look at. Use 'module:stage' for multi-stage modules
<code>[output_name]</code>	The name of the output from the module to browse. If blank, load the default output

Options

Option	Description
<code>--skip-invalid</code>	Skip over invalid documents, instead of showing the error that caused them to be invalid
<code>--formatter</code> <code>-f</code>	When browsing iterable corpora, fully qualified class name of a subclass of <code>DocumentBrowserFormatter</code> to use to determine what to output for each document. You may also choose from the named standard formatters for the datatype in question. Use <code>'-f help'</code> to see a list of available formatters

1.4.5 shell

Command-line tool subcommand

Open a shell to give access to the data output by a module.

Usage:

```
pimlico.sh [...] shell module_name [output_name] [-h]
```

Positional arguments

Arg	Description
<code>module_name</code>	The name (or number) of the module whose output to look at
<code>[output_name]</code>	The name of the output from the module to browse. If blank, load the default output

1.4.6 python

Command-line tool subcommand

Load the pipeline config and enter a Python interpreter with access to it in the environment.

Usage:

```
pimlico.sh [...] python [script] [-h] [-i]
```

Positional arguments

Arg	Description
<code>[script]</code>	Script file to execute. Omit to enter interpreter

Options

Option	Description
<code>-i</code>	Enter interactive shell after running script

1.4.7 reset

Command-line tool subcommand

Delete any output from the given module and restore it to unexecuted state.

Usage:

```
pimlico.sh [...] reset [modules [modules ...]] [-h] [-n]
```

Positional arguments

Arg	Description
[modules [modules ...]]	The names (or numbers) of the modules to reset, or 'all' to reset the whole pipeline

Options

Option	Description
-n,	Only reset the state of this module, even if it has dependent modules in an executed state, which could be invalidated by resetting and re-running this one
--no-deps	

1.4.8 clean

Command-line tool subcommand

Cleans up module output directories that have got left behind.

Often, when developing a pipeline incrementally, you try out some modules, but then remove them, or rename them to something else. The directory in the Pimlico output store that was created to contain their metadata, status and output data is then left behind and no longer associated with any module.

Run this command to check all storage locations for such directories. If it finds any, it prompts you to confirm before deleting them. (If there are things in the list that don't look like they were left behind by the sort of things mentioned above, don't delete them! I don't want you to lose your precious output data if I've made a mistake in this command.)

Note that the operation of this command is specific to the loaded pipeline variant. If you have multiple variants, make sure to select the one you want to clean with the general *-variant* option.

Usage:

```
pimlico.sh [...] clean [-h]
```

1.4.9 stores

Command-line tool subcommand

List Pimlico stores in use and the corresponding storage locations.

Usage:

```
pimlico.sh [...] stores [-h]
```

1.4.10 movestores

Command-line tool subcommand

Move a particular module's output from one storage location to another.

Usage:

```
pimlico.sh [...] movestores dest [modules [modules ...]] [-h]
```

Positional arguments

Arg	Description
dest	Name of destination store
[modules [modules ...]]	The names (or numbers) of the module whose output to move

1.4.11 unlock

Command-line tool subcommand

Forcibly remove an execution lock from a module. If a lock has ended up getting left on when execution exited prematurely, use this to remove it.

When a module starts running, it is locked to avoid making a mess of your output data by running the same module from another terminal, or some other silly mistake (I know, for some of us this sort of behaviour is frustratingly common).

Usually shouldn't be necessary, even if there's an error during execution, since the module should be unlocked when Pimlico exits, but occasionally (e.g. if you have to forcibly kill Pimlico during execution) the lock gets left on.

Usage:

```
pimlico.sh [...] unlock module_name [-h]
```

Positional arguments

Arg	Description
module_name	The name (or number) of the module to unlock

1.4.12 dump

Command-line tool subcommand

Dump the entire available output data from a given pipeline module to a tarball, so that it can easily be loaded into the same pipeline on another system. This is primarily to support spreading the execution of a pipeline between multiple machines, so that the output from a module can easily be transferred and loaded into a pipeline.

Dump to a tarball using this command, transfer the file between machines and then run the *load command* to import it there.

See also:

Running one pipeline on multiple computers: for a more detailed guide to transferring data across servers.

Usage:

```
pimlico.sh [...] dump [modules [modules ...]] [-h] [--output OUTPUT] [--inputs]
```

Positional arguments

Arg	Description
[modules [modules ...]]	Names or numbers of modules whose data to dump. If multiple are given, a separate file will be dumped for each

Options

Option	Description
--output -o	Path to directory to output to. Defaults to the current user's home directory
--inputs -i	Dump data for the modules corresponding to the inputs of the named modules, instead of those modules themselves. Useful for when you're preparing to run a module on a different machine, for getting all the necessary input data for a module

1.4.13 load

Command-line tool subcommand

Load the output data for a given pipeline module from a tarball previously created by the *dump* command (typically on another machine). This is primarily to support spreading the execution of a pipeline between multiple machines, so that the output from a module can easily be transferred and loaded into a pipeline.

Dump to a tarball using the *dump command*, transfer the file between machines and then run this command to import it there.

See also:

Running one pipeline on multiple computers: for a more detailed guide to transferring data across servers.

Usage:

```
pimlico.sh [...] load [paths [paths ...]] [-h] [--force-overwrite]
```

Positional arguments

Arg	Description
[paths [paths ...]]	Paths to dump files (tarballs) to load into the pipeline

Options

Option	Description
--force-overwrite -f	If data already exists for a module being imported, overwrite without asking. By default, the user will be prompted to check whether they want to overwrite

1.4.14 deps

Command-line tool subcommand

Output information about module dependencies.

Usage:

```
pimlico.sh [...] deps [modules [modules ...]] [-h]
```

Positional arguments

Arg	Description
[modules [modules ...]]	Check dependencies for named modules and install any that are automatically installable. Use 'all' to install dependencies for all modules

1.4.15 install

Command-line tool subcommand

Install missing dependencies.

Usage:

```
pimlico.sh [...] install [modules [modules ...]] [-h] [--trust-downloaded]
```

Positional arguments

Arg	Description
[modules [modules ...]]	Check dependencies for named modules and install any that are automatically installable. Use 'all' to install dependencies for all modules

Options

Option	Description
--trust-downloaded	If an archive file to be downloaded is found to be in the lib dir already, trust that it is the file we're after. By default, we only reuse archives we've just downloaded, so we know they came from the right URL, avoiding accidental name clashes
-t	

1.4.16 inputs

Command-line tool subcommand

Show the locations of the inputs of a given module. If the input datasets are available, their actual location is shown. Otherwise, all directories in which the data is being checked for are shown.

Usage:

```
pimlico.sh [...] inputs module_name [-h]
```

Positional arguments

Arg	Description
module_name	The name (or number) of the module to display input locations for

1.4.17 output

Command-line tool subcommand

Show the location where the given module's output data will be (or has been) stored.

Usage:

```
pimlico.sh [...] output module_name [-h]
```

Positional arguments

Arg	Description
module_name	The name (or number) of the module to display input locations for

1.4.18 newmodule

Command-line tool subcommand

Interactive tool to create a new module type, generating a skeleton for the module's code. Currently only works for certain module types. May be extended in future to help with creating a broader range of sorts of modules.

Usage:

```
pimlico.sh [...] newmodule [-h]
```

1.4.19 visualize

Command-line tool subcommand

(Not yet fully implemented!) Visualize the pipeline, with status information for modules.

Usage:

```
pimlico.sh [...] visualize [-h] [--all]
```

Options

Option	Description
--all, -a	Show all modules defined in the pipeline, not just those that can be executed

1.4.20 email

Command-line tool subcommand

Test email settings and try sending an email using them.

Usage:

```
pimlico.sh [...] email [-h]
```

1.5 API Documentation

API documentation for the main Pimlico codebase, excluding the *built-in Pimlico module types*.

1.5.1 pimlico package

Subpackages

pimlico.cli package

Subpackages

pimlico.cli.browser package

Subpackages

pimlico.cli.browser.tools package

Submodules

pimlico.cli.browser.tools.corpus module

Browser tool for iterable corpora.

browse_data (*reader, formatter, skip_invalid=False*)

class CorpusState (*corpus*)

Bases: object

Keep track of which document we're on.

next_document ()

skip (*n*)

class InputDialog (*text, input_edit*)

Bases: urwid.widget.WidgetWrap

A dialog that appears with an input

signals = ['close', 'cancel']

keypress (*size, k*)

class MessageDialog (*text, default=None*)

Bases: `urwid.widget.WidgetWrap`

A dialog that appears with a message

class InputPopupLauncher (*original_widget, text, input_edit, callback=None*)

Bases: `urwid.wimp.PopUpLauncher`

create_pop_up ()

Subclass must override this method and return a widget to be used for the pop-up. This method is called once each time the pop-up is opened.

get_pop_up_parameters ()

Subclass must override this method and have it return a dict, eg:

```
{ 'left':0, 'top':1, 'overlay_width':30, 'overlay_height':4 }
```

This method is called each time this widget is rendered.

skip_popup_launcher (*original_widget, text, default=None, callback=None*)

save_popup_launcher (*original_widget, text, default=None, callback=None*)

class MessagePopupLauncher (*original_widget, text*)

Bases: `urwid.wimp.PopUpLauncher`

create_pop_up ()

Subclass must override this method and return a widget to be used for the pop-up. This method is called once each time the pop-up is opened.

get_pop_up_parameters ()

Subclass must override this method and have it return a dict, eg:

```
{ 'left':0, 'top':1, 'overlay_width':30, 'overlay_height':4 }
```

This method is called each time this widget is rendered.

pimlico.cli.browser.tools.files module

browse_files (*reader*)

Browser tool for NamedFileCollections.

is_binary_string (*bytes*)

is_binary_file (*path*)

Try reading a bit of a file to work out whether it's a binary file or text

pimlico.cli.browser.tools.formatter module

The command-line iterable corpus browser displays one document at a time. It can display the raw data from the corpus files, which sometimes is sufficiently human-readable to not need any special formatting. It can also parse the data using its datatype and output text either from the datatype's standard unicode representation or, if the document datatype provides it, a special browser formatting of the data.

When viewing output data, particularly during debugging of modules, it can be useful to provide special formatting routines to the browser, rather than using or overriding the datatype's standard formatting methods. For example, you might want to pull out specific attributes for each document to get an overview of what's coming out.

The browser command accepts a command-line option that specifies a Python class to format the data. This class should be a subclass of `:class:~pimlico.cli.browser.formatter.DocumentBrowserFormatter` that accepts a datatype compatible with the datatype being browsed and provides a method to format each document. You can write these in your custom code and refer to them by their fully qualified class name.

class DocumentBrowserFormatter (*corpus_datatype*)

Bases: `object`

Base class for formatters used to post-process documents for display in the iterable corpus browser.

DATATYPE = DataPointType()

format_document (*doc*)

Format a single document and return the result as a string (or unicode, but it will be converted to ASCII for display).

Must be overridden by subclasses.

filter_document (*doc*)

Each doc is passed through this function directly after being read from the corpus. If None is returned, the doc is skipped. Otherwise, the result is used instead of the doc data. The default implementation does nothing.

class DefaultFormatter (*corpus_datatype*)

Bases: `pimlico.cli.browser.tools.formatter.DocumentBrowserFormatter`

Generic implementation of a browser formatter that's used if no other formatter is given.

DATATYPE = DataPointType()

format_document (*doc*)

Format a single document and return the result as a string (or unicode, but it will be converted to ASCII for display).

Must be overridden by subclasses.

class InvalidDocumentFormatter (*corpus_datatype*)

Bases: `pimlico.cli.browser.tools.formatter.DocumentBrowserFormatter`

Formatter that skips over all docs other than invalid results. Uses standard formatting for InvalidDocument information.

format_document (*doc*)

Format a single document and return the result as a string (or unicode, but it will be converted to ASCII for display).

Must be overridden by subclasses.

filter_document (*doc*)

Each doc is passed through this function directly after being read from the corpus. If None is returned, the doc is skipped. Otherwise, the result is used instead of the doc data. The default implementation does nothing.

typecheck_formatter (*formatted_doc_type, formatter_cls*)

Check that a document type is compatible with a particular formatter.

load_formatter (*datatype, formatter_name=None*)

Load a formatter specified by its fully qualified Python class name. If None, loads the default formatter. You may also specify a formatter by name, choosing from one of the standard ones that the formatted datatype gives.

Parameters

- **formatter_name** – class name, or class

- **datatype** – dataset that will be formatted
- **parse** – only used if the default formatter is loaded, determines *raw_data* (= *not parse*)

Returns instantiated formatter

Module contents

Submodules

pimlico.cli.browser.tool module

Tool for browsing datasets, reading from the data output by pipeline modules.

browse_cmd (*pipeline, opts*)
Command for main Pimlico CLI

Module contents

pimlico.cli.debug package

Submodules

pimlico.cli.debug.stepper module

class Stepper

Bases: `object`

Type that stores the state of the stepping process. This allows information and parameters to be passed around through the process and updated as we go. For example, if particular type of output is disabled by the user, a parameter can be updated here so we know not to output it later.

enable_step_for_pipeline (*pipeline*)

Prepares a pipeline to run in step mode, modifying modules and wrapping methods to supply the extra functionality.

This approach means that we don't have to consume extra computation time checking whether step mode is enabled during normal runs.

Parameters *pipeline* – instance of PipelineConfig

instantiate_output_datatype_decorator (*instantiate_output_datatype, module_name, output_names, stepper*)

wrap_tarred_corpus (*dtype, module_name, output_name, stepper*)

archive_iter_decorator (*archive_iter, module_name, output_name, stepper*)

get_input_decorator (*get_input, module_name, stepper*)

Decorator to wrap a module info's `get_input()` method so when know where inputs are being used.

option_message (*message_lines, stepper, options=None, stack_trace_option=True, category=None*)

Module contents

Extra-verbose debugging facility

Tools for very slowly and verbosely stepping through the processing that a given module does to debug it.

Enabled using the `-step` switch to the run command.

fmt_frame_info (*info*)

output_stack_trace (*frame=None*)

pimlico.cli.shell package

Submodules

pimlico.cli.shell.base module

class ShellCommand

Bases: `object`

Base class used to provide commands for exploring a particular datatype. A basic set of commands is provided for all datatypes, but specific datatype classes may provide their own, by overriding the `shell_commands` attribute.

commands = []

help_text = None

execute (*shell, *args, **kwargs*)

Execute the command. Get the dataset reader as `shell.data`.

Parameters

- **shell** – DataShell instance. Reader available as `shell.data`
- **args** – Args given by the user
- **kwargs** – Named args given by the user as `key=val`

class DataShell (*data, commands, *args, **kwargs*)

Bases: `cmd.Cmd`

Terminal shell for querying datatypes.

prompt = '>>> '

get_names ()

do_EOF (*line*)

Exits the shell

preloop ()

postloop ()

emptyline ()

Don't repeat the last command (default): ignore empty lines

default (*line*)

We use this to handle commands that can't be handled using the `do_` pattern. Also handles the default fallback, which is to execute Python.

`cmdloop` (*intro=None*)

exception ShellError

Bases: `exceptions.Exception`

pimlico.cli.shell.commands module

Basic set of shell commands that are always available.

class MetadataCmd

Bases: `pimlico.cli.shell.base.ShellCommand`

`commands` = ['metadata']

`help_text` = "Display the loaded dataset's metadata"

`execute` (*shell, *args, **kwargs*)

Execute the command. Get the dataset reader as shell.data.

Parameters

- `shell` – DataShell instance. Reader available as shell.data
- `args` – Args given by the user
- `kwargs` – Named args given by the user as key=val

class PythonCmd

Bases: `pimlico.cli.shell.base.ShellCommand`

`commands` = ['python', 'py']

`help_text` = "Run a Python interpreter using the current environment, including import ."

`execute` (*shell, *args, **kwargs*)

Execute the command. Get the dataset reader as shell.data.

Parameters

- `shell` – DataShell instance. Reader available as shell.data
- `args` – Args given by the user
- `kwargs` – Named args given by the user as key=val

pimlico.cli.shell.runner module

class ShellCLICmd

Bases: `pimlico.cli.subcommands.PimlicoCLISubcommand`

`command_name` = 'shell'

`command_help` = 'Open a shell to give access to the data output by a module'

`add_arguments` (*parser*)

`run_command` (*pipeline, opts*)

`launch_shell` (*data*)

Starts a shell to view and query the given datatype instance.

Module contents

Submodules

pimlico.cli.check module

class InstallCmd

Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

Install missing dependencies.

command_name = 'install'

command_help = 'Install missing module library dependencies'

add_arguments (*parser*)

run_command (*pipeline, opts*)

class DepsCmd

Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

Output information about module dependencies.

command_name = 'deps'

command_help = "List information about software dependencies: whether they're available"

add_arguments (*parser*)

run_command (*pipeline, opts*)

pimlico.cli.clean module

class CleanCmd

Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

Cleans up module output directories that have got left behind.

Often, when developing a pipeline incrementally, you try out some modules, but then remove them, or rename them to something else. The directory in the Pimlico output store that was created to contain their metadata, status and output data is then left behind and no longer associated with any module.

Run this command to check all storage locations for such directories. If it finds any, it prompts you to confirm before deleting them. (If there are things in the list that don't look like they were left behind by the sort of things mentioned above, don't delete them! I don't want you to lose your precious output data if I've made a mistake in this command.)

Note that the operation of this command is specific to the loaded pipeline variant. If you have multiple variants, make sure to select the one you want to clean with the general *-variant* option.

command_name = 'clean'

command_help = 'Remove all module directories that do not correspond to a module in the pipeline'

command_desc = 'Remove all module output directories that do not correspond to a module in the pipeline'

run_command (*pipeline, opts*)

pimlico.cli.loaddump module

class DumpCmd

Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

Dump the entire available output data from a given pipeline module to a tarball, so that it can easily be loaded into the same pipeline on another system. This is primarily to support spreading the execution of a pipeline between multiple machines, so that the output from a module can easily be transferred and loaded into a pipeline.

Dump to a tarball using this command, transfer the file between machines and then run the *load command* to import it there.

See also:

Running one pipeline on multiple computers: for a more detailed guide to transferring data across servers

```
command_name = 'dump'
```

```
command_help = 'Dump the entire available output data from a given pipeline module to a tarball'
```

```
command_desc = 'Dump the entire available output data from a given pipeline module to a tarball'
```

```
add_arguments(parser)
```

```
run_command(pipeline, opts)
```

class LoadCmd

Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

Load the output data for a given pipeline module from a tarball previously created by the *dump* command (typically on another machine). This is primarily to support spreading the execution of a pipeline between multiple machines, so that the output from a module can easily be transferred and loaded into a pipeline.

Dump to a tarball using the *dump command*, transfer the file between machines and then run this command to import it there.

See also:

Running one pipeline on multiple computers: for a more detailed guide to transferring data across servers

```
command_name = 'load'
```

```
command_help = "Load a module's output data from a tarball previously created by the dump command"
```

```
command_desc = "Load a module's output data from a tarball previously created by the dump command"
```

```
add_arguments(parser)
```

```
run_command(pipeline, opts)
```

pimlico.cli.locations module

class InputsCmd

Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

```
command_name = 'inputs'
```

```
command_help = 'Show the locations of the inputs of a given module. If the input data is not available, show the (expected) locations of the inputs of a given module'
```

```
command_desc = 'Show the (expected) locations of the inputs of a given module'
```

```
add_arguments(parser)
```

```
run_command(pipeline, opts)
```

class OutputCmdBases: *pimlico.cli.subcommands.PimlicoCLISubcommand***command_name** = 'output'**command_help** = "Show the location where the given module's output data will be (or has**add_arguments** (*parser*)**run_command** (*pipeline, opts*)**class ListStoresCmd**Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand***command_name** = 'stores'**command_help** = 'List Pimlico stores in use and the corresponding storage locations'**command_desc** = 'List named Pimlico stores'**run_command** (*pipeline, opts*)**class MoveStoresCmd**Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand***command_name** = 'movestores'**command_help** = "Move a particular module's output from one storage location to another"**command_desc** = 'Move data between stores'**add_arguments** (*parser*)**run_command** (*pipeline, opts*)**pimlico.cli.main module**Main command-line script for running Pimlico, typically called from *pimlico.sh*.

Provides access to many subcommands, acting as the primary interface to Pimlico's functionality.

class VariantsCmdBases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

List the available variants of a pipeline config

See *Pipeline variants* for more details.**command_name** = 'variants'**command_help** = 'List the available variants of a pipeline config'**add_arguments** (*parser*)**run_command** (*pipeline, opts*)**class UnlockCmd**Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

Forcibly remove an execution lock from a module. If a lock has ended up getting left on when execution exited prematurely, use this to remove it.

When a module starts running, it is locked to avoid making a mess of your output data by running the same module from another terminal, or some other silly mistake (I know, for some of us this sort of behaviour is frustratingly common).

Usually shouldn't be necessary, even if there's an error during execution, since the module should be unlocked when Pimlico exits, but occasionally (e.g. if you have to forcibly kill Pimlico during execution) the lock gets left on.

```
command_name = 'unlock'
command_help = "Forcibly remove an execution lock from a module. If a lock has ended up
command_desc = 'Forcibly remove an execution lock from a module'
add_arguments(parser)
run_command(pipeline, opts)
```

class BrowseCmd

Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

```
command_name = 'browse'
command_help = 'View the data output by a module'
add_arguments(parser)
run_command(pipeline, opts)
```

class VisualizeCmd

Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

```
command_name = 'visualize'
command_help = '(Not yet fully implemented!) Visualize the pipeline, with status inform
command_desc = 'Comming soon...visualize the pipeline in a pretty way'
add_arguments(parser)
run_command(pipeline, opts)
```

pimlico.cli.newmodule module

class NewModuleCmd

Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*

```
command_name = 'newmodule'
command_help = "Interactive tool to create a new module type, generating a skeleton fo
command_desc = 'Create a new module type'
run_command(pipeline, opts)
```

ask (*prompt*, *strip_space=True*)

pimlico.cli.pyshell module

class PimlicoPythonShellContext

Bases: *object*

A class used as a static global data structure to provide access to the loaded pipeline when running the Pimlico Python shell command.

This should never be used in any other context to pass around loaded pipelines or other global data. We don't do that sort of thing.

class PythonShellCmdBases: *pimlico.cli.subcommands.PimlicoCLISubcommand***command_name** = 'python'**command_help** = 'Load the pipeline config and enter a Python interpreter with access to**add_arguments** (*parser*)**run_command** (*pipeline, opts*)**get_pipeline** ()

This function may be used in scripts that are expected to be run exclusively from the Pimlico Python shell command (`python`) to get hold of the pipeline that was specified on the command line and loaded when the shell was started.

exception ShellContextErrorBases: *exceptions.Exception***pimlico.cli.reset module****class ResetCmd**Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand***command_name** = 'reset'**command_help** = 'Delete any output from the given module and restore it to unexecuted s**add_arguments** (*parser*)**run_command** (*pipeline, opts*)**pimlico.cli.run module****class RunCmd**Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand*Main command for executing Pimlico modules from the command line *run* command.**command_name** = 'run'**command_help** = 'Execute an individual pipeline module, or a sequence'**add_arguments** (*parser*)**run_command** (*pipeline, opts*)**pimlico.cli.status module****class StatusCmd**Bases: *pimlico.cli.subcommands.PimlicoCLISubcommand***command_name** = 'status'**command_help** = 'Output a module execution schedule for the pipeline and execution stat**add_arguments** (*parser*)**run_command** (*pipeline, opts*)**module_status_color** (*module*)

status_colored (*module*, *text=None*)

Colour the text according to the status of the given module. If text is not given, the module's name is returned.

module_status (*module*)

Detailed module status, shown when a specific module's status is requested.

pimlico.cli.subcommands module

class PimlicoCLISubcommand

Bases: `object`

Base class for defining subcommands to the main command line tool.

This allows us to split up subcommands, together with all their arguments/options and their functionality, since there are quite a lot of them.

Documentation of subcommands should be supplied in the following ways:

- Include help texts for positional args and options in the `add_arguments()` method. They will all be included in the doc page for the command.
- Write a very short description of what the command is for (a few words) in `command_desc`. This will be used in the summary table / TOC in the docs.
- Write a short description of what the command does in `command_help`. This will be available in command-line help and used as a fallback if you don't do the next point.
- Write a good guide to using the command (or at least say what it does) in the class' docstring (i.e. overriding this). This will form the bulk of the command's doc page.

`command_name = None`

`command_help = None`

`command_desc = None`

`add_arguments` (*parser*)

`run_command` (*pipeline*, *opts*)

pimlico.cli.testemail module

class EmailCmd

Bases: `pimlico.cli.subcommands.PimlicoCLISubcommand`

`command_name = 'email'`

`command_help = 'Test email settings and try sending an email using them'`

`run_command` (*pipeline*, *opts*)

pimlico.cli.util module

module_number_to_name (*pipeline*, *name*)

module_numbers_to_names (*pipeline*, *names*)

Convert module numbers to names, also handling ranges of numbers (and names) specified with "...". Any "..." will be filled in by the sequence of intervening modules.

Also, if an unexpanded module name is specified for a module that's been expanded into multiple corresponding to alternative parameters, all of the expanded module names are inserted in place of the unexpanded name.

format_execution_error (*error*)

Produce a string with lots of error output to help debug a module execution error.

Parameters **error** – the exception raised (ModuleExecutionError or ModuleInfoLoadError)

Returns formatted output

print_execution_error (*error*)

Module contents

pimlico.core package

Subpackages

pimlico.core.dependencies package

Submodules

pimlico.core.dependencies.base module

Base classes for defining software dependencies for module types and routines for fetching them.

class SoftwareDependency (*name, url=None, dependencies=None*)

Bases: object

Base class for all Pimlico module software dependencies.

available (*local_config*)

Return True if the dependency is satisfied, meaning that the software/library is installed and ready to use.

problems (*local_config*)

Returns a list of problems standing in the way of the dependency being available. If the list is empty, the dependency is taken to be installed and ready to use.

Overriding methods should call super method.

installable ()

Return True if it's possible to install this library automatically. If False, the user will have to install it themselves. Instructions for doing this may be provided by `installation_instructions()`, which will only generally be called if `installable()` returns False.

This might be the case, for example, if the software is not available to download freely, or if it requires a system-wide installation.

installation_instructions ()

Where a dependency can't be installed programmatically, we typically want to be able to output instructions for the user to tell them how to go about doing it themselves. Any subclass that doesn't provide an automatic installation routine should override this to provide instructions.

You may also provide this even if the class does provide automatic installation. For example, you might want to provide instructions for other ways to install the software, like a system-wide install. This instructions will be shown together with missing dependency information.

dependencies ()

Returns a list of instances of :class:SoftwareDependency subclasses representing this library's own dependencies. If the library is already available, these will never be consulted, but if it is to be installed, we will check first that all of these are available (and try to install them if not).

install (*local_config*, *trust_downloaded_archives=False*)

Should be overridden by any subclasses whose library is automatically installable. Carries out the actual installation.

You may assume that all dependencies returned by :method:dependencies have been satisfied prior to calling this.

all_dependencies ()

Recursively fetch all dependencies of this dependency (not including itself).

get_installed_version (*local_config*)

If available() returns True, this method should return a SoftwareVersion object (or subclass) representing the software's version.

The base implementation returns an object representing an unknown version number.

If available() returns False, the behaviour is undefined and may raise an error.

class SystemCommandDependency (*name*, *test_command*, ***kwargs*)

Bases: *pimlico.core.dependencies.base.SoftwareDependency*

Dependency that tests whether a command is available on the command line. Generally requires system-wide installation.

installable ()

Usually not automatically installable

problems (*local_config*)

Returns a list of problems standing in the way of the dependency being available. If the list is empty, the dependency is taken to be installed and ready to use.

Overriding methods should call super method.

exception InstallationError

Bases: *exceptions.Exception*

check_and_install (*deps*, *local_config*, *trust_downloaded_archives=False*)

Check whether dependencies are available and try to install those that aren't. Returns a list of dependencies that can't be installed.

install (*dep*, *local_config*, *trust_downloaded_archives=False*)

install_dependencies (*pipeline*, *modules=None*, *trust_downloaded_archives=True*)

Install dependencies for pipeline modules

Parameters

- **pipeline** –
- **modules** – list of module names, or None to install for all

Returns

recursive_deps (*dep*)

Collect all recursive dependencies of this dependency. Does a depth-first search so that everything comes later in the list than things it depends on.

pimlico.core.dependencies.core module

Basic Pimlico core dependencies

CORE_PIMLICO_DEPENDENCIES = [`PythonPackageSystemwideInstall<Pip>`, `PythonPackageOnPip<virtua`

Core dependencies required by the basic Pimlico installation, regardless of what pipeline is being processed.

These will be checked when Pimlico is run, using the same dependency-checking mechanism that Pimlico modules use, and installed automatically if they're not found.

pimlico.core.dependencies.java module

class JavaDependency (*name*, *classes*=[], *jars*=[], *class_dirs*=[], ***kwargs*)

Bases: `pimlico.core.dependencies.base.SoftwareDependency`

Base class for Java library dependencies.

In addition to the usual functionality provided by dependencies, subclasses of this provide contributions to the Java classpath in the form of directories of jar files.

The instance has a set of representative Java classes that the checker will try to load to check whether the library is available and functional. It will also check that all jar files exist.

Jar paths and class directory paths are assumed to be relative to the Java lib dir (`lib/java`), unless they are absolute paths.

Subclasses should provide `install()` and override `installable()` if it's possible to install them automatically.

problems (*local_config*)

Returns a list of problems standing in the way of the dependency being available. If the list is empty, the dependency is taken to be installed and ready to use.

Overriding methods should call super method.

installable ()

Return True if it's possible to install this library automatically. If False, the user will have to install it themselves. Instructions for doing this may be provided by `installation_instructions()`, which will only generally be called if `installable()` returns False.

This might be the case, for example, if the software is not available to download freely, or if it requires a system-wide installation.

get_classpath_components ()

class JavaJarsDependency (*name*, *jar_urls*, ***kwargs*)

Bases: `pimlico.core.dependencies.java.JavaDependency`

Simple way to define a Java dependency where the library is packaged up in a jar, or a series of jars. The jars should be given as a list of (name, url) pairs, where name is the filename the jar should have and url is a url from which it can be downloaded.

URLs may also be given in the form "url->member", where url is a URL to a tar.gz or zip archive and member is a member to extract from the archive. If the type of the file isn't clear from the URL (i.e. if it doesn't have ".zip" or ".tar.gz" in it), specify the intended extension in the form "[ext]url->member", where ext is "tar.gz" or "zip".

installable ()

Return True if it's possible to install this library automatically. If False, the user will have to install it themselves. Instructions for doing this may be provided by `installation_instructions()`, which will only generally be called if `installable()` returns False.

This might be the case, for example, if the software is not available to download freely, or if it requires a system-wide installation.

install (*local_config*, *trust_downloaded_archives=False*)

Should be overridden by any subclasses whose library is automatically installable. Carries out the actual installation.

You may assume that all dependencies returned by `:method:dependencies` have been satisfied prior to calling this.

class PimlicoJavaLibrary (*name*, *classes=[]*, *additional_jars=[]*)

Bases: *pimlico.core.dependencies.java.JavaDependency*

Special type of Java dependency for the Java libraries provided with Pimlico. These are packages up in jars and stored in the build dir.

check_java_dependency (*class_name*, *classpath=None*)

Utility to check that a java class is able to be loaded.

check_java ()

Check that the JVM executable can be found. Raises a `DependencyError` if it can't be found or can't be run.

get_classpath (*deps*, *as_list=False*)

Given a list of `JavaDependency` subclass instances, returns all the components of the classpath that will make sure that the dependencies are available.

If *as_list=True*, returned as a list. Get the full classpath by `"".join(x)` on the list. If *as_list=False*, returns classpath string.

get_module_classpath (*module*)

Builds a classpath that includes all of the classpath elements specified by Java dependencies of the given module. These include the dependencies from `get_software_dependencies()` and also any dependencies of the datatype.

Used to ensure that Java modules that depend on particular jars or classes get all of those files included on their classpath when Java is run.

class Py4JSoftwareDependency

Bases: *pimlico.core.dependencies.java.JavaDependency*

Java component of Py4J. Use this one as the main dependency, as it depends on the Python component and will install that first if necessary.

dependencies ()

Returns a list of instances of `:class:SoftwareDependency` subclasses representing this library's own dependencies. If the library is already available, these will never be consulted, but if it is to be installed, we will check first that all of these are available (and try to install them if not).

jars

installable ()

Return True if it's possible to install this library automatically. If False, the user will have to install it themselves. Instructions for doing this may be provided by `installation_instructions()`, which will only generally be called if `installable()` returns False.

This might be the case, for example, if the software is not available to download freely, or if it requires a system-wide installation.

install (*local_config*, *trust_downloaded_archives=False*)

Should be overridden by any subclasses whose library is automatically installable. Carries out the actual installation.

You may assume that all dependencies returned by `:method:dependencies` have been satisfied prior to calling this.

pimlico.core.dependencies.python module

Tools for Python library dependencies.

Provides superclasses for Python library dependencies and a selection of commonly used dependency instances.

class PythonPackageDependency (*package, name, **kwargs*)

Bases: *pimlico.core.dependencies.base.SoftwareDependency*

Base class for Python dependencies. Provides import checks, but no installation routines. Subclasses should either provide `install()` or `installation_instructions()`.

The import checks do not (as of 0.6rc) actually import the package, as this may have side-effects that are difficult to account for, causing odd things to happen when you check multiple times, or try to import later. Instead, it just checks whether the package finder is about to locate the package. This doesn't guarantee that the import will succeed.

problems (*local_config*)

Returns a list of problems standing in the way of the dependency being available. If the list is empty, the dependency is taken to be installed and ready to use.

Overriding methods should call super method.

import_package ()

Try importing `package_name`. By default, just uses `__import__`. Allows subclasses to allow for special import behaviour.

Should raise an *ImportError* if import fails.

get_installed_version (*local_config*)

Tries to import a `__version__` variable from the package, which is a standard way to define the package version.

class PythonPackageSystemwideInstall (*package_name, name, pip_package=None, apt_package=None, yum_package=None, **kwargs*)

Bases: *pimlico.core.dependencies.python.PythonPackageDependency*

Dependency on a Python package that needs to be installed system-wide.

installable ()

Return True if it's possible to install this library automatically. If False, the user will have to install it themselves. Instructions for doing this may be provided by `installation_instructions()`, which will only generally be called if `installable()` returns False.

This might be the case, for example, if the software is not available to download freely, or if it requires a system-wide installation.

installation_instructions ()

Where a dependency can't be installed programmatically, we typically want to be able to output instructions for the user to tell them how to go about doing it themselves. Any subclass that doesn't provide an automatic installation routine should override this to provide instructions.

You may also provide this even if the class does provide automatic installation. For example, you might want to provide instructions for other ways to install the software, like a system-wide install. This instructions will be shown together with missing dependency information.

class PythonPackageOnPip (*package, name=None, pip_package=None, **kwargs*)

Bases: *pimlico.core.dependencies.python.PythonPackageDependency*

Python package that can be installed via pip. Will be installed in the virtualenv if not available.

installable ()

Return True if it's possible to install this library automatically. If False, the user will have to install it

themselves. Instructions for doing this may be provided by `installation_instructions()`, which will only generally be called if `installable()` returns `False`.

This might be the case, for example, if the software is not available to download freely, or if it requires a system-wide installation.

install (*local_config*, *trust_downloaded_archives=False*)

Should be overridden by any subclasses whose library is automatically installable. Carries out the actual installation.

You may assume that all dependencies returned by `:method:dependencies` have been satisfied prior to calling this.

get_installed_version (*local_config*)

Tries to import a `__version__` variable from the package, which is a standard way to define the package version.

safe_import_bs4 ()

BS can go very slowly if it tries to use `chardet` to detect input encoding. Remove `chardet` and `cchardet` from the Python modules, so that import fails and it doesn't try to use them. This prevents it getting stuck on reading long input files.

class BeautifulSoupDependency

Bases: `pimlico.core.dependencies.python.PythonPackageOnPip`

Test import with special BS import behaviour.

import_package ()

Try importing `package_name`. By default, just uses `__import__`. Allows subclasses to allow for special import behaviour.

Should raise an `ImportError` if import fails.

class NLTKResource (*name*, *url=None*, *dependencies=None*)

Bases: `pimlico.core.dependencies.base.SoftwareDependency`

Check for and install NLTK resources, using NLTK's own downloader.

problems (*local_config*)

Returns a list of problems standing in the way of the dependency being available. If the list is empty, the dependency is taken to be installed and ready to use.

Overriding methods should call super method.

installable ()

Return `True` if it's possible to install this library automatically. If `False`, the user will have to install it themselves. Instructions for doing this may be provided by `installation_instructions()`, which will only generally be called if `installable()` returns `False`.

This might be the case, for example, if the software is not available to download freely, or if it requires a system-wide installation.

install (*local_config*, *trust_downloaded_archives=False*)

Should be overridden by any subclasses whose library is automatically installable. Carries out the actual installation.

You may assume that all dependencies returned by `:method:dependencies` have been satisfied prior to calling this.

dependencies ()

Returns a list of instances of `:class:SoftwareDependency` subclasses representing this library's own dependencies. If the library is already available, these will never be consulted, but if it is to be installed, we will check first that all of these are available (and try to install them if not).

pimlico.core.dependencies.versions module

class SoftwareVersion (*string_id*)

Bases: object

Base class for representing version numbers / IDs of software. Different software may use different conventions to represent its versions, so it may be necessary to subclass this class to provide the appropriate parsing and comparison of versions.

compare_dotted_versions (*version0*, *version1*)

Comparison function for reasonably standard version numbers, with subversions to any level of nesting specified by dots.

Module contents

pimlico.core.external package

Submodules

pimlico.core.external.java module

call_java (*class_name*, *args*=[], *classpath*=None)

start_java_process (*class_name*, *args*=[], *java_args*=[], *wait*=0.1, *classpath*=None)

class Py4JInterface (*gateway_class*, *port*=None, *python_port*=None, *gateway_args*=[], *pipeline*=None, *print_stdout*=True, *print_stderr*=True, *env*={}, *system_properties*={}, *java_opts*=[], *timeout*=10.0, *prefix_classpath*=None)

Bases: object

start (*timeout*=None, *port_output_prefix*=None)

Start a Py4J gateway server in the background on the given port, which will then be used for communicating with the Java app.

If a port has been given, it is assumed that the gateway accepts a `-port` option. Likewise with `python_port` and a `-python-port` option.

If `timeout` is given, it overrides any `timeout` given in the constructor or specified in local config.

new_client ()

stop ()

clear_output_queues ()

no_retry_gateway (***kwargs*)

A wrapper around the constructor of `JavaGateway` that produces a version of it that doesn't retry on errors. The default gateway keeps retrying and outputting millions of errors if the server goes down, which makes responding to interrupts horrible (as the server might die before the Python process gets the interrupt).

TODO This isn't working: it just gets worse when I use my version!

gateway_client_to_running_server (*port*)

launch_gateway (*gateway_class*='py4j.GatewayServer', *args*=[], *javaopts*=[], *redirect_stdout*=None, *redirect_stderr*=None, *daemonize_redirect*=True, *env*={}, *port_output_prefix*=None, *startup_timeout*=10.0, *prefix_classpath*=None)

Our own more flexible version of Py4J's `launch_gateway`.

get_redirect_func (*redirect*)

class OutputConsumer (*redirects, stream, *args, **kwargs*)

Bases: `threading.Thread`

Thread that consumes output Modification of Py4J's OutputConsumer to allow multiple redirects.

remove_temporary_redirects ()

run ()

Method representing the thread's activity.

You may override this method in a subclass. The standard run() method invokes the callable object passed to the object's constructor as the target argument, if any, with sequential and keyword arguments taken from the args and kwargs arguments, respectively.

output_p4j_error_info (*command, returncode, stdout, stderr*)

make_py4j_errors_safe (*fn*)

Decorator for functions/methods that call Py4J. Py4J's exceptions include information that gets retrieved from the Py4J server when they're displayed. This is a problem if the server is not longer running and raises another exception, making the whole situation very confusing.

If you wrap your function with this, Py4JJavaErrors will be replaced by our own exception type Py4JSafeJavaError, containing some of the information about the Java exception if possible.

exception Py4JSafeJavaError (*java_exception=None, str=None*)

Bases: `exceptions.Exception`

exception DependencyCheckerError

Bases: `exceptions.Exception`

exception JavaProcessError

Bases: `exceptions.Exception`

Module contents

Tools for calling external (non-Python) tools.

pimlico.core.modules package

Subpackages

pimlico.core.modules.map package

Submodules

pimlico.core.modules.map.filter module

Todo: Continue updating this for the new datatype system. I've got partway, but the reader is still far from finished

class DocumentMapOutputTypeWrapper (**args, **kwargs*)

Bases: `object`

TODO: This was from the old datatypes system

Remove it once you've replicated the key bits in the reader class below.

```
non_filter_datatype = None
```

```
wrapped_module_info = None
```

```
output_name = None
```

```
archive_iter (subsample=None, start_after=None)
```

Provides an iterator just like TarredCorpus, but instead of iterating over data read from disk, gets it on the fly from the input datatype.

```
data_ready ()
```

Ready to supply this data as soon as all the wrapper module's inputs are ready to produce their data.

```
class FilterModuleOutputReader (datatype, setup, pipeline, **kwargs)
```

```
Bases: pimlico.datatypes.corpora.grouped.GroupedCorpusReader
```

A custom reader that is used for the output of a filter module, producing documents on the fly.

```
metadata = {}
```

```
extract_file (archive_name, filename)
```

Extract an individual file by archive name and filename. This is not an efficient way of extracting a lot of files. The typical use case of a grouped corpus is to iterate over its files, which is much faster.

```
list_archive_iter ()
```

```
archive_iter (start_after=None, skip=None, name_filter=None)
```

Iterate over corpus archive by archive, yielding for each document the archive name, the document name and the document itself.

Parameters

- **name_filter** – if given, should be a callable that takes two args, an archive name and document name, and returns True if the document should be yielded and False if it should be skipped. This can be preferable to filtering the yielded documents, as it skips all document pre-processing for skipped documents, so speeds up things like random subsampling of a corpus, where the document content never needs to be read in skipped cases
- **start_after** – skip over the first portion of the corpus, until the given document is reached. Should be specified as a pair (archive name, doc name)
- **skip** – skips over the first portion of the corpus, until this number of documents have been seen

Setup

alias of FilterModuleOutputReaderSetup

```
process_setup ()
```

Do any processing of the setup object (e.g. retrieving values and setting attributes on the reader) that should be done when the reader is instantiated.

```
wrap_module_info_as_filter (module_info_instance)
```

Create a filter module from a document map module so that it gets executed on the fly to provide its outputs as input to later modules. Can be applied to any document map module simply by adding *filter=T* to its config.

This function is called when *filter=T* is given.

Todo: Under the new datatype system, this should be done differently. Don't wrap datatypes, but instead use the actual output datatypes (taken from the wrapped module type's output) and instead create custom readers that gets instantiated when fetching the module's output readers.

I've created the test pipeline filter_tokenize for testing this.

Parameters `module_info_instance` – basic module info to wrap the outputs of

Returns a new non-executable ModuleInfo whose outputs are produced on the fly and will be identical to the outputs of the wrapper module.

pimlico.core.modules.map.multiproc module

Document map modules can in general be easily parallelized using multiprocessing. This module provides implementations of a pool and base worker processes that use multiprocessing, making it dead easy to implement a parallelized module, simply by defining what should be done on each document.

In particular, use `:fun:multiprocessing_executor_factory` wherever possible.

class MultiprocessingMapProcess (*input_queue, output_queue, exception_queue, executor, docs_per_batch=1*)

Bases: `multiprocessing.process.Process`, `pimlico.core.modules.map.DocumentMapProcessMixin`

A base implementation of document map parallelization using multiprocessing. Note that not all document map modules will want to use this: e.g. if you call a background service that provides parallelization itself (like the CoreNLP module) there's no need for multiprocessing in the Python code.

notify_no_more_inputs ()

Called when there aren't any more inputs to come.

run ()

Method to be run in sub-process; can be overridden in sub-class

class MultiprocessingMapPool (*executor, processes*)

Bases: `pimlico.core.modules.map.DocumentProcessorPool`

A base implementation of document map parallelization using multiprocessing.

PROCESS_TYPE = None

SINGLE_PROCESS_TYPE = None

start_worker ()

static create_queue (*maxsize=None*)

shutdown ()

notify_no_more_inputs ()

empty_all_queues ()

class MultiprocessingMapModuleExecutor (*module_instance_info, **kwargs*)

Bases: `pimlico.core.modules.map.DocumentMapModuleExecutor`

POOL_TYPE = None

create_pool (*processes*)

Should return an instance of the pool to be used for document processing. Should generally be a subclass of DocumentProcessorPool.

Always called after preprocess().

postprocess (*error=False*)

Allows subclasses to define a finishing procedure to be called after corpus processing if finished.

multiprocessing_executor_factory (*process_document_fn*, *preprocess_fn=None*, *post-process_fn=None*, *worker_set_up_fn=None*, *worker_tear_down_fn=None*, *batch_docs=None*, *multiprocessing_single_process=False*)

Factory function for creating an executor that uses the multiprocessing-based implementations of document-map pools and worker processes. This is an easy way to implement a parallelizable executor, which is suitable for a large number of module types.

process_document_fn should be a function that takes the following arguments (unless *batch_docs* is given):

- the worker process instance (allowing access to things set during setup)
- archive name
- document name
- the rest of the args are the document itself, from each of the input corpora

If *preprocess_fn* is given, it is called from the main process once before execution begins, with the executor as an argument.

If *postprocess_fn* is given, it is called from the main process at the end of execution, including on the way out after an error, with the executor as an argument and a kwarg *error* which is True if execution failed.

If *worker_set_up_fn* is given, it is called within each worker before execution begins, with the worker process instance as an argument. Likewise, *worker_tear_down_fn* is called from within the worker process before it exits.

Alternatively, you can supply a worker type, a subclass of `:class:MultiprocessingMapProcess`, as the first argument. If you do this, *worker_set_up_fn* and *worker_tear_down_fn* will be ignored.

If *batch_docs* is not None, *process_document_fn* is treated differently. Instead of supplying the *process_document()* of the worker, it supplies a *process_documents()*. The second argument is a list of tuples, each of which is assumed to be the args to *process_document()* for a single document. In this case, *docs_per_batch* is set on the worker processes, so that the given number of docs are collected from the input and passed into *process_documents()* at once.

By default, if only a single process is needed, we use the threaded implementation of a map process instead of multiprocessing. If this doesn't work out in your case, for some reason, specify *multiprocessing_single_process=True* and a multiprocessing process will be used even when only creating one.

pimlico.core.modules.map.singleproc module

Sometimes the simple multiprocessing-based approach to map module parallelization just isn't suitable. This module provides an equivalent set of implementations and convenience functions that don't use multiprocessing, but conform to the pool-based execution pattern by creating a single-thread pool.

class SingleThreadMapModuleExecutor (*module_instance_info*, ***kwargs*)

Bases: `pimlico.core.modules.map.threaded.ThreadingMapModuleExecutor`

create_pool (*processes*)

Should return an instance of the pool to be used for document processing. Should generally be a subclass of `DocumentProcessorPool`.

Always called after `preprocess()`.

single_process_executor_factory (*process_document_fn*, *preprocess_fn=None*, *post-process_fn=None*, *worker_set_up_fn=None*, *worker_tear_down_fn=None*, *batch_docs=None*)

Factory function for creating an executor that uses the single-process implementations of document-map pools and workers. This is an easy way to implement a non-parallelized executor

`process_document_fn` should be a function that takes the following arguments:

- the executor instance (allowing access to things set during setup)
- archive name
- document name
- the rest of the args are the document itself, from each of the input corpora

If `proprocess_fn` is given, it is called once before execution begins, with the executor as an argument.

If `postprocess_fn` is given, it is called at the end of execution, including on the way out after an error, with the executor as an argument and a kwarg `error` which is True if execution failed.

pimlico.core.modules.map.threaded module

Just like multiprocessing, but using threading instead. If you're not sure which you should use, it's probably multiprocessing.

class ThreadingMapThread (*input_queue, output_queue, exception_queue, executor*)

Bases: `threading.Thread`, `pimlico.core.modules.map.DocumentMapProcessMixin`

notify_no_more_inputs ()

Called when there aren't any more inputs to come.

run ()

Method representing the thread's activity.

You may override this method in a subclass. The standard `run()` method invokes the callable object passed to the object's constructor as the target argument, if any, with sequential and keyword arguments taken from the args and kwargs arguments, respectively.

shutdown (*timeout=3.0*)

class ThreadingMapPool (*executor, processes*)

Bases: `pimlico.core.modules.map.DocumentProcessorPool`

THREAD_TYPE = None

start_worker ()

static create_queue (*maxsize=None*)

shutdown ()

class ThreadingMapModuleExecutor (*module_instance_info, **kwargs*)

Bases: `pimlico.core.modules.map.DocumentMapModuleExecutor`

POOL_TYPE = None

create_pool (*processes*)

Should return an instance of the pool to be used for document processing. Should generally be a subclass of `DocumentProcessorPool`.

Always called after `preprocess()`.

postprocess (*error=False*)

Allows subclasses to define a finishing procedure to be called after corpus processing if finished.

threading_executor_factory (*process_document_fn, preprocess_fn=None, postprocess_fn=None, worker_set_up_fn=None, worker_tear_down_fn=None*)

Factory function for creating an executor that uses the threading-based implementations of document-map pools and worker processes.

`process_document_fn` should be a function that takes the following arguments:

- the worker process instance (allowing access to things set during setup)
- archive name
- document name
- the rest of the args are the document itself, from each of the input corpora

If `proprocess_fn` is given, it is called from the main thread once before execution begins, with the executor as an argument.

If `postprocess_fn` is given, it is called from the main thread at the end of execution, including on the way out after an error, with the executor as an argument and a kwarg `error` which is True if execution failed.

If `worker_set_up_fn` is given, it is called within each worker before execution begins, with the worker thread instance as an argument. Likewise, `worker_tear_down_fn` is called from within the worker thread before it exits.

Alternatively, you can supply a worker type, a subclass of `:class:ThreadingMapThread`, as the first argument. If you do this, `worker_set_up_fn` and `worker_tear_down_fn` will be ignored.

Module contents

class DocumentMapModuleInfo (*module_name, pipeline, **kwargs*)

Bases: `pimlico.core.modules.base.BaseModuleInfo`

Abstract module type that maps each document in turn in a corpus. It produces a single output document for every input.

Subclasses should specify the input types, which should all be subclasses of `TarredCorpus`, and output types, the first of which (i.e. default) should also be a subclass of `TarredCorpus`. The base class deals with iterating over the input(s) and writing the outputs to a new `TarredCorpus`. The subclass only needs to implement the mapping function applied to each document (in its executor).

`module_outputs = [('documents', grouped_corpus)]`

`input_corpora`

`get_writers (append=False)`

`get_detailed_status ()`

Returns a list of strings, containing detailed information about the module's status that is specific to the module type. This may include module-specific information about execution status, for example.

Subclasses may override this to supply useful (human-readable) information specific to the module type. They should called the super method.

`document (output_name=None, **kwargs)`

Instantiate a document of the output type for the given output name (or number), or the default output.

Convenience utility to avoid having to look up the output data point type to do this.

class DocumentMapModuleExecutor (*module_instance_info, **kwargs*)

Bases: `pimlico.core.modules.base.BaseModuleExecutor`

Base class for executors for document map modules. Subclasses should provide the behaviour for each individual document by defining a pool (and worker processes) to handle the documents as they're fed into it.

Note that in most cases it won't be necessary to override the pool and worker base classes yourself. Unless you need special behaviour, use the standard implementations and factory functions.

Although the pattern of execution for all document map modules is based on parallel processing (creating a pool, spawning worker processes, etc), this doesn't mean that all such modules have to be parallelizable. If you have no reason not to parallelize, it's recommended that you do (with single-process execution as a special case). However, sometimes parallelizing isn't so simple: in these cases, consider using the tools in `:mod:.singleproc`.

preprocess ()

Allows subclasses to define a set-up procedure to be called before corpus processing begins.

postprocess (*error=False*)

Allows subclasses to define a finishing procedure to be called after corpus processing if finished.

create_pool (*processes*)

Should return an instance of the pool to be used for document processing. Should generally be a subclass of `DocumentProcessorPool`.

Always called after `preprocess()`.

retrieve_processing_status ()

update_processing_status (*docs_completed, archive_name, filename*)

execute ()

Run the actual module execution.

May return `None`, in which case it's assumed to have fully completed. If a string is returned, it's used as an alternative module execution status. Used, e.g., by multi-stage modules that need to be run multiple times.

skip_invalid (*fn*)

Decorator to apply to document map executor `process_document()` methods where you want to skip doing any processing if any of the input documents are invalid and just pass through the error information.

Be careful not to confuse this with the `process_document()` methods on datatypes. You don't need a decorator on them to skip invalid documents, as it's not called on them anyway.

skip_invalids (*fn*)

Decorator to apply to document map executor `process_documents()` methods where you want to skip doing any processing if any of the input documents are invalid and just pass through the error information.

invalid_doc_on_error (*fn*)

Decorator to apply to `process_document()` methods that causes all exceptions to be caught and an `InvalidDocument` to be returned as the result, instead of letting the error propagate up and call a halt to the whole corpus processing.

invalid_docs_on_error (*fn*)

Decorator to apply to `process_documents()` methods that causes all exceptions to be caught and an `InvalidDocument` to be returned as the result for every input document.

class ProcessOutput (*archive, filename, data*)

Bases: `object`

Wrapper for all result data coming out from a worker.

class InputQueueFeeder (*input_queue, iterator, complete_callback=None*)

Bases: `threading.Thread`

Background thread to read input documents from an iterator and feed them onto an input queue for worker processes/threads.

get_next_output_document ()

check_invalid (*archive, filename*)

Checks whether a given document was invalid in the input. Once the check has been performed, the item is removed from the list, for efficiency, so this should only be called once per document.

run ()

Method representing the thread's activity.

You may override this method in a subclass. The standard run() method invokes the callable object passed to the object's constructor as the target argument, if any, with sequential and keyword arguments taken from the args and kwargs arguments, respectively.

check_for_error ()

Can be called from the main thread to check whether an error has occurred in this thread and raise a suitable exception if so

shutdown (timeout=3.0)

Cancel the feeder, if it's still feeding and stop the thread. Call only after you're sure you no longer need anything from any of the queues. Waits for the thread to end.

Call from the main thread (that created the feeder) only.

class DocumentProcessorPool (processes)

Bases: object

Base class for pools that provide an easy implementation of parallelization for document map modules. Defines the core interface for pools.

If you're using multiprocessing, you'll want to use the multiprocessing-specific subclass.

notify_no_more_inputs ()

static create_queue (maxsize=None)

May be overridden by subclasses to provide different implementations of a Queue. By default, uses the multiprocessing queue type. Whatever is returned, it should implement the interface of Queue.Queue.

shutdown ()

empty_all_queues ()

class DocumentMapProcessMixin (input_queue, output_queue, exception_queue, docs_per_batch=1)

Bases: object

Mixin/base class that should be implemented by all worker processes for document map pools.

set_up ()

Called when the process starts, before it starts accepting documents.

process_document (archive, filename, *docs)

process_documents (doc_tuples)

Batched version of process_document(). Default implementation just calls process_document() on each document, but if you want to group documents together and process multiple at once, you can override this method and make sure the docs_per_batch is set > 1.

Each item in the list of doc tuples should be a tuple of the positional args to process_document() – i.e. archive_name, filename, doc_from_corpus1, [doc_from_corpus2, ...]

tear_down ()

Called from within the process after processing is complete, before exiting.

notify_no_more_inputs ()

Called when there aren't any more inputs to come.

exception WorkerStartupError (*args, **kwargs)

Bases: exceptions.Exception

```
exception WorkerShutdownError (*args, **kwargs)
    Bases: exceptions.Exception
```

Submodules

pimlico.core.modules.base module

This module provides base classes for Pimlico modules.

The procedure for creating a new module is the same whether you're contributing a module to the core set in the Pimlico codebase or a standalone module in your own codebase, or for a specific pipeline.

A Pimlico module is identified by the full Python-path to the Python package that contains it. This package should be laid out as follows:

- The module's metadata is defined by a class in `info.py` called `ModuleInfo`, which should inherit from `BaseModuleInfo` or one of its subclasses.
- The module's functionality is provided by a class in `execute.py` called `ModuleExecutor`, which should inherit from `BaseModuleExecutor`.

The exec Python module will not be imported until an instance of the module is to be run. This means that you can import dependencies and do any necessary initialization at the point where it's executed, without worrying about incurring the associated costs (and dependencies) every time a pipeline using the module is loaded.

```
class BaseModuleInfo (module_name, pipeline, inputs={}, options={}, optional_outputs=[],
                      docstring="", include_outputs=[], alt_expanded_from=None,
                      alt_param_settings=[], module_variables={})
```

Bases: `object`

Abstract base class for all pipeline modules' metadata.

module_type_name = None

module_readable_name = None

module_options = {}

Specifies a list of (name, datatype class) pairs for inputs that are always required

module_inputs = []

Specifies a list of (name, datatype class) pairs for optional inputs. The module's execution may vary depending on what is provided. If these are not given, `None` is returned from `get_input()`

module_optional_inputs = []

Specifies a list of (name, datatype class) pairs for outputs that are always written

module_optional_outputs = []

Whether the module should be executed Typically `True` for almost all modules, except input modules (though some of them may also require execution) and filters

module_executable = True

If specified, this `ModuleExecutor` class will be used instead of looking one up in the exec Python module

module_executor_override = None

Usually `None`. In the case of stages of a multi-stage module, stores a pointer to the main module.

main_module = None

module_outputs = []

Specifies a list of (name, datatype class) pairs for outputs that are written only if they're specified in the "output" option or used by another module

load_executor()

Loads a `ModuleExecutor` for this Pimlico module. Usually, this just involves calling `load_module_executor()`, but the default executor loading may be overridden for a particular module type by overriding this function. It should always return a subclass of `ModuleExecutor`, unless there's an error.

classmethod get_key_info_table()

When generating module docs, the table at the top of the page is produced by calling this method. It should return a list of two-item lists (title + value). Make sure to include the super-class call if you override this to add in extra module-specific info.

metadata_filename**get_metadata()****set_metadata_value** (*attr, val*)**set_metadata_values** (*val_dict*)**status****execution_history_path****add_execution_history_record** (*line*)

Output a single line to the file that stores the history of module execution, so we can trace what we've done.

execution_history

Get the entire recorded execution history for this module. Returns an empty string if no history has been recorded.

input_names

All required inputs, first, then all supplied optional inputs

output_names**classmethod process_module_options** (*opt_dict*)

Parse the options in a dictionary (probably from a config file), checking that they're valid for this model type.

Parameters *opt_dict* – dict of options, keyed by option name

Returns dict of options

classmethod extract_input_options (*opt_dict, module_name=None, previous_module_name=None, module_expansions={}*)

Given the config options for a module instance, pull out the ones that specify where the inputs come from and match them up with the appropriate input names.

The inputs returned are just names as they come from the config file. They are split into module name and output name, but they are not in any way matched up with the modules they connect to or type checked.

Parameters

- **module_name** – name of the module being processed, for error output. If not given, the name isn't included in the error.
- **previous_module_name** – name of the previous module in the order given in the config file, allowing a single-input module to default to connecting to this if the input connection wasn't given
- **module_expansions** – dictionary mapping module names to a list of expanded module names, where expansion has been performed as a result of alternatives in the param-

ters. Provided here so that the unexpanded names may be used to refer to the whole list of module names, where a module takes multiple inputs on one input parameter

Returns dictionary of inputs

static get_extra_outputs_from_options (*options*)

Normally, which optional outputs get produced by a module depend on the ‘output’ option given in the config file, plus any outputs that get used by subsequent modules. By overriding this method, module types can add extra outputs into the list of those to be included, conditional on other options.

It also receives the processed dictionary of inputs, so that the additional outputs can depend on what is fed into the input.

E.g. the `corenlp` module include the ‘annotations’ output if annotators are specified, so that the user doesn’t need to give both options.

provide_further_outputs ()

Called during instantiation, once inputs and options are available, to add a further list of module outputs that are dependent on inputs or options.

get_module_output_dir (*absolute=False, short_term_store=None*)

Gets the path to the base output dir to be used by this module, relative to the storage base dir. When outputting data, the storage base dir will always be the short term store path, but when looking for the output data other base paths might be explored, including the long term store.

Kwarg `short_term_store` is included for backward compatibility, but outputs a deprecation warning.

Parameters **absolute** – if True, return absolute path to output dir in output store

Returns path, relative to store base path, or if `absolute=True` absolute path to output dir

get_absolute_output_dir (*output_name*)

The simplest way to get hold of the directory to use to output data to for a given output. This is the usual way to get an output directory for an output writer.

The directory is an absolute path to a location in the Pimlico output storage location.

Parameters **output_name** – the name of an output

Returns the absolute path to the output directory to use for the named output

get_output_dir (*output_name, absolute=False, short_term_store=None*)

Kwarg `short_term_store` is included for backward compatibility, but outputs a deprecation warning.

Parameters

- **absolute** – return an absolute path in the storage location used for output. If False (default), return a relative path, specified relative to the root of the Pimlico store used. This allows multiple stores to be searched for output
- **output_name** – the name of an output

Returns the path to the output directory to use for the named output, which may be relative to the root of the Pimlico store in use (default) or an absolute path in the output store, depending on *absolute*

get_output_datatype (*output_name=None*)

Get the datatype of a named output, or the default output. Returns an instance of the relevant Pimlico-Datatype subclass. This can be used for typechecking and also for getting a reader for the output data, once it’s ready, by supplying it with the path to the data.

To get a reader for the output data, use `get_output()`.

Parameters **output_name** – output whose datatype to retrieve. Default output if not specified

Returns**instantiate_output_datatype** (*output_name, output_datatype*)

Subclasses may want to override this to provide special behaviour for instantiating particular outputs' datatypes.

Deprecated since version new: datatypes

Roughly replaced by `instantiate_output_reader()`, but many of the use cases can be covered in other ways now

output_ready (*output_name=None*)

Check whether the named output is ready to be read from one of its possible storage locations.

Parameters **output_name** – output to check, or default output if not given

Returns False if data is not ready to be read

instantiate_output_reader_setup (*output_name, datatype*)

Produce a reader setup instance that will be used to prepare this reader. This provides functionality like checking that the data is ready to be read before the reader is instantiated.

The standard implementation uses the datatype's methods to get its standard reader setup and reader, but some modules may need to override this to provide other readers.

output_name is provided so that overriding methods' behaviour can be conditioned on which output is being fetched.

instantiate_output_reader (*output_name, datatype, pipeline, module=None*)

Prepare a reader for a particular output. The default implementation is very simple, but subclasses may override this for cases where the normal process of creating readers has to be modified.

Parameters

- **output_name** – output to produce a reader for
- **datatype** – the datatype for this output, already inferred

get_output_reader_setup (*output_name=None*)**get_output** (*output_name=None*)

Get a reader corresponding to one of the outputs of the module. The reader will be that which corresponds to the output's declared datatype and will read the data from any of the possible locations where it can be found.

If the data is not available in any location, raises a `DataNotReadyError`.

To check whether the data is ready without calling this, call `output_ready()`.

get_output_writer (*output_name=None, **kwargs*)

Get a writer instance for the given output. Kwargs will be passed through to the writer and used to specify metadata and writer params.

Parameters

- **output_name** – output to get writer for, or default output if left
- **kwargs** –

Returns**is_multiple_input** (*input_name=None*)

Returns True if the named input (or default input if no name is given) is a `MultipleInputs` input, False otherwise. If it is, `get_input()` will return a list, otherwise it will return a single datatype.

get_input_module_connection (*input_name=None, always_list=False*)

Get the ModuleInfo instance and output name for the output that connects up with a named input (or the first input) on this module instance. Used by `get_input()` – most of the time you probably want to use that to get the instantiated datatype for an input.

If the input type was specified with `MultipleInputs`, meaning that we’re expecting an unbounded number of inputs, this is a list. Otherwise, it’s a single (module, output_name) pair. If `always_list=True`, in this latter case we return a single-item list.

get_input_datatype (*input_name=None, always_list=False*)

Get a list of datatype instances corresponding to one of the inputs to the module. If an input name is not given, the first input is returned.

If the input type was specified with `MultipleInputs`, meaning that we’re expecting an unbounded number of inputs, this is a list. Otherwise, it’s a single datatype.

get_input_reader_setup (*input_name=None, always_list=False*)

Get reader setup for one of the inputs to the module. Looks up the corresponding output from another module and uses that module’s metadata to get that output’s instance. If an input name is not given, the first input is returned.

If the input type was specified with `MultipleInputs`, meaning that we’re expecting an unbounded number of inputs, this is a list. Otherwise, it’s a single datatype instance. If `always_list=True`, in this latter case we return a single-item list.

If the requested input name is an optional input and it has not been supplied, returns `None`.

You can get a reader for the input, once the data is ready to be read, by calling `get_reader()` on the setup object. Or use `get_input()` on the module.

get_input (*input_name=None, always_list=False*)

Get a reader for one of the inputs to the module. Should only be called once the input data is ready to read. It’s therefore fine to call this from a module executor, since data availability has already been checked by this point.

If the input type was specified with `MultipleInputs`, meaning that we’re expecting an unbounded number of inputs, this is a list. Otherwise, it’s a single datatype instance. If `always_list=True`, in this latter case we return a single-item list.

If the requested input name is an optional input and it has not been supplied, returns `None`.

input_ready (*input_name=None*)

Check whether the data is ready to go corresponding to the named input.

Parameters `input_name` – input to check

Returns True if input is ready

all_inputs_ready ()

Check `input_ready()` on all inputs.

Returns True if all input datatypes are ready to be used

classmethod is_filter ()

missing_module_data ()

Reports missing data not associated with an input dataset.

Calling `missing_data()` reports any problems with input data associated with a particular input to this module. However, modules may also rely on data that does not come from one of their inputs. This happens primarily (perhaps solely) when a module option points to a data source. This might be the case with any module, but is particularly common among input reader modules, which have no inputs, but read data according to their options.

Returns list of problems

missing_data (*input_names=None, assume_executed=[], assume_failed=[], allow_preliminary=False*)

Check whether all the input data for this module is available. If not, return a list strings indicating which outputs of which modules are not available. If it's all ready, returns an empty list.

To check specific inputs, give a list of input names. To check all inputs, don't specify *input_names*. To check the default input, give *input_names=[None]*. If not checking a specific input, also checks non-input data (see *missing_module_data()*).

If *assume_executed* is given, it should be a list of module names which may be assumed to have been executed at the point when this module is executed. Any outputs from those modules will be excluded from the input checks for this module, on the assumption that they will have become available, even if they're not currently available, by the time they're needed.

If *assume_failed* is given, it should be a list of module names which should be assumed to have failed. If we rely on data from the output of one of them, instead of checking whether it's available we simply assume it's not.

Why do this? When running multiple modules in sequence, if one fails it is possible that its output datasets look like complete datasets. For example, a partially written iterable corpus may look like a perfectly valid corpus, which happens to be smaller than it should be. After the execution failure, we may check other modules to see whether it's possible to run them. Then we need to know not to trust the output data from the failed module, even if it looks valid.

If *allow_preliminary=True*, for any inputs that are multiple inputs and have multiple connections to previous modules, consider them to be satisfied if at least one of their inputs is ready. The normal behaviour is to require all of them to be ready, but in a preliminary run this requirement is relaxed.

classmethod is_input ()

dependencies

Returns list of names of modules that this one depends on for its inputs.

get_transitive_dependencies ()

Transitive closure of *dependencies*.

Returns list of names of modules that this one recursively (transitively) depends on for its inputs.

typecheck_inputs ()

typecheck_input (input_name)

Typecheck a single input. *typecheck_inputs ()* calls this and is used for typechecking of a pipeline. This method returns the (or the first) satisfied input requirement, or raises an exception if typechecking failed, so can be handy separately to establish which requirement was met.

The result is always a list, but will contain only one item unless the input is a multiple input.

get_software_dependencies ()

Check that all software required to execute this module is installed and locatable. This is separate to metadata config checks, so that you don't need to satisfy the dependencies for all modules in order to be able to run one of them. You might, for example, want to run different modules on different machines. This is called when a module is about to be executed and each of the dependencies is checked.

Returns a list of instances of subclasses of `:class:~pimlico.core.dependencies.base.SoftwareDependency`, representing the libraries that this module depends on.

Take care when providing dependency classes that you don't put any import statements at the top of the Python module that will make loading the dependency type itself dependent on runtime dependencies. You'll want to run import checks by putting import statements within this method.

You should call the super method for checking superclass dependencies.

get_input_software_dependencies ()

Collects library dependencies from the input datatypes to this module, which will need to be satisfied for the module to be run.

Unlike *get_software_dependencies ()*, it shouldn't need to be overridden by subclasses, since it just collects the results of getting dependencies from the datatypes.

get_output_software_dependencies ()

Collects library dependencies from the output datatypes to this module, which will need to be satisfied for the module to be run.

Unlike *get_input_software_dependencies ()*, it may not be the case that all of these dependencies strictly need to be satisfied before the module can be run. It could be that a datatype can be written without satisfying all the dependencies needed to read it. However, we assume that dependencies of all output datatypes must be satisfied in order to run the module that writes them, since this is usually the case, and these are checked before running the module.

Unlike *get_software_dependencies ()*, it shouldn't need to be overridden by subclasses, since it just collects the results of getting dependencies from the datatypes.

check_ready_to_run ()

Called before a module is run, or if the 'check' command is called. This will only be called after all library dependencies have been confirmed ready (see *:method:get_software_dependencies*).

Essentially, this covers any module-specific checks that used to be in *check_runtime_dependencies()* other than library installation (e.g. checking models exist).

Always call the super class' method if you override.

Returns a list of (name, description) pairs, where the name identifies the problem briefly and the description explains what's missing and (ideally) how to fix it.

reset_execution ()

Remove all output data and metadata from this module to make a fresh start, as if it's never been executed.

May be overridden if a module has some side effect other than creating/modifying things in its output directory(/ies), but overridden methods should always call the super method. Occasionally this is necessary, but most of the time the base implementation is enough.

get_detailed_status ()

Returns a list of strings, containing detailed information about the module's status that is specific to the module type. This may include module-specific information about execution status, for example.

Subclasses may override this to supply useful (human-readable) information specific to the module type. They should called the super method.

classmethod module_package_name ()

The package name for the module, which is used to identify it in config files. This is the package containing the *info.py* in which the *ModuleInfo* is defined.

get_execution_dependency_tree ()

Tree of modules that will be executed when this one is executed. Where this module depends on filters, the tree goes back through them to find what they depend on (since they will be executed simultaneously)

get_all_executed_modules ()

Returns a list of all the modules that will be executed when this one is (including itself). This is the current module (if executable), plus any filters used to produce its inputs.

lock_path

lock()

Mark the module as locked, so that it cannot be executed. Called when execution begins, to ensure that you don't end up executing the same module twice simultaneously.

unlock()

Remove the execution lock on this module.

is_locked()

Returns True if the module is currently locked from execution

get_new_log_filename (*name='error'*)

Returns an absolute path that can be used to output a log file for this module. This is used for outputting error logs. It will always return a filename that doesn't currently exist, so can be used multiple times to output multiple logs.

collect_unexecuted_dependencies (*modules*)

Given a list of modules, checks through all the modules that they depend on to put together a list of modules that need to be executed so that the given list will be left in an executed state. The list includes the modules themselves, if they're not fully executed, and unexecuted dependencies of any unexecuted modules (recursively).

Parameters **modules** – list of ModuleInfo instances

Returns list of ModuleInfo instances that need to be executed

collect_runnable_modules (*pipeline, preliminary=False*)

Look for all unexecuted modules in the pipeline to find any that are ready to be executed. Keep collecting runnable modules, including those that will become runnable once we've run earlier ones in the list, to produce a list of a sequence of modules that could be set running now.

Parameters **pipeline** – pipeline config

Returns ordered list of runnable modules. Note that it must be run in this order, as some might depend on earlier ones in the list

satisfies_typecheck (*provided_type, type_requirements*)

Interface to Pimlico's standard type checking (see *check_type*) that returns a boolean to say whether type checking succeeded or not.

check_type (*provided_type, type_requirements*)

Type-checking algorithm for making sure outputs from modules connect up with inputs that they satisfy the requirements for.

class BaseModuleExecutor (*module_instance_info, stage=None, debug=False, force_rerun=False*)

Bases: object

Abstract base class for executors for Pimlico modules. These are classes that actually do the work of executing the module on given inputs, writing to given output locations.

execute()

Run the actual module execution.

May return None, in which case it's assumed to have fully completed. If a string is returned, it's used as an alternative module execution status. Used, e.g., by multi-stage modules that need to be run multiple times.

exception ModuleInfoLoadError (**args, **kwargs*)

Bases: exceptions.Exception

exception ModuleExecutorLoadError

Bases: exceptions.Exception

exception ModuleTypeError

Bases: exceptions.Exception

exception TypeCheckError

Bases: `exceptions.Exception`

exception DependencyError (*message, stderr=None, stdout=None*)

Bases: `exceptions.Exception`

Raised when a module's dependencies are not satisfied. Generally, this means a dependency library needs to be installed, either on the local system or (more often) by calling the appropriate make target in the lib directory.

load_module_executor (*path_or_info*)

Utility for loading the executor class for a module from its full path. More or less just a wrapper around an import, with some error checking. Locates the executor by a standard procedure that involves checking for an "execute" python module alongside the info's module.

Note that you shouldn't generally use this directly, but instead call the `load_executor()` method on a module info (which will call this, unless special behaviour has been defined).

Parameters `path` – path to Python package containing the module

Returns class

load_module_info (*path*)

Utility to load the metadata for a Pimlico pipeline module from its package Python path.

Parameters `path` –

Returns

pimlico.core.modules.execute module

Runtime execution of modules

This module provides the functionality to check that Pimlico modules are ready to execute and execute them. It is used by the `run` command.

check_and_execute_modules (*pipeline, module_names, force_rerun=False, debug=False, log=None, all_deps=False, check_only=False, exit_on_error=False, preliminary=False, email=None*)

Main method called by the `run` command that first checks a pipeline, checks all pre-execution requirements of the modules to be executed and then executes each of them. The most common case is to execute just one module, but a sequence may be given.

Parameters

- **exit_on_error** – drop out if a `ModuleExecutionError` occurs in any individual module, instead of continuing to the next module that can be run
- **pipeline** – loaded `PipelineConfig`
- **module_names** – list of names of modules to execute in the order they should be run
- **force_rerun** – execute modules, even if they're already marked as complete
- **debug** – output debugging info
- **log** – logger, if you have one you want to reuse
- **all_deps** – also include unexecuted dependencies of the given modules
- **check_only** – run all checks, but stop before executing. Used for `check` command

Returns

check_modules_ready (*pipeline, modules, log, preliminary=False*)

Check that a module is ready to be executed. Always called before execution begins.

Parameters

- **pipeline** – loaded PipelineConfig
- **modules** – loaded ModuleInfo instances, given in the order they’re going to be executed. For each module, it’s assumed that those before it in the list have already been run when it is run.
- **log** – logger to output to

Returns If *preliminary=True*, list of problems that were ignored by allowing preliminary run. Otherwise, None – we raise an exception when we first encounter a problem

execute_modules (*pipeline, modules, log, force_rerun=False, debug=False, exit_on_error=False, preliminary=False, email=None*)

format_execution_dependency_tree (*tree*)

send_final_report_email (*pipeline, error_modules, success_modules, skipped_modules, all_modules*)

send_module_report_email (*pipeline, module, short_error, long_error*)

exception ModuleExecutionError (**args, **kwargs*)

Bases: `exceptions.Exception`

exception ModuleNotReadyError (**args, **kwargs*)

Bases: `pimlico.core.modules.execute.ModuleExecutionError`

exception ModuleAlreadyCompletedError (**args, **kwargs*)

Bases: `pimlico.core.modules.execute.ModuleExecutionError`

exception StopProcessing

Bases: `exceptions.Exception`

pimlico.core.modules.inputs module

Base classes and utilities for input modules in a pipeline.

```
class InputModuleInfo (module_name, pipeline, inputs={}, options={}, optional_outputs=[],
                        docstring="", include_outputs=[], alt_expanded_from=None,
                        alt_param_settings=[], module_variables={})
```

Bases: `pimlico.core.modules.base.BaseModuleInfo`

Base class for input modules. These don’t get executed in general, they just provide a way to iterate over input data.

You probably don’t want to subclass this. It’s usually simplest to define a datatype for reading the input data and then just specify its class as the module’s type. This results in a subclass of this module info being created dynamically to read that data.

Note that `module_executable` is typically set to `False` and the base class does this. However, some input modules need to be executed before the input is usable, for example to collect stats about the input data.

```
module_type_name = 'input'
```

```
module_executable = False
```

input_module_factory (*datatype*)

Create an input module class to load a given datatype.

This is used by the pipeline config loader to create a suitable module type when the config has a datatype as a module type. It loads data from the given directory exactly as if Pimlico had itself output a dataset of the specified type to that directory.

The main use for this is loading prepared datasets in test pipelines. It is also useful if you have output from some other Pimlico pipeline that you just want to load as it is and use as input to a module.

It is not for loading input data from external sources. See *iterable_input_reader* for creating normal input modules.

iterable_input_reader (*input_module_options*, *data_point_type*, *data_ready_fn*, *len_fn*, *iter_fn*, *module_type_name=None*, *module_readable_name=None*, *software_dependencies=None*, *execute_count=False*, *no_group=False*)

Factory for creating an input reader module info. This is a (typically) non-executable module that has no inputs. It reads its data from some external location, using the given module options. The resulting dataset is a GroupedCorpus, with the given document type.

This is the normal way to create input reader modules.

The returned class is a subclass of *BaseModuleInfo*. It is typically used like this, within a Pimlico module's *info.py*:

```
ModuleInfo = iterable_input_reader(  
    {  
        # ... module options ...  
    },  
    DataPointType(),  
    data_ready_function,  
    len_function,  
    iter_function,  
    "my_module_name"  
)
```

If *execute_count=True*, the module will be an executable module and the execution will simply count the number of documents in the corpus and store the count. This should be used if counting the documents in the dataset is not completely trivial and quick (e.g. if you need to read through the data itself, rather than something like counting files in a directory or checking metadata). It is common for this to be the only processing that needs to be done on the dataset before using it. The *len_fn* is used to count the documents in the module's execution phase.

If the counting method returns a pair of integers, instead of just a single integer, they are taken to be the total number of documents in the corpus and the number of valid documents (i.e. the number that will be produce an InvalidDocument). In this case, the valid documents count is also stored in the metadata, as *valid_documents*.

Reader options are available at read time from the reader setup instance's *reader_options* attribute, also available from the reader instance as *reader.options*.

Note: Producing an IterableCorpus used to be the default behaviour. However, since we almost always want to convert to a GroupedCorpus immediately after reading, the default behaviour is now to do the grouping as part of the reading process and produce a GroupedCorpus straight away. If you want to regroup for some reason, you can, of course, still do that with the resulting GroupedCorpus.

If you need a plain IterableCorpus as output, you can use *no_group=True* when calling this factory, which will produce the old behaviour.

Parameters

- **input_module_options** – dictionary defining the module options for the input module, which will be provided to all the functions
- **data_point_type** – a data point type for the individual documents that will be produced. They do not need to be read in using this type’s reading functionality, which will later be used for storing and reading the documents, but can be produced by some other means.
- **data_ready_fn** – function that takes the processed options given to the module in the config file and returns True if the data is ready to read, False otherwise. If `execute_count` is used, the data will be considered unread until the count has been run, even if this function returns True.
- **iter_fn** – function that takes a reader instance and returns a generator to iterate over the documents of the corpus. Like any `IterableCorpus`, it should yield pairs of (`doc_name`, `doc`). Reader options are available as `reader.setup.reader_options`.
- **len_fn** – function that takes the processed options given to the module in the config file and returns the number of docs
- **module_type_name** –
- **module_readable_name** –
- **software_dependencies** – a list of software dependencies that the module-info will return when `get_software_dependencies()` is called, or a function that takes the module-info instance and returns such a list. If left blank, no dependencies are returned.
- **execute_count** – make an executable module that counts the data to get its length (num docs)
- **no_group** – by default, the output datatype is a `GroupedCorpus`. If True, use an `IterableCorpus` instead without grouping documents into archives.

Returns module info class

pimlico.core.modules.multistage module

class MultistageModuleInfo (*module_name*, *pipeline*, ***kwargs*)

Bases: `pimlico.core.modules.base.BaseModuleInfo`

Base class for multi-stage modules. You almost certainly don’t want to override this yourself, but use the factory method instead. It exists mainly for providing a way of identifying multi-stage modules.

module_executable = True

stages = None

typecheck_inputs ()

Overridden to check internal output-input connections as well as the main module’s inputs.

get_software_dependencies ()

Check that all software required to execute this module is installed and locatable. This is separate to metadata config checks, so that you don’t need to satisfy the dependencies for all modules in order to be able to run one of them. You might, for example, want to run different modules on different machines. This is called when a module is about to be executed and each of the dependencies is checked.

Returns a list of instances of subclasses of `:class:~pimlico.core.dependencies.base.SoftwareDependency`, representing the libraries that this module depends on.

Take care when providing dependency classes that you don't put any import statements at the top of the Python module that will make loading the dependency type itself dependent on runtime dependencies. You'll want to run import checks by putting import statements within this method.

You should call the super method for checking superclass dependencies.

get_input_software_dependencies()

Collects library dependencies from the input datatypes to this module, which will need to be satisfied for the module to be run.

Unlike *get_software_dependencies()*, it shouldn't need to be overridden by subclasses, since it just collects the results of getting dependencies from the datatypes.

check_ready_to_run()

Called before a module is run, or if the 'check' command is called. This will only be called after all library dependencies have been confirmed ready (see :method:get_software_dependencies).

Essentially, this covers any module-specific checks that used to be in *check_runtime_dependencies()* other than library installation (e.g. checking models exist).

Always call the super class' method if you override.

Returns a list of (name, description) pairs, where the name identifies the problem briefly and the description explains what's missing and (ideally) how to fix it.

get_detailed_status()

Returns a list of strings, containing detailed information about the module's status that is specific to the module type. This may include module-specific information about execution status, for example.

Subclasses may override this to supply useful (human-readable) information specific to the module type. They should call the super method.

reset_execution()

Remove all output data and metadata from this module to make a fresh start, as if it's never been executed.

May be overridden if a module has some side effect other than creating/modifying things in its output directory(/ies), but overridden methods should always call the super method. Occasionally this is necessary, but most of the time the base implementation is enough.

classmethod get_key_info_table()

Add the stages into the key info table.

get_next_stage()

If there are more stages to be executed, returns a pair of the module info and stage definition. Otherwise, returns (None, None)

status

is_locked()

Returns True if the module is currently locked from execution

multistage_module(*multistage_module_type_name*, *module_stages*, *use_stage_option_names=False*, *module_readable_name=None*)

Factory to build a multi-stage module type out of a series of stages, each of which specifies a module type for the stage. The stages should be a list of *ModuleStage* objects.

class ModuleStage(*name*, *module_info_cls*, *connections=None*, *output_connections=None*, *option_connections=None*, *use_stage_option_names=False*)

Bases: object

A single stage in a multi-stage module.

If no explicit input connections are given, the default input to this module is connected to the default output from the previous.

Connections can be given as a list of `ModuleConnections`.

Output connections specify that one of this module's outputs should be used as an output from the multi-stage module. Optional outputs for the multi-stage module are not currently supported (though could in theory be added later). This should be a list of `ModuleOutputConnections`. If none are given for any of the stages, the module will have a single output, which is the default output from the last stage.

Option connections allow you to specify the names that are used for the multistage module's options that get passed through to this stage's module options. Simply specify a dict for `option_connections` where the keys are names module options for this stage and the values are the names that should be used for the multistage module's options.

You may map multiple options from different stages to the same option name for the multistage module. This will result in the same option value being passed through to both stages. Note that help text, option type, option processing, etc will be taken from the first stage's option (in case the two options aren't identical).

Options not explicitly mapped to a name will use the name `<stage_name>_<option_name>`. If `use_stage_option_names=True`, this prefix will not be added: the stage's option names will be used directly as the option name of the multistage module. Note that there is a danger of clashing option names with this behaviour, so only do it if you know the stages have distinct option names (or should share their values where the names overlap).

class ModuleConnection

Bases: `object`

class InternalModuleConnection (*input_name, output_name=None, previous_module=None*)

Bases: `pimlico.core.modules.multistage.ModuleConnection`

Connection between the output of one module in the multi-stage module and the input to another.

May specify the name of the previous module that a connection should be made to. If this is not given, the previous module in the sequence will be assumed.

If `output_name=None`, connects to the default output of the previous module.

class ModuleInputConnection (*stage_input_name=None, main_input_name=None*)

Bases: `pimlico.core.modules.multistage.ModuleConnection`

Connection of a sub-module's input to an input to the multi-stage module.

If `main_input_name` is not given, the name for the input to the multistage module will be identical to the stage input name. This might lead to unintended behaviour if multiple inputs end up with the same name, so you can specify a different name if necessary to avoid clashes.

If multiple inputs (e.g. from different stages) are connected to the same main input name, they will take input from the same previous module output. Nothing clever is done to unify the type requirements, however: the first stage's type requirement is used for the main module's input.

If `stage_input_name` is not given, the module's default input will be connected.

class ModuleOutputConnection (*stage_output_name=None, main_output_name=None*)

Bases: `object`

Specifies the connection of a sub-module's output to the multi-stage module's output. Works in a similar way to `ModuleInputConnection`.

exception MultistageModulePreparationError

Bases: `exceptions.Exception`

pimlico.core.modules.options module

Utilities and type processors for module options.

opt_type_help (*help_text*)

Decorator to add help text to functions that are designed to be used as module option processors. The help text will be used to describe the type in documentation.

opt_type_example (*example_text*)

Decorate to add an example value to function that are designed to be used as module option processors. The given text will be used in module docs as an example of how to specify the option in a config file.

format_option_type (*t*)

str_to_bool (*string*)

Convert a string value to a boolean in a sensible way. Suitable for specifying booleans as options.

Parameters **string** – input string

Returns boolean value

choose_from_list (*options, name=None*)

Utility for option processors to limit the valid values to a list of possibilities.

comma_separated_list (*item_type=<type 'str'>, length=None*)

Option processor type that accepts comma-separated lists of strings. Each value is then parsed according to the given *item_type* (default: string).

comma_separated_strings (*string*)

json_string (*string*)

json_dict (*string*)

JSON dicts, with or without {}s

process_module_options (*opt_def, opt_dict, module_type_name*)

Utility for processing runtime module options. Called from module base class.

Also used when loading a dataset's datatype from datatype options specified in a config file.

Parameters

- **opt_def** – dictionary defining available options
- **opt_dict** – dictionary of option values
- **module_type_name** – name for error output

Returns dictionary of processed options

exception **ModuleOptionParseError**

Bases: `exceptions.Exception`

Module contents

Core functionality for loading and executing different types of pipeline module.

Submodules

pimlico.core.config module

Reading of pipeline config from a file into the data structure used to run and manipulate the pipeline's data.

```
class PipelineConfig(name, pipeline_config, local_config, filename=None, variant='main', available_variants=[], log=None, all_filenames=None, module_aliases={}, local_config_sources=None)
```

Bases: object

Main configuration for a pipeline, read in from a config file.

For details on how to write config files that get read by this class, see [Pipeline config](#).

modules

List of module names, in the order they were specified in the config file.

module_dependencies

Dictionary mapping a module name to a list of the names of modules that it depends on for its inputs.

module_dependents

Opposite of `module_dependencies`. Returns a mapping from module names to a list of modules the depend on the module.

```
get_dependent_modules(module_name, recurse=False, exclude=[])
```

Return a list of the names of modules that depend on the named module for their inputs.

If `exclude` is given, we don't perform a recursive call on any of the modules in the list. For each item we recurse on, we extend the exclude list in the recursive call to include everything found so far (in other recursive calls). This avoids unnecessary recursion in complex pipelines.

If `exclude=None`, it is also passed through to recursive calls as `None`. Its default value of `[]` avoids excessive recursion from the top-level call, by allowing things to be added to the exclusion list for recursive calls.

Parameters `recurse` – include all transitive dependents, not just those that immediately depend on the module.

```
append_module(module_info)
```

Add a moduleinfo to the end of the pipeline. This is mainly for use while loaded a pipeline from a config file.

```
get_module_schedule()
```

Work out the order in which modules should be executed. This is an ordering that respects dependencies, so that modules are executed after their dependencies, but otherwise follows the order in which modules were specified in the config.

Returns list of module names

```
reset_all_modules()
```

Resets the execution states of all modules, restoring the output dirs as if nothing's been run.

```
path_relative_to_config(path)
```

Get an absolute path to a file/directory that's been specified relative to a config file (usually within the config file).

Parameters `path` – relative path

Returns absolute path

```
short_term_store
```

For backwards compatibility: returns output path

long_term_store

For backwards compatibility: return storage location 'long' if it exists, else first storage location

named_storage_locations

store_names

output_path

static load (*filename*, *local_config=None*, *variant='main'*, *override_local_config={}*,
only_override_config=False)

Main function that loads a pipeline from a config file.

Parameters

- **filename** – file to read config from
- **local_config** – location of local config file, where we'll read system-wide config. Usually not specified, in which case standard locations are searched. When loading programmatically, you might want to give this
- **variant** – pipeline variant to load
- **override_local_config** – extra configuration values to override the system-wide config
- **only_override_config** – don't load local config from files, just use that given in *override_local_config*. Used for loading test pipelines

Returns

static load_local_config (*filename=None*, *override={}*, *only_override=False*)

Load local config parameters. These are usually specified in a *.pimlico* file, but may be overridden by other config locations, on the command line, or elsewhere programmatically.

If *only_override=True*, don't load any files, just use the values given in *override*. The various locations for local config files will not be checked (which usually happens when *filename=None*). This is not useful for normal pipeline loading, but is used for loading test pipelines.

static empty (*local_config=None*, *override_local_config={}*, *override_pipeline_config={}*,
only_override_config=False)

Used to programmatically create an empty pipeline. It will contain no modules, but provides a gateway to system info, etc and can be used in place of a real Pimlico pipeline.

Parameters

- **local_config** – filename to load local config from. If not given, the default locations are searched
- **override_local_config** – manually override certain local config parameters. Dict of parameter values
- **only_override_config** – don't load any files, just use the values given in *override*. The various locations for local config files will not be checked (which usually happens when *filename=None*). This is not useful for normal pipeline loading, but is used for loading test pipelines.

Returns the *PipelineConfig* instance

find_data_path (*path*, *default=None*)

Given a path to a data dir/file relative to a data store, tries taking it relative to various store base dirs. If it exists in a store, that absolute path is returned. If it exists in no store, return *None*. If the path is already an absolute path, nothing is done to it.

Searches all the specified storage locations.

Parameters

- **path** – path to data, relative to store base
- **default** – usually, return None if no data is found. If default is given, return the path relative to the named storage location if no data is found. Special value “output” returns path relative to output location, whichever of the storage locations that might be

Returns absolute path to data, or None if not found in any store

find_data_store (*path*, *default=None*)

Like *find_data_path()*, searches through storage locations to see if any of them include the data that lives at this relative path. This method returns the name of the store in which it was found.

Parameters

- **path** – path to data, relative to store base
- **default** – usually, return None if no data is found. If default is given, return the path relative to the named storage location if no data is found. Special value “output” returns path relative to output location, whichever of the storage locations that might be

Returns name of store

find_data (*path*, *default=None*)

Given a path to a data dir/file relative to a data store, tries taking it relative to various store base dirs. If it exists in a store, that absolute path is returned. If it exists in no store, return None. If the path is already an absolute path, nothing is done to it.

Searches all the specified storage locations.

Parameters

- **path** – path to data, relative to store base
- **default** – usually, return None if no data is found. If default is given, return the path relative to the named storage location if no data is found. Special value “output” returns path relative to output location, whichever of the storage locations that might be

Returns (store, path), where store is the name of the store used and path is absolute path to data, or None for both if not found in any store

get_data_search_paths (*path*)

Like *find_all_data_paths()*, but returns a list of all absolute paths which this data path could correspond to, whether or not they exist.

Parameters **path** – relative path within Pimlico directory structures

Returns list of string

step

enable_step ()

Enable super-verbose, interactive step mode.

::seealso:

```
Module :mod:pimlico.cli.debug
    The debug module defines the behaviour of step mode.
```

exception PipelineConfigParseError (**args*, ***kwargs*)

Bases: `exceptions.Exception`

General problems interpreting pipeline config

exception PipelineStructureError

Bases: `exceptions.Exception`

Fundamental structural problems in a pipeline.

exception PipelineCheckError (*cause, *args, **kwargs*)

Bases: `exceptions.Exception`

Error in the process of explicitly checking a pipeline for problems.

preprocess_config_file (*filename, variant='main', initial_vars={}*)

Workhorse of the initial part of config file reading. Deals with all of our custom stuff for pipeline configs, such as preprocessing directives and includes.

Parameters

- **filename** – file from which to read main config
- **variant** – name of a variant to load. The default (*main*) loads the main variant, which always exists
- **initial_vars** – variable assignments to make available for substitution. This will be added to by any *vars* sections that are read.

Returns tuple: raw config dict; list of variants that could be loaded; final vars dict; list of filenames that were read, including included files; dict of docstrings for each config section

check_for_cycles (*pipeline*)

Basic cyclical dependency check, always run on pipeline before use.

check_release (*release_str*)

Check a release name against the current version of Pimlico to determine whether we meet the requirement.

check_pipeline (*pipeline*)

Checks a pipeline over for metadata errors, cycles, module typing errors and other problems. Called every time a pipeline is loaded, to check the whole pipeline's metadata is in order.

Raises a *PipelineCheckError* if anything's wrong.

get_dependencies (*pipeline, modules, recursive=False, sources=False*)

Get a list of software dependencies required by the subset of modules given.

If recursive=True, dependencies' dependencies are added to the list too.

Parameters

- **pipeline** –
- **modules** – list of modules to check. If None, checks all modules

print_missing_dependencies (*pipeline, modules*)

Check runtime dependencies for a subset of modules and output a table of missing dependencies.

Parameters

- **pipeline** –
- **modules** – list of modules to check. If None, checks all modules

Returns True if no missing dependencies, False otherwise

print_dependency_leaf_problems (*dep, local_config*)

pimlico.core.logs module

`get_log_file` (*name*)

Returns the path to a log file that may be used to output helpful logging info. Typically used to output verbose error information if something goes wrong. The file can be found in the Pimlico log dir.

Parameters **name** – identifier to distinguish from other logs

Returns path

pimlico.core.paths module

`abs_path_or_model_dir_path` (*path, model_type*)

Module contents

pimlico.datatypes package

Subpackages

pimlico.datatypes.corpora package

Submodules

pimlico.datatypes.corpora.base module

`class CountInvalidCmd`

Bases: `pimlico.cli.shell.base.ShellCommand`

Data shell command to count up the number of invalid docs in a tarred corpus. Applies to any iterable corpus.

commands = ['invalid']

help_text = 'Count the number of invalid documents in this dataset'

execute (*shell, *args, **kwargs*)

Execute the command. Get the dataset reader as shell.data.

Parameters

- **shell** – DataShell instance. Reader available as shell.data
- **args** – Args given by the user
- **kwargs** – Named args given by the user as key=val

`data_point_type_opt` (*text*)

`class IterableCorpus` (**args, **kwargs*)

Bases: `pimlico.datatypes.base.PimlicoDatatype`

Superclass of all datatypes which represent a dataset that can be iterated over document by document (or datapoint by datapoint - what exactly we're iterating over may vary, though documents are most common).

This is an abstract base class and doesn't provide any mechanisms for storing documents or organising them on disk in any way. Many input modules will override this to provide a reader that iterates over the documents

directly, according to `IterableCorpus`' interface. The main subclass of this used within pipelines is `GroupedCorpus`, which provides an interface for iterating over groups of documents and a storage mechanism for grouping together documents in archives on disk.

May be used as a type requirement, but remember that it is not possible to create a reader from this type directly: use a subtype, like `GroupedCorpus`, instead.

The actual type of the data depends on the type given as the first argument, which should be an instance of `DataPointType` or a subclass: it could be, e.g. `coref` output, etc. Information about the type of individual documents is provided by `data_point_type` and this is used in type checking.

Note that the data point type is the first datatype option, so can be given as the first positional arg when instantiating an iterable corpus subtype:

```
corpus_type = GroupedCorpus(RawTextDocumentType())
corpus_reader = corpus_type("... base dir path ...")
```

At creation time, length should be provided in the metadata, denoting how many documents are in the dataset.

datatype_name = 'iterable_corpus'

shell_commands = [<pimlico.datatypes.corpora.base.CountInvalidCmd object>]

datatype_options = {'data_point_type': {'default': `DataPointType()`, 'type': <functi

run_browser (*reader*, *opts*)

Launches a browser interface for reading this datatype, browsing the data provided by the given reader.

Not all datatypes provide a browser. For those that don't, this method should raise a `NotImplementedError`.

opts provides the argparse options from the command line.

This tool used to be only available for iterable corpora, but now it's possible for any datatype to provide a browser. `IterableCorpus` provides its own browser, as before, which uses one of the data point type's formatters to format documents.

Reader

alias of `IterableCorpusReader`

Writer

alias of `IterableCorpusWriter`

check_type (*supplied_type*)

Override type checking to require that the supplied type have a document type that is compatible with (i.e. a subclass of) the document type of this class.

type_checking_name ()

Supplies a name for this datatype to be used in type-checking error messages. Default implementation just provides the class name. Classes that override `check_supplied_type()` may want to override this too.

full_datatype_name ()

Returns a string/unicode name for the datatype that includes relevant sub-type information. The default implementation just uses the attribute `datatype_name`, but subclasses may have more detailed information to add. For example, iterable corpus types also supply information about the data-point type.

pimlico.datatypes.corpora.data_points module

Document types used to represent datatypes of individual documents in an `IterableCorpus` or subtype.

class `DataPointType`

Bases: `object`

Base data-point type for iterable corpora. All iterable corpora should have data-point types that are subclasses of this.

Every data point type has a corresponding document class, which can be accessed as *MyDataPointType.Document*. When overriding data point types, you can define a nested *Document* class, with no base class, to override parts of the document class' functionality or add new methods, etc. This will be used to automatically create the *Document* class for the data point type.

Note: For now, data point types don't have a way of specifying options (like main datatypes do). I'm not sure whether this is needed, so I'm leaving it out for now. If it is needed, an additional datatype option can be added to iterable corpora that allows you to specify data point type options for when a datatype is being loaded using a config file.

formatters = []

List of (name, cls_path) pairs specifying a standard set of formatters that the user might want to choose from to view a dataset of this type. The user is not restricted to this set, but can easily choose these by name, instead of specifying a class path themselves. The first in the list is the default used if no formatter is specified. Falls back to DefaultFormatter if empty

metadata_defaults = {}

Metadata keys that should be written for this data point type, with default values and strings documenting the meaning of the parameter. Used for writers for this data point type. See *Writer*.

name

is_type_for_doc (*doc*)

Check whether the given document is of this type, or a subclass of this one.

reader_init (*reader*)

Called when a reader is initialized. May be overridden to perform any tasks specific to the data point type that need to be done before the reader starts producing data points.

The super *reader_init()* should be called. This takes care of making reader metadata available in the *metadata* attribute of the data point type instance.

writer_init (*writer*)

Called when a writer is initialized. May be overridden to perform any tasks specific to the data point type that should be done before documents start getting written.

The super *writer_init()* should be called. This takes care of updating the writer's metadata from anything in the instance's *metadata* attribute, for any keys given in the data point type's *metadata_defaults*.

classmethod full_class_name ()

The fully qualified name of the class for this data point type, by which it is referenced in config files. Used in docs

class Document (*data_point_type*, *raw_data=None*, *internal_data=None*)

Bases: *object*

The abstract superclass of all documents.

You do not need to subclass or instantiate these yourself: subclasses are created automatically to correspond to each document type. You can add functionality to a datapoint type's document by creating a nested *Document* class. This will inherit from the parent datapoint type's document. This happens automatically - you don't need to do it yourself and shouldn't inherit from anything:

```
class MyDataPointType(DataPointType):
    class Document:
```

(continues on next page)

(continued from previous page)

```
# Override document things here
# Add your own methods, properties, etc for getting data from the_
↔document
```

A data point type's constructed document class is available as *MyDataPointType.Document*.

Each document type should provide a method to convert from raw data (a unicode string) to the internal representation (an arbitrary dictionary) called *raw_to_internal()*, and another to convert the other way called *internal_to_raw()*. Both forms of the data are available using the properties *raw_data* and *internal_data*, and these methods are called as necessary to convert back and forth.

This is to avoid unnecessary conversions. For example, if the raw data is supplied and then only the raw data is ever used (e.g. passing the document straight through and writing out to disk), we want to avoid converting back and forth.

A subtype should then supply methods or properties (typically using the *cached_property* decorator) to provide access to different parts of the data. See the many built-in document types for examples of doing this.

You should not generally need to override the *__init__* method. You may, however, wish to override *internal_available()* or *raw_available()*. These are called as soon as the internal data or raw data, respectively, become available, which may be at instantiation or after conversion. This can be useful if there are bits of computation that you want to do on the basis of one of these and then store to avoid repeated computation.

keys = []

Specifies the keys that a document has in its internal data. Subclasses should specify their keys. The internal data fields corresponding to these can be accessed as attributes of the document.

raw_to_internal (*raw_data*)

Take a unicode string containing the raw data for a document, read in from disk, and produce a dictionary containing all the processed data in the document's internal format.

You will often want to call the super method and replace values or add to the dictionary. Whatever you do, make sure that all the internal data that the super type provides is also provided here, so that all of its properties and methods work.

internal_to_raw (*internal_data*)

Take a dictionary containing all the document's data in its internal format and produce a unicode string containing all that data, which can be written out to disk.

raw_available ()

Called as soon as the raw data becomes available, either at instantiation or conversion.

internal_available ()

Called as soon as the internal data becomes available, either at instantiation or conversion.

raw_data

internal_data

class InvalidDocument

Bases: *pimlico.datatypes.corpora.data_points.DataPointType*

Widely used in Pimlico to represent an empty document that is empty not because the original input document was empty, but because a module along the way had an error processing it. Document readers/writers should generally be robust to this and simply pass through the whole thing where possible, so that it's always possible to work out, where one of these pops up, where the error occurred.

Document

alias of `InvalidDocumentDocument`

class RawDocumentType

Bases: `pimlico.datatypes.corpora.data_points.DataPointType`

Base document type. All document types for grouped corpora should be subclasses of this.

It may be used itself as well, where documents are just treated as raw data, though most of the time it will be appropriate to use subclasses to provide more information and processing operations specific to the datatype.

Document

alias of `RawDocument`

class TextDocumentType

Bases: `pimlico.datatypes.corpora.data_points.RawDocumentType`

Documents that contain text, most often human-readable documents from a textual corpus. Most often used as a superclass for other, more specific, document types.

This type does not special processing, since the storage format is already a unicode string, which is fine for raw text. However, it serves to indicate that the document represents text (not just any old raw data).

The property `text` provides the text, which is, for this base type, just the raw data. However, subclasses will override this, since their raw data will contain information other than the raw text.

Document

alias of `TextDocument`

class RawTextDocumentType

Bases: `pimlico.datatypes.corpora.data_points.TextDocumentType`

Subclass of `TextDocumentType` used to indicate that the text hasn't been processed (tokenized, etc). Note that text that has been tokenized, parsed, etc does not use subclasses of this type, so they will not be considered compatible if this type is used as a requirement.

Document

alias of `RawTextDocument`

exception DataPointError

Bases: `exceptions.Exception`

pimlico.datatypes.corpora.floats module

Corpora consisting of lists of ints. These data point types are useful, for example, for encoding text or other sequence data as integer IDs. They are designed to be fast to read.

class FloatListsDocumentType

Bases: `pimlico.datatypes.corpora.data_points.RawDocumentType`

Corpus of float list data: each doc contains lists of float. Unlike `IntegerTableDocumentCorpus`, they are not all constrained to have the same length. The downside is that the storage format (and probably I/O speed) isn't quite as efficient. It's still better than just storing ints as strings or JSON objects.

The floats are stored as C double, which use 8 bytes. At the moment, we don't provide any way to change this. An alternative would be to use C floats, losing precision but (almost) halving storage size.

metadata_defaults = {'bytes': (8, 'Number of bytes to use to represent each int. Defa

reader_init (*reader*)

Called when a reader is initialized. May be overridden to perform any tasks specific to the data point type that need to be done before the reader starts producing data points.

The super *reader_init()* should be called. This takes care of making reader metadata available in the *metadata* attribute of the data point type instance.

writer_init (*writer*)

Called when a writer is initialized. May be overridden to perform any tasks specific to the data point type that should be done before documents start getting written.

The super *writer_init()* should be called. This takes care of updating the writer's metadata from anything in the instance's *metadata* attribute, for any keys given in the data point type's *metadata_defaults*.

Document

alias of `FloatListsDocument`

class FloatListDocumentType

Bases: `pimlico.datatypes.corpora.data_points.RawDocumentType`

Corpus of float data: each doc contains a single sequence of floats.

The floats are stored as C doubles, using 8 bytes each.

reader_init (*reader*)

Called when a reader is initialized. May be overridden to perform any tasks specific to the data point type that need to be done before the reader starts producing data points.

The super *reader_init()* should be called. This takes care of making reader metadata available in the *metadata* attribute of the data point type instance.

writer_init (*writer*)

Called when a writer is initialized. May be overridden to perform any tasks specific to the data point type that should be done before documents start getting written.

The super *writer_init()* should be called. This takes care of updating the writer's metadata from anything in the instance's *metadata* attribute, for any keys given in the data point type's *metadata_defaults*.

Document

alias of `FloatListDocument`

class FloatListsFormatter (*corpus_datatype*)

Bases: `pimlico.cli.browser.tools.formatter.DocumentBrowserFormatter`

DATATYPE

alias of `FloatListsDocumentType`

format_document (*doc*)

Format a single document and return the result as a string (or unicode, but it will be converted to ASCII for display).

Must be overridden by subclasses.

class VectorDocumentType

Bases: `pimlico.datatypes.corpora.data_points.RawDocumentType`

Like `FloatListDocumentType`, but each document has the same number of float values.

Each document contains a single list of floats and each one has the same length. That is, each document is one vector.

The floats are stored as C doubles, using 8 bytes each.

formatters = [('vector', 'pimlico.datatypes.floats.VectorFormatter')]

reader_init (*reader*)

Called when a reader is initialized. May be overridden to perform any tasks specific to the data point type that need to be done before the reader starts producing data points.

The super `reader_init()` should be called. This takes care of making reader metadata available in the `metadata` attribute of the data point type instance.

writer_init (*reader*)

Called when a writer is initialized. May be overridden to perform any tasks specific to the data point type that should be done before documents start getting written.

The super `writer_init()` should be called. This takes care of updating the writer's metadata from anything in the instance's `metadata` attribute, for any keys given in the data point type's `metadata_defaults`.

Document

alias of `VectorDocument`

class VectorFormatter (*corpus_datatype*)

Bases: `pimlico.cli.browser.tools.formatter.DocumentBrowserFormatter`

DATATYPE

alias of `VectorDocumentType`

format_document (*doc*)

Format a single document and return the result as a string (or unicode, but it will be converted to ASCII for display).

Must be overridden by subclasses.

pimlico.datatypes.corpora.grouped module

class GroupedCorpus (**args, **kwargs*)

Bases: `pimlico.datatypes.corpora.base.IterableCorpus`

datatype_name = 'grouped_corpus'

document_preprocessors = []

Reader

alias of `GroupedCorpusReader`

Writer

alias of `GroupedCorpusWriter`

class AlignedGroupedCorpora (*readers*)

Bases: `object`

Iterator for iterating over multiple corpora simultaneously that contain the same files, grouped into archives in the same way. This is the standard utility for taking multiple inputs to a Pimlico module that contain different data but for the same corpus (e.g. output of different tools).

archive_iter (*start_after=None, skip=None, name_filter=None*)

class GroupedCorpusWithTypeFromInput (*input_name=None*)

Bases: `pimlico.datatypes.base.DynamicOutputDatatype`

Dynamic datatype that produces a `GroupedCorpus` with a document datatype that is the same as the input's document/data-point type.

If the input name is not given, uses the first input.

Unlike `CorpusWithTypeFromInput`, this does not infer whether the result should be a grouped corpus or not: it always is. The input should be an iterable corpus (or subtype, including grouped corpus), and that's where the datatype will come from.

datatype_name = 'grouped corpus with input doc type'

get_base_datatype_class ()

If it's possible to say before the instance of a ModuleInfo is available what base datatype will be produced, implement this to return the class. By default, it returns None.

If this information is available, it will be used in documentation.

get_datatype (*module_info*)

class CorpusWithTypeFromInput (*input_name=None*)

Bases: *pimlico.datatypes.base.DynamicOutputDatatype*

Infer output corpus' data-point type from the type of an input. Passes the data point type through. Similar to *GroupedCorpusWithTypeFromInput*, but more flexible.

If the input is a grouped corpus, so is the output. Otherwise, it's just an IterableCorpus.

Handles the case where the input is a multiple input. Tries to find a common data point type among the inputs. They must have the same data point type, or all must be subtypes of one of them. (In theory, we could find the most specific common ancestor and use that as the output type, but this is not currently implemented and is probably not worth the trouble.)

Input name may be given. Otherwise, the default input is used.

datatype_name = 'corpus with data-point from input'

get_datatype (*module_info*)

exception CorpusAlignmentError

Bases: *exceptions.Exception*

exception GroupedCorpusIterationError

Bases: *exceptions.Exception*

pimlico.datatypes.corpora.ints module

Corpora consisting of lists of ints. These data point types are useful, for example, for encoding text or other sequence data as integer IDs. They are designed to be fast to read.

class IntegerListsDocumentType

Bases: *pimlico.datatypes.corpora.data_points.RawDocumentType*

Corpus of integer list data: each doc contains lists of ints. Unlike *IntegerTableDocumentType*, they are not all constrained to have the same length. The downside is that the storage format (and I/O speed) isn't quite as good. It's still better than just storing ints as strings or JSON objects.

By default, the ints are stored as C longs, which use 4 bytes. If you know you don't need ints this big, you can choose 1 or 2 bytes, or even 8 (long long). By default, the ints are unsigned, but they may be signed.

metadata_defaults = {'bytes': (8, 'Number of bytes to use to represent each int. Defa

reader_init (*reader*)

Called when a reader is initialized. May be overridden to perform any tasks specific to the data point type that need to be done before the reader starts producing data points.

The super *reader_init()* should be called. This takes care of making reader metadata available in the *metadata* attribute of the data point type instance.

writer_init (*writer*)

Called when a writer is initialized. May be overridden to perform any tasks specific to the data point type that should be done before documents start getting written.

The super `writer_init()` should be called. This takes care of updating the writer's metadata from anything in the instance's `metadata` attribute, for any keys given in the data point type's `metadata_defaults`.

struct

length_struct

Document

alias of IntegerListsDocument

class IntegerListDocumentType

Bases: `pimlico.datatypes.corpora.data_points.RawDocumentType`

Corpus of integer data: each doc contains a single sequence of ints.

Like IntegerListsDocumentType, but each document is treated as a single list of integers.

By default, the ints are stored as C longs, which use 4 bytes. If you know you don't need ints this big, you can choose 1 or 2 bytes, or even 8 (long long). By default, the ints are unsigned, but they may be signed.

metadata_defaults = {'bytes': (8, 'Number of bytes to use to represent each int. Defa

reader_init (*reader*)

Called when a reader is initialized. May be overridden to perform any tasks specific to the data point type that need to be done before the reader starts producing data points.

The super `reader_init()` should be called. This takes care of making reader metadata available in the `metadata` attribute of the data point type instance.

writer_init (*writer*)

Called when a writer is initialized. May be overridden to perform any tasks specific to the data point type that should be done before documents start getting written.

The super `writer_init()` should be called. This takes care of updating the writer's metadata from anything in the instance's `metadata` attribute, for any keys given in the data point type's `metadata_defaults`.

struct

Document

alias of IntegerListDocument

pimlico.datatypes.corpora.table module

Corpora where each document is a table, i.e. a list of lists, where each row has the same length and each column has a single datatype. This is designed to be fast to read, but is not a very flexible datatype.

get_struct (*bytes, signed, row_length*)

class IntegerTableDocumentType

Bases: `pimlico.datatypes.corpora.data_points.RawDocumentType`

Corpus of tabular integer data: each doc contains rows of ints, where each row contains the same number of values. This allows a more compact representation, which doesn't require converting the ints to strings or scanning for line ends, so is quite a bit quicker and results in much smaller file sizes. The downside is that the files are not human-readable.

By default, the ints are stored as C longs, which use 4 bytes. If you know you don't need ints this big, you can choose 1 or 2 bytes, or even 8 (long long). By default, the ints are unsigned, but they may be signed.

metadata_defaults = {'bytes': (8, 'Number of bytes to use to represent each int. Defa

reader_init (*reader*)

Called when a reader is initialized. May be overridden to perform any tasks specific to the data point type that need to be done before the reader starts producing data points.

The super *reader_init()* should be called. This takes care of making reader metadata available in the *metadata* attribute of the data point type instance.

writer_init (*writer*)

Called when a writer is initialized. May be overridden to perform any tasks specific to the data point type that should be done before documents start getting written.

The super *writer_init()* should be called. This takes care of updating the writer's metadata from anything in the instance's *metadata* attribute, for any keys given in the data point type's *metadata_defaults*.

Document

alias of IntegerTableDocument

pimlico.datatypes.corpora.tokenized module

class TokenizedDocumentType

Bases: *pimlico.datatypes.corpora.data_points.TextDocumentType*

Specialized data point type for documents that have had tokenization applied. It does very little processing - the main reason for its existence is to allow modules to require that a corpus has been tokenized before it's given as input.

Each document is a list of sentences. Each sentence is a list of words.

Document

alias of TokenizedDocument

class CharacterTokenizedDocumentType

Bases: *pimlico.datatypes.corpora.tokenized.TokenizedDocumentType*

Simple character-level tokenized corpus. The text isn't stored in any special way, but is represented when read internally just as a sequence of characters in each sentence.

If you need a more sophisticated way to handle character-type (or any non-word) units within each sequence, see *SegmentedLinesDocumentType*.

Document

alias of CharacterTokenizedDocument

class SegmentedLinesDocumentType

Bases: *pimlico.datatypes.corpora.tokenized.TokenizedDocumentType*

Document consisting of lines, each split into elements, which may be characters, words, or whatever. Rather like a tokenized corpus, but doesn't make the assumption that the elements (words in the case of a tokenized corpus) don't include spaces.

You might use this, for example, if you want to train character-level models on a text corpus, but don't use strictly single-character units, perhaps grouping together certain short character sequences.

Uses the character / to separate elements in the raw data. If a / is found in an element, it is stored as *@slash@*, so this string is assumed not to be used in any element (which seems reasonable enough, generally).

Document

alias of SegmentedLinesDocument

Module contents

Submodules

pimlico.datatypes.arrays module

Wrappers around Numpy arrays and Scipy sparse matrices.

class NumpyArray (*args, **kwargs)

Bases: *pimlico.datatypes.files.NamedFileCollection*

datatype_name = 'numpy_array'

get_software_dependencies ()

Get a list of all software required to **read** this datatype. This is separate to metadata config checks, so that you don't need to satisfy the dependencies for all modules in order to be able to run one of them. You might, for example, want to run different modules on different machines. This is called when a module is about to be executed and each of the dependencies is checked.

Returns a list of instances of subclasses of :class:`~pimlico.core.dependencies.base.SoftwareDependency`, representing the libraries that this module depends on.

Take care when providing dependency classes that you don't put any import statements at the top of the Python module that will make loading the dependency type itself dependent on runtime dependencies. You'll want to run import checks by putting import statements within this method.

You should call the super method for checking superclass dependencies.

Note that there may be different software dependencies for **writing** a datatype using its *Writer*. These should be specified using *get_writer_software_dependencies()*.

Reader

alias of NumpyArrayReader

Writer

alias of NumpyArrayWriter

class ScipySparseMatrix (*args, **kwargs)

Bases: *pimlico.datatypes.files.NamedFileCollection*

Wrapper around Scipy sparse matrices. The matrix loaded is always in COO format – you probably want to convert to something else before using it. See scipy docs on sparse matrix conversions.

datatype_name = 'scipy_sparse_array'

get_software_dependencies ()

Get a list of all software required to **read** this datatype. This is separate to metadata config checks, so that you don't need to satisfy the dependencies for all modules in order to be able to run one of them. You might, for example, want to run different modules on different machines. This is called when a module is about to be executed and each of the dependencies is checked.

Returns a list of instances of subclasses of :class:`~pimlico.core.dependencies.base.SoftwareDependency`, representing the libraries that this module depends on.

Take care when providing dependency classes that you don't put any import statements at the top of the Python module that will make loading the dependency type itself dependent on runtime dependencies. You'll want to run import checks by putting import statements within this method.

You should call the super method for checking superclass dependencies.

Note that there may be different software dependencies for **writing** a datatype using its *Writer*. These should be specified using *get_writer_software_dependencies()*.

Reader

alias of `ScipySparseMatrixReader`

Writer

alias of `ScipySparseMatrixWriter`

pimlico.datatypes.base module

Datatypes provide interfaces for reading and writing datasets. They provide different ways of reading in or iterating over datasets and different ways to write out datasets, as appropriate to the datatype. They are used by Pimlico to typecheck connections between modules to make sure that the output from one module provides a suitable type of data for the input to another. They are then also used by the modules to read in their input data coming from earlier in a pipeline and to write out their output data, to be passed to later modules.

As much as possible, Pimlico pipelines should use standard datatypes to connect up the output of modules with the input of others. Most datatypes have a lot in common, which should be reflected in their sharing common base classes. Input modules take care of reading in data from external sources and they provide access to that data in a way that is identified by a Pimlico datatype.

Instances of subclasses of *PimlicoDatatype* represent the type of datasets and are used for typechecking in a pipeline. Each datatype has an associated *Reader* class, accessed by *datatype_cls.Reader*. These are created automatically and can be instantiated via the datatype (by calling it). They are all subclasses of *PimlicoDatatype.Reader*. It is these readers that are used within a pipeline to read a dataset output by an earlier module. In some cases, other readers may be used: for example, input modules provide standard datatypes at their outputs, but use special readers to provide access to the external data.

A similar reflection of the datatype hierarchy is used for dataset writers, which are used to write the outputs from modules, to be passed to subsequent modules. These are created automatically, just like readers, and are all subclasses of *PimlicoDatatype.Writer*. You can get a datatype's standard writer class via *datatype_cls.Writer*. Some datatypes might not provide a writer, but most do.

class PimlicoDatatype (*args, **kwargs)

Bases: object

The abstract superclass of all datatypes. Provides basic functionality for identifying where data should be stored and such.

Datatypes are used to specify the routines for reading the output from modules, via their reader class.

module is the `ModuleInfo` instance for the pipeline module that this datatype was produced by. It may be `None`, if the datatype wasn't instantiated by a module. It is not required to be set if you're instantiating a datatype in some context other than module output. It should generally be set for input datatypes, though, since they are treated as being created by a special input module.

Creating a new datatype

This is the typical process for creating a new datatype. Of course, some datatypes do more and some of the following is not always necessary, but it's a good guide for reference.

1. Create the datatype class, which may subclass *PimlicoDatatype* or some other existing datatype.
2. Specify a *datatype_name* as a class attribute.
3. Specify software dependencies for reading the data, if any, by overriding *get_software_dependencies()* (calling the super method as well).

4. Specify software dependencies for writing the data, if any that are not among the reading dependencies, by overriding `get_writer_software_dependencies()`
5. Define a nested *Reader* class to add any methods to the reader for this datatype. The data should be read from the directory given by its `data_dir`. It should provide methods for getting different bits of the data, iterating over it, or whatever is appropriate.
6. Define a nested *Setup* class within the reader with a `data_ready(base_dir)` method to check whether the data in `base_dir` is ready to be read using the reader. If all that this does is check the existence of particular filenames or paths within the data dir, you can instead implement the *Setup* class' `get_required_paths()` method to return the paths relative to the data dir.
7. Define a nested *Writer* class in the datatype to add any methods to the writer for this datatype. The data should be written to the path given by its `data_dir`. Provide methods that the user can call to write things to the dataset. Required elements of the dataset should be specified as a list of strings as the `required_tasks` attribute and ticked off as written using `task_complete()`
8. You may want to specify:
 - `datatype_options`: an `OrderedDict` of option definitions
 - `shell_commands`: a list of shell commands associated with the datatype

```
datatype_options = {}
```

Override to provide shell commands specific to this datatype. Should include the superclass' list.

```
shell_commands = []
```

```
datatype_name = 'base_datatype'
```

Options specified in the same way as module options that control the nature of the datatype. These are not things to do with reading of specific datasets, for which the dataset's metadata should be used. These are things that have an impact on typechecking, such that options on the two checked datatypes are required to match for the datatypes to be considered compatible.

They should always be an ordered dict, so that they can be specified using positional arguments as well as kwargs and config parameters.

```
get_software_dependencies ()
```

Get a list of all software required to **read** this datatype. This is separate to metadata config checks, so that you don't need to satisfy the dependencies for all modules in order to be able to run one of them. You might, for example, want to run different modules on different machines. This is called when a module is about to be executed and each of the dependencies is checked.

Returns a list of instances of subclasses of `:class:~pimlico.core.dependencies.base.SoftwareDependency`, representing the libraries that this module depends on.

Take care when providing dependency classes that you don't put any import statements at the top of the Python module that will make loading the dependency type itself dependent on runtime dependencies. You'll want to run import checks by putting import statements within this method.

You should call the super method for checking superclass dependencies.

Note that there may be different software dependencies for **writing** a datatype using its *Writer*. These should be specified using `get_writer_software_dependencies()`.

```
get_writer_software_dependencies ()
```

Get a list of all software required to **write** this datatype using its *Writer*. This works in a similar way to `get_software_dependencies()` (for the *Reader*) and the dependencies will be check before the writer is instantiated.

It is assumed that all the reader's dependencies also apply to the writer, so this method only needs to specify any additional dependencies the writer has.

You should call the super method for checking superclass dependencies.

Todo: Call `get_writer_software_dependencies` before instantiating writer

get_writer (*base_dir*, *pipeline*, *module=None*, ***kwargs*)

Instantiate a writer to write data to the given base dir.

Kwargs are passed through to the writer and used to specify initial metadata and writer params.

Parameters

- **base_dir** – output dir to write dataset to
- **pipeline** – current pipeline
- **module** – module name (optional, for debugging only)

Returns instance of the writer subclass corresponding to this datatype

classmethod instantiate_from_options (*options={}*)

Given string options e.g. from a config file, perform option processing and instantiate datatype

classmethod datatype_full_class_name ()

The fully qualified name of the class for this datatype, by which it is reference in config files. Generally, datatypes don't need to override this, but type requirements that take the place of datatypes for type checking need to provide it.

check_type (*supplied_type*)

Method used by datatype type-checking algorithm to determine whether a supplied datatype (given as an instance of a subclass of `PimlicoDatatype`) is compatible with the present datatype, which is being treated as a type requirement.

Typically, the present class is a type requirement on a module input and *supplied_type* is the type provided by a previous module's output.

The default implementation simply checks whether *supplied_type* is a subclass of the present class. Subclasses may wish to impose different or additional checks.

Parameters **supplied_type** – type provided where the present class is required, or datatype instance

Returns True if the check is successful, False otherwise

type_checking_name ()

Supplies a name for this datatype to be used in type-checking error messages. Default implementation just provides the class name. Classes that override `check_supplied_type()` may want to override this too.

full_datatype_name ()

Returns a string/unicode name for the datatype that includes relevant sub-type information. The default implementation just uses the attribute *datatype_name*, but subclasses may have more detailed information to add. For example, iterable corpus types also supply information about the data-point type.

run_browser (*reader*, *opts*)

Launches a browser interface for reading this datatype, browsing the data provided by the given reader.

Not all datatypes provide a browser. For those that don't, this method should raise a `NotImplementedError`.

opts provides the argparse options from the command line.

This tool used to be only available for iterable corpora, but now it's possible for any datatype to provide a browser. `IterableCorpus` provides its own browser, as before, which uses one of the data point type's formatters to format documents.

Readeralias of `PimlicoDatatypeReader`**Writer**alias of `PimlicoDatatypeWriter`**class DynamicOutputDatatype**Bases: `object`

Types of module outputs may be specified as an instance of a subclass of `PimlicoDatatype`, or alternatively as an instance of `DynamicOutputType`. In this case, `get_datatype()` is called when the output datatype is needed, passing in the module info instance for the module, so that a specialized datatype can be produced on the basis of options, input types, etc.

The dynamic type must provide certain pieces of information needed for typechecking.

datatype_name = None**get_datatype** (*module_info*)**get_base_datatype_class** ()

If it's possible to say before the instance of a `ModuleInfo` is available what base datatype will be produced, implement this to return the class. By default, it returns `None`.

If this information is available, it will be used in documentation.

class DynamicInputDatatypeRequirementBases: `object`

Types of module inputs may be given as an instance of a subclass of `PimlicoDatatype`, a tuple of datatypes, or an instance a `DynamicInputDatatypeRequirement` subclass. In this case, `check_type(supplied_type)` is called during typechecking to check whether the type that we've got conforms to the input type requirements.

Additionally, if `datatype_doc_info` is provided, it is used to represent the input type constraints in documentation.

datatype_doc_info = None**check_type** (*supplied_type*)**type_checking_name** ()

Supplies a name for this datatype to be used in type-checking error messages. Default implementation just provides the class name. Subclasses may want to override this too.

class MultipleInputs (*datatype_requirements*)Bases: `object`

An input datatype that can be used as an item in a module's inputs, which lets the module accept an unbounded number of inputs, all satisfying the same datatype requirements. When writing the inputs in a config file, they can be specified as a comma-separated list of the usual type of specification (module name, with optional output name). Each item in the list must point to a datatype that satisfies the type-checking.

The list may also include (or entirely consist of) a base module name from the pipeline that has been expanded into multiple modules according to alternative parameters (the type separated by vertical bars, see [Multiple parameter values](#)). Use the notation `*name`, where `name` is the base module name, to denote all of the expanded module names as inputs. These are treated as if you'd written out all of the expanded module names separated by commas.

In a config file, if you need the same input specification to be repeated multiple times in a list, instead of writing it out explicitly you can use a multiplier to repeat it `N` times by putting `*N` after it. This is particularly useful when `N` is the result of expanding module variables, allowing the number of times an input is repeated to depend on some modvar expression.

When `get_input()` is called on the module, instead of returning a single datatype, a list of datatypes is returned.

```
exception DatatypeLoadError  
    Bases: exceptions.Exception  
exception DatatypeWriteError  
    Bases: exceptions.Exception
```

pimlico.datatypes.core module

Some basic core datatypes that are commonly used for passing simple data, like strings and dicts, through pipelines.

```
class Dict (*args, **kwargs)  
    Bases: pimlico.datatypes.base.PimlicoDatatype
```

Simply stores a Python dict, pickled to disk. All content in the dict should be pickleable.

```
datatype_name = 'dict'
```

```
Reader  
    alias of DictReader
```

```
Writer  
    alias of DictWriter
```

```
class StringList (*args, **kwargs)  
    Bases: pimlico.datatypes.base.PimlicoDatatype
```

Simply stores a Python list of strings, written out to disk in a readable form. Not the most efficient format, but if the list isn't humungous it's OK (e.g. storing vocabularies).

```
datatype_name = 'string_list'
```

```
Reader  
    alias of StringListReader
```

```
Writer  
    alias of StringListWriter
```

pimlico.datatypes.dictionary module

This module implements the concept of a Dictionary – a mapping between words and their integer ids.

The implementation is based on Gensim, because Gensim is wonderful and there's no need to reinvent the wheel. We don't use Gensim's data structure directly, because it's unnecessary to depend on the whole of Gensim just for one data structure.

However, it is possible to retrieve a Gensim dictionary directly from the Pimlico data structure if you need to use it with Gensim.

```
class Dictionary (*args, **kwargs)  
    Bases: pimlico.datatypes.base.PimlicoDatatype
```

Dictionary encapsulates the mapping between normalized words and their integer ids. This class is responsible for reading and writing dictionaries.

`DictionaryData` is the data structure itself, which is very closely related to Gensim's dictionary.

```
datatype_name = 'dictionary'
```

Readeralias of `DictionaryReader`**Writer**alias of `DictionaryWriter`**pimlico.datatypes.embeddings module**

Datatypes to store embedding vectors, together with their words.

The main datatype here, *Embeddings*, is the main datatype that should be used for passing embeddings between modules.

We also provide a simple file collection datatype that stores the files used by Tensorflow, for example, as input to the Tensorflow Projector. Modules that need data in this format can use this datatype, which makes it easy to convert from other formats.

```
class Embeddings (*args, **kwargs)
```

Bases: `pimlico.datatypes.base.PimlicoDatatype`

Datatype to store embedding vectors, together with their words. Based on Gensim's `KeyedVectors` object, but adapted for use in Pimlico and so as not to depend on Gensim. (This means that this can be used more generally for storing embeddings, even when we're not depending on Gensim.)

Provides a method to map to Gensim's `KeyedVectors` type for compatibility.

Doesn't provide all of the functionality of `KeyedVectors`, since the main purpose of this is for storage of vectors and other functionality, like similarity computations, can be provided by utilities or by direct use of Gensim.

```
datatype_name = 'embeddings'
```

```
get_software_dependencies ()
```

Get a list of all software required to **read** this datatype. This is separate to metadata config checks, so that you don't need to satisfy the dependencies for all modules in order to be able to run one of them. You might, for example, want to run different modules on different machines. This is called when a module is about to be executed and each of the dependencies is checked.

Returns a list of instances of subclasses of `:class:~pimlico.core.dependencies.base.SoftwareDependency`, representing the libraries that this module depends on.

Take care when providing dependency classes that you don't put any import statements at the top of the Python module that will make loading the dependency type itself dependent on runtime dependencies. You'll want to run import checks by putting import statements within this method.

You should call the super method for checking superclass dependencies.

Note that there may be different software dependencies for **writing** a datatype using its *Writer*. These should be specified using `get_writer_software_dependencies()`.

```
get_writer_software_dependencies ()
```

Get a list of all software required to **write** this datatype using its *Writer*. This works in a similar way to `get_software_dependencies()` (for the *Reader*) and the dependencies will be check before the writer is instantiated.

It is assumed that all the reader's dependencies also apply to the writer, so this method only needs to specify any additional dependencies the writer has.

You should call the super method for checking superclass dependencies.

Todo: Call `get_writer_software_dependencies` before instantiating writer

Reader

alias of `EmbeddingsReader`

Writer

alias of `EmbeddingsWriter`

class `TSVVecFiles` (*args, **kwargs)

Bases: `pimlico.datatypes.files.NamedFileCollection`

Embeddings stored in TSV files. This format is used by Tensorflow and can be used, for example, as input to the Tensorflow Projector.

It's just a TSV file with each vector on a row, and another metadata TSV file with the names associated with the points and the counts. The counts are not necessary, so the metadata can be written without them if necessary.

datatype_name = `'tsv_vec_files'`

Reader

alias of `TSVVecFilesReader`

Writer

alias of `TSVVecFilesWriter`

pimlico.datatypes.features module

class `ScoredRealFeatureSets` (*args, **kwargs)

Bases: `pimlico.datatypes.files.NamedFileCollection`

Sets of features, where each feature has an associated real number value, and each set (i.e. data point) has a score.

This is suitable as training data for a multidimensional regression.

Stores a dictionary of feature types and uses integer IDs to refer to them in the data storage.

Todo: Add unit test for `ScoredReadFeatureSets`

datatype_name = `'scored_real_feature_sets'`

browse_file (*reader, filename*)

Return text for a particular file in the collection to show in the browser. By default, just reads in the file's data and returns it, but subclasses might want to override this (perhaps conditioned on the filename) to format the data readably.

Parameters

- **reader** –
- **filename** –

Returns file data to show

Reader

alias of `ScoredRealFeatureSetsReader`

Writer

alias of `ScoredRealFeatureSetsWriter`

pimlico.datatypes.files module

File collections and files.

There used to be an `UnnamedFileCollection`, which has been removed in the move to the new datatype system. It used to be used mostly for input datatypes, which don't exist any more. There may still be a use for this, though, so I may be added in future.

class NamedFileCollection (*args, **kwargs)

Bases: `pimlico.datatypes.base.PimlicoDatatype`

Datatypes that stores a fixed collection of files, which have fixed names (or at least names that can be determined from the class). Very many datatypes fall into this category. Overriding this base class provides them with some common functionality, including the possibility of creating a union of multiple datatypes.

The datatype option `filenames` should specify a list of filenames contained by the datatype. For typechecking, the provided type must have at least all the filenames of the type requirement, though it may include more.

All files are contained in the datatypes data directory. If files are stored in subdirectories, this may be specified in the list of filenames using `/` `s`. (Always use forward slashes, regardless of the operating system.)

datatype_name = 'named_file_collection'

datatype_options = {'filenames': {'default': [], 'type': <function _fn at 0x7ff2eba...

check_type (supplied_type)

Method used by datatype type-checking algorithm to determine whether a supplied datatype (given as an instance of a subclass of `PimlicoDatatype`) is compatible with the present datatype, which is being treated as a type requirement.

Typically, the present class is a type requirement on a module input and `supplied_type` is the type provided by a previous module's output.

The default implementation simply checks whether `supplied_type` is a subclass of the present class. Subclasses may wish to impose different or additional checks.

Parameters `supplied_type` – type provided where the present class is required, or datatype instance

Returns True if the check is successful, False otherwise

browse_file (reader, filename)

Return text for a particular file in the collection to show in the browser. By default, just reads in the file's data and returns it, but subclasses might want to override this (perhaps conditioned on the filename) to format the data readably.

Parameters

- **reader** –
- **filename** –

Returns file data to show

run_browser (reader, opts)

All `NamedFileCollections` provide a browser that just lets you see a list of the files and view them, in the case of text files.

Subclasses may override the way individual files are shown by overriding `browse_file()`.

Reader

alias of `NamedFileCollectionReader`

Writer

alias of NamedFileCollectionWriter

```
class NamedFile (*args, **kwargs)
```

Bases: *pimlico.datatypes.files.NamedFileCollection*

Like NamedFileCollection, but always has exactly one file.

The filename is given as the *filename* datatype option, which can also be given as the first init arg: *NamedFile("myfile.txt")*.

Since NamedFile is a subtype of NamedFileCollection, it also has a “filenames” option. It is ignored if the *filename* option is given, and otherwise must have exactly one item.

```
datatype_name = 'named_file'
```

```
datatype_options = {'filename': {'help': "The file's name"}, 'filenames': {'default
```

Reader

alias of NamedFileReader

Writer

alias of NamedFileWriter

```
class FilesInput (min_files=1)
```

Bases: *pimlico.datatypes.base.DynamicInputDatatypeRequirement*

```
datatype_doc_info = 'A file collection containing at least one file (or a given specif
```

```
check_type (supplied_type)
```

FileInput

alias of *pimlico.datatypes.files.FilesInput*

```
class TextFile (*args, **kwargs)
```

Bases: *pimlico.datatypes.files.NamedFile*

Simple dataset containing just a single utf-8 encoded text file.

```
datatype_name = 'text_document'
```

```
datatype_options = {'filename': {'default': 'data.txt', 'help': "The file's name. T
```

Reader

alias of TextFileReader

Writer

alias of TextFileWriter

pimlico.datatypes.gensim module

```
class GensimLdaModel (*args, **kwargs)
```

Bases: *pimlico.datatypes.base.PimlicoDatatype*

```
datatype_name = 'lda_model'
```

```
get_software_dependencies ()
```

Get a list of all software required to **read** this datatype. This is separate to metadata config checks, so that you don't need to satisfy the dependencies for all modules in order to be able to run one of them. You might, for example, want to run different modules on different machines. This is called when a module is about to be executed and each of the dependencies is checked.

Returns a list of instances of subclasses of `:class:~pimlico.core.dependencies.base.SoftwareDependency`, representing the libraries that this module depends on.

Take care when providing dependency classes that you don't put any import statements at the top of the Python module that will make loading the dependency type itself dependent on runtime dependencies. You'll want to run import checks by putting import statements within this method.

You should call the super method for checking superclass dependencies.

Note that there may be different software dependencies for **writing** a datatype using its *Writer*. These should be specified using `get_writer_software_dependencies()`.

Reader

alias of `GensimLdaModelReader`

Writer

alias of `GensimLdaModelWriter`

pimlico.datatypes.sklearn module

class SklearnModel (*args, **kwargs)

Bases: `pimlico.datatypes.files.NamedFile`

Saves and loads scikit-learn models using the library's joblib functions.

See [the sklearn docs for more details](#)

datatype_name = 'sklearn_model'

get_software_dependencies ()

Get a list of all software required to **read** this datatype. This is separate to metadata config checks, so that you don't need to satisfy the dependencies for all modules in order to be able to run one of them. You might, for example, want to run different modules on different machines. This is called when a module is about to be executed and each of the dependencies is checked.

Returns a list of instances of subclasses of `:class:~pimlico.core.dependencies.base.SoftwareDependency`, representing the libraries that this module depends on.

Take care when providing dependency classes that you don't put any import statements at the top of the Python module that will make loading the dependency type itself dependent on runtime dependencies. You'll want to run import checks by putting import statements within this method.

You should call the super method for checking superclass dependencies.

Note that there may be different software dependencies for **writing** a datatype using its *Writer*. These should be specified using `get_writer_software_dependencies()`.

Reader

alias of `SklearnModelReader`

Writer

alias of `SklearnModelWriter`

Module contents

load_datatype (path, options={})

Try loading a datatype class for a given path. Raises a `DatatypeLoadError` if it's not a valid datatype path. Also looks up class names of builtin datatypes and datatype names.

Options are unprocessed strings that will be processed using the datatype's option definitions.

pimlico.test package

Submodules

pimlico.test.pipeline module

Pipeline tests

Pimlico modules and datatypes cannot always be easily tested with unit tests and where they can it's often not easy to work out how to write the tests in a neatly packaged way. Instead, modules can package up tests in the form of a small pipeline that comes with a tiny dataset to use as input. The pipeline can be run in a test environment, where software dependencies are installed and local config is prepared to store output and so on.

This way of providing tests also has the advantage that modules at the same time provide a demo (or several) of how to use them – how pipeline config should look and what sort of input data to use.

class TestPipeline (*pipeline, run_modules, log*)

Bases: object

static load_pipeline (*path, storage_root*)

Load a test pipeline from a config file.

Path may be absolute, or given relative to Pimlico test data directory (`PIMLICO_ROOT/test/data`)

get_uninstalled_dependencies ()

test_all_modules ()

test_input_module (*module_name*)

test_module_execution (*module_name*)

run_test_pipeline (*path, module_names, log, no_clean=False*)

Run a test pipeline, loading the pipeline config from a given path (which may be relative to the Pimlico test data directory) and running each of the named modules, including any of those modules' dependencies.

Any software dependencies not already available that can be installed automatically will be installed in the current environment. If there are unsatisfied dependencies that can't be automatically installed, an error will be raised.

If any of the modules name explicitly is an input dataset, it is loaded and `data_ready()` is checked. If it is an `IterableCorpus`, it is tested simply by iterating over the full corpus.

run_test_suite (*pipelines_and_modules, log, no_clean=False*)

Parameters pipeline_and_modules – list of (pipeline, modules) pairs, where pipeline is a path to a config file and modules a list of module names to test

clear_storage_dir ()

exception TestPipelineRunError

Bases: `exceptions.Exception`

pimlico.test.suite module

Module contents

pimlico.utils package

Subpackages

pimlico.utils.docs package

Submodules

pimlico.utils.docs.commandgen module

pimlico.utils.docs.modulegen module

pimlico.utils.docs.rest module

make_table (*grid, header=None*)

table_div (*col_widths, header_flag=False*)

normalize_cell (*string, length*)

pimlico.utils.docs.testgen module

Module contents

trim_docstring (*docstring*)

Submodules

pimlico.utils.communicate module

timeout_process (**args, **kws*)

Context manager for use in a *with* statement. If the with block hasn't completed after the given number of seconds, the process is killed.

Parameters **proc** – process to kill if timeout is reached before end of block

Returns

terminate_process (*proc, kill_time=None*)

Ends a process started with subprocess. Tries killing, then falls back on terminating if it doesn't work.

Parameters

- **kill_time** – time to allow the process to be killed before falling back on terminating
- **proc** – Popen instance

Returns

class StreamCommunicationPacket (*data*)

Bases: object

length

encode ()

static read (*stream*)

exception StreamCommunicationError

Bases: exceptions.Exception

pimlico.utils.core module

multiwith (**args, **kws*)

Taken from contextlib's nested(). We need the variable number of context managers that this function allows.

is_identifier (*ident*)

Determines if string is valid Python identifier.

remove_duplicates (*lst, key=<function <lambda>>*)

Remove duplicate values from a list, keeping just the first one, using a particular key function to compare them.

infinite_cycle (*iterable*)

Iterate infinitely over the given iterable.

Watch out for calling this on a generator or iter: they can only be iterated over once, so you'll get stuck in an infinite loop with no more items yielded once you've gone over it once.

You may also specify a callable, in which case it will be called each time to get a new iterable/iterator. This is useful in the case of generator functions.

Parameters **iterable** – iterable or generator to loop over indefinitely

import_member (*path*)

Import a class, function, or other module member by its fully-qualified Python name.

Parameters **path** – path to member, including full package path and class/function/etc name

Returns cls

split_seq (*seq, separator, ignore_empty_final=False*)

Iterate over a sequence and group its values into lists, separated in the original sequence by the given value. If *on* is callable, it is called on each element to test whether it is a separator. Otherwise, elements that are equal to *on* are treated as separators.

Parameters

- **seq** – sequence to divide up
- **separator** – separator or separator test function
- **ignore_empty_final** – by default, if there's a separator at the end, the last sequence yielded is empty. If `ignore_empty_final=True`, in this case the last empty sequence is dropped

Returns iterator over subsequences

split_seq_after (*seq, separator*)

Somewhat like `split_seq`, but starts a new subsequence after each separator, without removing the separators. Each subsequence therefore ends with a separator, except the last one if there's no separator at the end.

Parameters

- **seq** – sequence to divide up
- **separator** – separator or separator test function

Returns iterator over subsequences

chunk_list (*lst, length*)

Divides a list into chunks of max *length* length.

class cached_property (*func*)

Bases: `object`

A property that is only computed once per instance and then replaces itself with an ordinary attribute. Deleting the attribute resets the property.

Often useful in Pimlico datatypes, where it can be time-consuming to load data, but we can't do it once when the datatype is first loaded, since the data might not be ready at that point. Instead, we can access the data, or particular parts of it, using properties and easily cache the result.

Taken from: <https://github.com/bottlepy/bottle>

pimlico.utils.email module

Email sending utilities

Configure email sending functionality by adding the following fields to your Pimlico local config file:

email_sender From-address for all sent emails

email_recipients To-addresses, separated by commas. All notification emails will be sent to all recipients

email_host (optional) Hostname of your SMTP server. Defaults to *localhost*

email_username (optional) Username to authenticate with your SMTP server. If not given, it is assumed that no authentication is required

email_password (optional) Password to authenticate with your SMTP server. Must be supplied if *username* is given

class EmailConfig (*sender=None, recipients=None, host=None, username=None, password=None*)

Bases: `object`

classmethod from_local_config (*local_config*)

send_pimlico_email (*subject, content, local_config, log*)

Primary method for sending emails from Pimlico. Tries to send an email with the given content, using the email details found in the local config. If something goes wrong, an error is logged on the given log.

Parameters

- **subject** – email subject
- **content** – email text (may be unicode)
- **local_config** – local config dictionary
- **log** – logger to log errors to (and info if the sending works)

send_text_email (*email_config, subject, content=None*)

exception EmailError

Bases: `exceptions.Exception`

pimlico.utils.filesystem module

dirsize (*path*)

Recursively compute the size of the contents of a directory.

Parameters *path* –

Returns size in bytes

format_file_size (*bytes*)

copy_dir_with_progress (*source_dir*, *target_dir*, *move=False*)

Utility for moving/copying a large directory and displaying a progress bar showing how much is copied.

Note that the directory is first copied, then the old directory is removed, if *move=True*.

Parameters

- **source_dir** –
- **target_dir** –

Returns

move_dir_with_progress (*source_dir*, *target_dir*)

new_filename (*directory*, *initial_filename='tmp_file'*)

Generate a filename that doesn't already exist.

retry_open (*filename*, *errno=[13]*, *retry_schedule=[2, 10, 30, 120, 300]*, ***kwargs*)

Try opening a file, using the builtin `open()` function. If an `IOError` is raised and its *errno* is in the given list, wait a moment then retry. Keeps doing this, waiting a bit longer each time, hoping that the problem will go away.

Once too many attempts have been made, outputs a message and waits for user input. This means the user can fix the problem (e.g. renew credentials) and pick up where execution left off. If they choose not to, the original error will be raised

Default list of *errno*s is just `[13]` – permission denied.

Use *retry_schedule* to customize the lengths of time waited between retries. Default: 2s, 10s, 30s, 2m, 5m, then give up.

Additional *kwargs* are pass on to `open()`.

extract_from_archive (*archive_filename*, *members*, *target_dir*, *preserve_dirs=True*)

Extract a file or files from an archive, which may be a tarball or a zip file (determined by the file extension).

extract_archive (*archive_filename*, *target_dir*, *preserve_dirs=True*)

Extract all files from an archive, which may be a tarball or a zip file (determined by the file extension).

pimlico.utils.format module

multiline_tablate (*table*, *widths*, ***kwargs*)

title_box (*title_text*)

Make a nice big pretty title surrounded by a box.

pimlico.utils.linguistic module

strip_punctuation (*s*, *split_words=True*)

pimlico.utils.logging module

get_console_logger (*name*, *debug=False*)

Convenience function to make it easier to create new loggers.

Parameters

- **name** – logging system logger name
- **debug** – whether to use DEBUG level. By default, uses INFO

Returns

pimlico.utils.network module

get_unused_local_port ()

Find a local port that's not currently being used, which we'll be able to bind a service to once this function returns.

get_unused_local_ports (*n*)

Find a number of local ports not currently in use. Binds each port found before looking for the next one. If you just called `get_unused_local_port()` multiple times, you'd get to same answer coming back.

pimlico.utils.pipes module

qget (*queue*, **args*, ***kwargs*)

Wrapper that calls the `get()` method of a queue, catching EINTR interrupts and retrying. Recent versions of Python have this built in, but with earlier versions you can end up having processes die while waiting on queue output because an EINTR has received (which isn't necessarily a problem).

Parameters

- **queue** –
- **args** – args to pass to queue's `get()`
- **kwargs** – kwargs to pass to queue's `get()`

Returns

class OutputQueue (*out*)

Bases: `object`

Direct a readable output (e.g. pipe from a subprocess) to a queue. Returns the queue. Output is added to the queue one line at a time. To perform a non-blocking read call `get_nowait()` or `get(timeout=T)`

get_nowait ()

get (*timeout=None*)

get_available ()

Don't block. Just return everything that's available in the queue.

pimlico.utils.pos module

pos_tag_to_ptb (*tag*)

see :doc:pos_pos_tags_to_ptb

pos_tags_to_ptb (*tags*)

Takes a list of POS tags and checks they're all in the PTB tagset. If they're not, tries mapping them according to CCGBank's special version of the tagset. If that doesn't work, raises a `NonPTBTagError`.

exception NonPTBTagError

Bases: `exceptions.Exception`

pimlico.utils.probability module

limited_shuffle (*iterable, buffer_size, rand_generator=None*)

Some algorithms require the order of data to be randomized. An obvious solution is to put it all in a list and shuffle, but if you don't want to load it all into memory that's not an option. This method iterates over the data, keeping a buffer and choosing at random from the buffer what to put next. It's less shuffled than the simpler solution, but limits the amount of memory used at any one time to the buffer size.

limited_shuffle_numpy (*iterable, buffer_size, randint_buffer_size=1000*)

Identical behaviour to `limited_shuffle()`, but uses Numpy's random sampling routines to generate a large number of random integers at once. This can make execution a bit bursty, but overall tends to speed things up, as we get the random sampling over in one big call to Numpy.

batched_randint (*low, high=None, batch_size=1000*)

Infinite iterable that produces random numbers in the given range by calling Numpy now and then to generate lots of random numbers at once and then yielding them one by one. Faster than sampling one at a time.

Parameters

- **a** – lowest number in range
- **b** – highest number in range
- **batch_size** – number of ints to generate in one go

sequential_document_sample (*corpus, start=None, shuffle=None, sample_rate=None*)

Wrapper around a `pimlico.datatypes.tar.TarredCorpus` to draw infinite samples of documents from the corpus, by iterating over the corpus (looping infinitely), yielding documents at random. If `sample_rate` is given, it should be a float between 0 and 1, specifying the rough proportion of documents to sample. A lower value spreads out the documents more on average.

Optionally, the samples are shuffled within a limited scope. Set `shuffle` to the size of this scope (higher will shuffle more, but need to buffer more samples in memory). Otherwise (`shuffle=0`), they will appear in the order they were in the original corpus.

If `start` is given, that number of documents will be skipped before drawing any samples. Set `start=0` to start at the beginning of the corpus. By default (`start=None`) a random point in the corpus will be skipped to before beginning.

sequential_sample (*iterable, start=0, shuffle=None, sample_rate=None*)

Draw infinite samples from an iterable, by iterating over it (looping infinitely), yielding items at random. If `sample_rate` is given, it should be a float between 0 and 1, specifying the rough proportion of documents to sample. A lower value spreads out the documents more on average.

Optionally, the samples are shuffled within a limited scope. Set `shuffle` to the size of this scope (higher will shuffle more, but need to buffer more samples in memory). Otherwise (`shuffle=0`), they will appear in the order they were in the original corpus.

If `start` is given, that number of documents will be skipped before drawing any samples. Set `start=0` to start at the beginning of the corpus. Note that setting this to a high number can result in a slow start-up, if iterating over the items is slow.

Note: If you're sampling documents from a *TarredCorpus*, it's better to use `sequential_document_sample()`, since it makes use of *TarredCorpus*'s built-in features to do the skipping and sampling more efficiently.

subsample (*iterable, sample_rate*)

Subsample the given iterable at a given rate, between 0 and 1.

pimlico.utils.progress module

get_progress_bar (*maxval, counter=False, title=None, start=True*)

Simple utility to build a standard progress bar, so I don't have to think about this each time I need one. Starts the progress bar immediately.

`start` is no longer used, included only for backwards compatibility.

class SafeProgressBar (*maxval=None, widgets=None, term_width=None, poll=1, left_justify=True, fd=None*)

Bases: `progressbar.progressbar.ProgressBar`

Override basic progress bar to wrap `update()` method with a couple of extra features.

1. You don't need to call `start()` – it will be called when the first update is received. This is good for processes that have a bit of a start-up lag, or where starting to iterate might generate some other output.
2. An error is not raised if you update with a value higher than `maxval`. It's the most annoying thing ever if you run a long process and the whole thing fails near the end because you slightly miscalculated `maxval`.

update (*value=None*)

Updates the `ProgressBar` to a new value.

increment ()

class DummyFileDescriptor

Bases: `object`

Passed in to `ProgressBar` instead of a file descriptor (e.g. `stderr`) to ensure that nothing gets output.

read (*size=None*)

readLine (*size=None*)

write (*s*)

close ()

class NonOutputtingProgressBar (**args, **kwargs*)

Bases: `pimlico.utils.progress.SafeProgressBar`

Behaves like `ProgressBar`, but doesn't output anything.

class LittleOutputtingProgressBar (**args, **kwargs*)

Bases: `pimlico.utils.progress.SafeProgressBar`

Behaves like `ProgressBar`, but doesn't output much. Instead of constantly redrawing the progress bar line, it outputs a simple progress message every time it hits the next 10% mark.

If running on a terminal, this will update the line, as with a normal progress bar. If piping to a file, this will just print a new line occasionally, so won't fill up your file with thousands of progress updates.

start ()

Starts measuring time, and prints the bar at 0%.

It returns self so you can use it like this: >>> pbar = ProgressBar().start() >>> for i in range(100): ... # do something ... pbar.update(i+1) ... >>> pbar.finish()

finish ()

Puts the ProgressBar bar in the finished state.

slice_progress (*iterable, num_items, title=None*)

class ProgressBarIter (*iterable, title=None*)

Bases: object

pimlico.utils.strings module

truncate (*s, length, ellipsis=u'...'*)

similarities (*targets, reference*)

Compute string similarity of each of a list of targets to a given reference string. Uses *difflib.SequenceMatcher* to compute similarity.

Parameters

- **reference** – compare all strings to this one
- **targets** – list of targets to measure similarity of

Returns list of similarity values

sorted_by_similarity (*targets, reference*)

Return target list sorted by similarity to the reference string. See :func:similarities for similarity measurement.

pimlico.utils.system module

Lowish-level system operations

set_proc_title (*title*)

Tries to set the current process title. This is very system-dependent and may not always work.

If it's available, we use the *setproctitle* package, which is the most reliable way to do this. If not, we try doing it by loading *libc* and calling *prctl* ourselves. This is not reliable and only works on Unix systems. If neither of these works, we give up and return False.

If you want to increase the chances of this working (e.g. your process titles don't seem to be getting set by Pimlico and you'd like them to), try installing *setproctitle*, either system-wide or in Pimlico's virtualenv.

@return: True if the process succeeds, False if there's an error

pimlico.utils.timeout module

timeout (*func, args=(), kwargs={}, timeout_duration=1, default=None*)

pimlico.utils.urwid module

Some handy Urwid utilities.

Take care only to import this where we already have a dependency on Urwid, e.g. in the browser implementation modules.

Some of these are taken pretty exactly from Urwid examples.

Todo: Not got these things working yet, but they'll be useful in the long run

exception DialogExit

Bases: `exceptions.Exception`

class DialogDisplay (*original_widget, text, height=0, width=0, body=None*)

Bases: `urwid.wimp.PopUpLauncher`

`palette = [('body', 'black', 'light gray', 'standout'), ('border', 'black', 'dark blue`

`add_buttons (buttons)`

`button_press (button)`

`on_exit (exitcode)`

class ListDialogDisplay (*original_widget, text, height, width, constr, items, has_default*)

Bases: `pimlico.utils.urwid.DialogDisplay`

`unhandled_key (size, k)`

`on_exit (exitcode)`

Print the tag of the item selected.

msgbox (*original_widget, text, height=0, width=0*)

options_dialog (*original_widget, text, options, height=0, width=0, *items*)

yesno_dialog (*original_widget, text, height=0, width=0, *items*)

pimlico.utils.web module

download_file (*url, target_file*)

Module contents

Submodules

pimlico.cfg module

Global config

Various global variables. Access as follows:

```
from pimlico import cfg
```

```
# Set global config parameter cfg.parameter = "Value" # Use parameter print cfg.parameter
```

There are some global variables in `pimlico` (in the `__init__.py`) that probably should be moved here, but I'm leaving them for now. At the moment, none of those are ever written from outside that file (i.e. think of them as constants, rather than config), so the only reason to move them is to keep everything in one place.

Module contents

The Pimlico Processing Toolkit (PIpelled Modular LInguistic COrpus processing) is a toolkit for building pipelines made up of linguistic processing tasks to run on large datasets (corpora). It provides a wrappers around many existing, widely used NLP (Natural Language Processing) tools.

`install_core_dependencies()`

1.6 Module test pipelines

Test pipelines provide a special sort of unit testing for Pimlico.

Pimlico is distributed with a set of test pipeline config files, each just a small pipeline with a couple of modules in it. Each is designed to test the use of a particular one of Pimlico's builtin module types, or some combination of a smaller number of them.

1.6.1 Available pipelines

normalize

This is one of the test pipelines included in Pimlico's repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=normalize
release=latest

# Take input from a prepared Pimlico dataset
[europarl]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=TokenizedDocumentType
dir=%(test_data_dir)s/datasets/corpora/tokenized

[norm]
type=pimlico.modules.text.normalize
case=lower
remove_empty=T
```

Modules

The following Pimlico module types are used in this pipeline:

- `normalize`

simple_tokenize

This is one of the test pipelines included in Pimlico’s repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=simple_tokenize
release=latest

# Take input from a prepared Pimlico dataset
[europarl]
type=pimlico.datatypes.corpora.GroupedCorpus
# This corpus is actually tokenized text, but we treat it as raw text and apply the_
↔simple tokenizer
data_point_type=RawTextDocumentType
dir=%(test_data_dir)s/datasets/corpora/tokenized

[tokenize]
type=pimlico.modules.text.simple_tokenize
```

Modules

The following Pimlico module types are used in this pipeline:

- *simple_tokenize*

raw_text_files_test

This is one of the test pipelines included in Pimlico’s repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=raw_text_files_test
release=latest

# Read in some Europarl raw files
[europarl]
type=pimlico.modules.input.text.raw_text_files
files=%(test_data_dir)s/datasets/europarl_en_raw/*
```

fasttext_input_test

This is one of the test pipelines included in Pimlico’s repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=fasttext_input_test
release=latest

# Read in some vectors
[vectors]
type=pimlico.modules.input.embeddings.fasttext
path=%(test_data_dir)s/input_data/fasttext/wiki.en_top50.vec
```

Modules

The following Pimlico module types are used in this pipeline:

- *fasttext*

glove_input_test

This is one of the test pipelines included in Pimlico's repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=glove_input_test
release=latest

# Read in some vectors
[vectors]
type=pimlico.modules.input.embeddings.glove
path=%(test_data_dir)s/input_data/glove/glove.small.300d.txt
```

Modules

The following Pimlico module types are used in this pipeline:

- *glove*

opennlp_tokenize

This is one of the test pipelines included in Pimlico's repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```

[pipeline]
name=opennlp_tokenize
release=latest

# Prepared tarred corpus
[europarl]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=RawTextDocumentType
dir=%(test_data_dir)s/datasets/text_corpora/europarl

# There's a problem with the tests here
# Pimlico still has a clunky old Makefile-based system for installing model data for_
↔modules
# The tests don't know that this needs to be done before the pipeline can be run
# This is why this test is not in the main suite, but a special OpenNLP one
[tokenize]
type=pimlico.modules.opennlp.tokenize
token_model=en-token.bin
sentence_model=en-sent.bin

```

Modules

The following Pimlico module types are used in this pipeline:

- `tokenize`

store

This is one of the test pipelines included in Pimlico's repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```

[pipeline]
name=store
release=latest

# Read in some Europarl raw files
[europarl]
type=pimlico.modules.input.text.raw_text_files
files=%(test_data_dir)s/datasets/europarl_en_raw/*
encoding=utf8

# Group works as a filter module, so its output is not stored.
# This pipeline shows how you can store the output from such a
# module for static use by later modules.
# In this exact case, you don't gain anything by doing that, since
# the grouping filter is fast, but sometimes it could be desirable
# with other filters
[group]
type=pimlico.modules.corpora.group

```

(continues on next page)

```
[store]
type=pimlico.modules.corpora.store
```

Modules

The following Pimlico module types are used in this pipeline:

- *group*
- *store*

concat

This is one of the test pipelines included in Pimlico's repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=concat
release=latest

# Take input from some prepared Pimlico datasets
[europarl1]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=RawTextDocumentType
dir=%(test_data_dir)s/datasets/text_corpora/europarl

[europarl2]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=RawTextDocumentType
dir=%(test_data_dir)s/datasets/text_corpora/europarl2

[concat]
type=pimlico.modules.corpora.concat
input_corpora=europarl1,europarl2

[output]
type=pimlico.modules.corpora.format
```

Modules

The following Pimlico module types are used in this pipeline:

- *concat*
- *format*

group

This is one of the test pipelines included in Pimlico’s repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=group
release=latest

# Read in some Europarl raw files
[europarl]
type=pimlico.modules.input.text.raw_text_files
files=%(test_data_dir)s/datasets/europarl_en_raw/*
encoding=utf8

[group]
type=pimlico.modules.corpora.group

[output]
type=pimlico.modules.corpora.format
```

Modules

The following Pimlico module types are used in this pipeline:

- *group*
- *format*

vocab_builder

This is one of the test pipelines included in Pimlico’s repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=vocab_builder
release=latest

# Take input from a prepared Pimlico dataset
[europarl]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=TokenizedDocumentType
dir=%(test_data_dir)s/datasets/corpora/tokenized

[vocab]
type=pimlico.modules.corpora.vocab_builder
```

(continues on next page)

(continued from previous page)

```
threshold=2
limit=500
```

Modules

The following Pimlico module types are used in this pipeline:

- *vocab_builder*

subset

This is one of the test pipelines included in Pimlico's repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=subset
release=latest

# Take input from a prepared Pimlico dataset
[europarl]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=RawTextDocumentType
dir=%(test_data_dir)s/datasets/text_corpora/europarl

[subset]
type=pimlico.modules.corpora.subset
size=1
offset=2

[output]
type=pimlico.modules.corpora.format
```

Modules

The following Pimlico module types are used in this pipeline:

- *subset*
- *format*

interleave

This is one of the test pipelines included in Pimlico's repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=interleave
release=latest

# Take input from some prepared Pimlico datasets
[europarl1]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=RawTextDocumentType
dir=%(test_data_dir)s/datasets/text_corpora/europarl

[europarl2]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=RawTextDocumentType
dir=%(test_data_dir)s/datasets/text_corpora/europarl2

[interleave]
type=pimlico.modules.corpora.interleave
input_corpora=europarl1,europarl2

[output]
type=pimlico.modules.corpora.format
```

Modules

The following Pimlico module types are used in this pipeline:

- *interleave*
- *format*

stats

This is one of the test pipelines included in Pimlico's repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=stats
release=latest

# Take input from a prepared Pimlico dataset
[europarl]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=TokenizedDocumentType
dir=%(test_data_dir)s/datasets/corpora/tokenized
```

(continues on next page)

```
[stats]
type=pimlico.modules.corpora.corpus_stats
```

Modules

The following Pimlico module types are used in this pipeline:

- *corpus_stats*

list_filter

This is one of the test pipelines included in Pimlico's repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=list_filter
release=latest

# Take input from a prepared Pimlico dataset
[europarl]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=RawTextDocumentType
dir=%(test_data_dir)s/datasets/text_corpora/europarl

[filename_list]
type=StringList
dir=%(test_data_dir)s/datasets/europarl_filename_list

# Use the filename list to filter the documents
# This should leave 3 documents (of original 5)
[europarl_filtered]
type=pimlico.modules.corpora.list_filter
input_corpus=europarl
input_list=filename_list
```

Modules

The following Pimlico module types are used in this pipeline:

- *list_filter*

split

This is one of the test pipelines included in Pimlico's repository. See *Module test pipelines* for more details.

Config file

The complete config file for this test pipeline:

```
[pipeline]
name=split
release=latest

# Take input from a prepared Pimlico dataset
[europarl]
type=pimlico.datatypes.corpora.GroupedCorpus
data_point_type=RawTextDocumentType
dir=%(test_data_dir)s/datasets/text_corpora/europarl

[split]
type=pimlico.modules.corpora.split
set1_size=2
```

Modules

The following Pimlico module types are used in this pipeline:

- *split*

1.6.2 Input data

Pimlico also comes with all the data necessary to run the pipelines. They all use very small datasets, so that they don't take long to run and can be easily distributed.

Some of the datasets are raw data, of the sort you might find in a distributed corpus, and these are used to test input readers for that type of data. Most, however, are stored in one of Pimlico's datatype formats, exactly as they were output from some other module (most often from another test pipeline), so that they can be read in to test one module in isolation.

1.6.3 Usage examples

In addition to providing unit testing for core Pimlico modules, test pipelines also function as a source of examples of each module's usage. They are for that reason linked to from the module's documentation, so that example usages can be easily found where available.

1.6.4 Running

To run test pipelines, you can use the script `test_pipeline.sh` in Pimlico's bin directory, e.g.:

```
./test_pipeline.sh ../test/data/pipelines/corpora/concat.conf output
```

This will load a single test pipeline from the given config file and execute the module named `output`.

There are also some suites of tests, specified as CSV files giving a number of config files and module names to execute for each. To run the main suite of test pipelines for Pimlico's core modules, run:

```
./all_test_pipelines.sh
```

1.7 Future plans

Various things I plan to add to Pimlico in the futures. For a summary, see *Pimlico Wishlist*.

1.7.1 Pimlico Wishlist

Things I plan to add to Pimlico.

- Further modules:
 - *CherryPicker* for coreference resolution
 - *Berkeley Parser* for fast constituency parsing
 - *Reconcile* coref. Seems to incorporate upstream NLP tasks. Would want to interface such that we can reuse output from other modules and just do coref.
- **Pipeline graph visualizations:** *Outputting pipeline diagrams*. Maybe an interactive GUI to help with viewing large pipelines
- See [issue list on Github](#) for other specific plans
- Big redesign of datatype implementation is [documented as a Github project](#)

Todos

The following to-dos appear elsewhere in the docs. They are generally bits of the documentation I've not written yet, but am aware are needed.

Todo: Continue updating this for the new datatype system. I've got partway, but the reader is still far from finished

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/src/python/pimlico/core/modules.py` of `pimlico.core.modules.map.filter`, line 1.)

Todo: Under the new datatype system, this should be done differently. Don't wrap datatypes, but instead use the actual output datatypes (taken from the wrapped module type's output) and instead create custom readers that gets instantiated when fetching the module's output readers.

I've created the test pipeline filter_tokenize for testing this.

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/src/python/pimlico/core/modules.py` of `pimlico.core.modules.map.filter.wrap_module_info_as_filter`, line 7.)

Todo: Call `get_writer_software_dependencies` before instantiating writer

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/src/python/pimlico/datatypes.py` of `pimlico.datatypes.base.PimlicoDatatype.get_writer_software_dependencies`, line 10.)

Todo: Call `get_writer_software_dependencies` before instantiating writer

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/src/python/pimlico/datatypes.py` of `pimlico.datatypes.embeddings.Embeddings.get_writer_software_dependencies`, line 10.)

Todo: Add unit test for ScoredReadFeatureSets

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/src/python/pimlico/datatypes/features/ScoredRealFeatureSets, line 9.)

Todo: Not got these things working yet, but they'll be useful in the long run

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/src/python/pimlico/utils/urwid of pimlico.utils.urwid, line 8.)

Todo: Describe how module dependencies are defined for different types of deps

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/core/dependencies.rst, line 73.)

Todo: Include some examples from the core modules of how deps are defined and some special cases of software fetching

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/core/dependencies.rst, line 80.)

Todo: Write documentation for this

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/core/module_structure.rst, line 9.)

Todo: Filter module guide needs to be updated for new datatypes. This section is currently completely wrong – **ignore it!** This is quite a substantial change.

The difficulty of describing what you need to do here suggests we might want to provide some utilities to make this easier!

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/guides/filters.rst, line 31.)

Todo: Write a guide to building document map modules.

For now, the skeletons below are a useful starting point, but there should be a more fulsome explanation here of what document map modules are all about and how to use them.

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/guides/map_module.rst, line 5.)

Todo: Document map module guides needs to be updated for new datatypes.

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/guides/map_module.rst`, line 12.)

Todo: Module writing guide needs to be updated for new datatypes.

In particular, the executor example and datatypes in the module definition need to be updated.

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/guides/module.rst`, line 23.)

Todo: Setup guide has a lot that needs to be updated for the new datatypes system. I've updated up to **Getting input**.

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/guides/setup.rst`, line 5.)

Todo: Continue writing from here

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/guides/setup.rst`, line 110.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.module.rst`, line 25.)

Todo: Update to new datatypes system and add test pipelines

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.module.rst`, line 36.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.module.rst`, line 30.)

Todo: Write description of vocab mapper module

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.module.rst`, line 12.)

Todo: Add test pipeline and test

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.module.rst`, line 16.)

Todo: Document this module

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 16.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 20.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 22.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 23.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 23.)

Todo: Document this module

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 16.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 20.)

Todo: Document this module

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 16.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 20.)

Todo: Document this module

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 16.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 20.)

Todo: Document this module

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 16.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 20.)

Todo: Add test pipeline and test

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 15.)

Todo: Add test pipeline and test

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 19.)

Todo: Add test pipeline. This is slightly difficult, as we need a small FastText binary file, which is harder to produce, since you can't easily just truncate a big file.

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 27.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 20.)

Todo: Document this module

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 16.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 20.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 19.)

Todo: Document this module

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 16.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 20.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 19.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 26.)

Todo: Document this module

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 16.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 20.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod` line 22.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod` line 19.)

Todo: Document this module

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod` line 16.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod` line 20.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod` line 25.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod` line 21.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod` line 18.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod` line 26.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in `/home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod` line 51.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 24.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 21.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 18.)

Todo: Update to new datatypes system and add test pipeline

(The [original entry](#) is located in /home/docs/checkouts/readthedocs.org/user_builds/pimlico/checkouts/latest/docs/modules/pimlico.mod line 24.)

1.7.2 Berkeley Parser

<https://github.com/slavpetrov/berkeleyparser>

Java constituency parser. Pre-trained models are also provided in the Github repo.

Probably no need for a Java wrapper here. The parser itself accepts input on stdin and outputs to stdout, so just use a subprocess with pipes.

1.7.3 Cherry Picker

Coreference resolver

<http://www.hlt.utdallas.edu/~altaf/cherrypicker/>

Requires NER, POS tagging and constituency parsing to be done first. Tools for all of these are included in the Cherry Picker codebase, but we just need a wrapper around the Cherry Picker tool itself to be able to feed these annotations in from other modules and perform coref.

Write a Java wrapper and interface with it using Py4J, as with OpenNLP.

1.7.4 Outputting pipeline diagrams

Once pipeline config files get big, it can be difficult to follow what's going on in them, especially if the structure is more complex than just a linear pipeline. A useful feature would be the ability to display/output a visualization of the pipeline as a flow graph.

It looks like the easiest way to do this will be to construct a DOT graph using Graphviz/Pydot and then output the diagram using Graphviz.

<http://www.graphviz.org>

<https://pypi.python.org/pypi/pydot>

Building the graph should be pretty straightforward, since the mapping from modules to nodes is fairly direct.

We could also add extra information to the nodes, like current execution status.

- genindex
- search

C

- `pimlico.cfg`, 187
- `pimlico.cli`, 123
 - `browser`, 114
 - `tool`, 114
 - `tools`, 114
 - `corpus`, 111
 - `files`, 112
 - `formatter`, 112
 - `check`, 117
 - `clean`, 117
 - `debug`, 115
 - `stepper`, 114
 - `loaddump`, 118
 - `locations`, 118
 - `main`, 119
 - `newmodule`, 120
 - `pyshell`, 120
 - `reset`, 121
 - `run`, 121
 - `shell`, 117
 - `base`, 115
 - `commands`, 116
 - `runner`, 116
 - `status`, 121
 - `subcommands`, 122
 - `testemail`, 122
 - `util`, 122
- `pimlico.core`, 157
 - `config`, 153
 - `dependencies`, 129
 - `base`, 123
 - `core`, 125
 - `java`, 125
 - `python`, 127
 - `versions`, 129
 - `external`, 130
 - `java`, 129
 - `logs`, 157

- `pimlico.core.modules`, 152
 - `base`, 138
 - `execute`, 146
 - `inputs`, 147
 - `map`, 135
 - `filter`, 130
 - `map`, 132
 - `singleproc`, 133
 - `threaded`, 134
 - `multistage`, 149
 - `options`, 152
 - `paths`, 157

d

- `pimlico.datatypes`, 177
 - `arrays`, 167
 - `base`, 168
 - `core`, 172
 - `corpora`, 167
 - `base`, 157
 - `data_points`, 158
 - `floats`, 161
 - `grouped`, 163
 - `ints`, 164
 - `table`, 165
 - `tokenized`, 166
 - `dictionary`, 172
 - `embeddings`, 173
 - `features`, 174
 - `files`, 175
 - `gensim`, 176
 - `sklearn`, 177

m

- `pimlico.modules`, 40
 - `candc`, 40
 - `corenlp`, 41
 - `corpora`, 43
 - `concat`, 43

pimlico.modules.corpora.corpus_stats, 44
 pimlico.modules.corpora.format, 45
 pimlico.modules.corpora.group, 46
 pimlico.modules.corpora.interleave, 47
 pimlico.modules.corpora.list_filter, 49
 pimlico.modules.corpora.split, 49
 pimlico.modules.corpora.store, 51
 pimlico.modules.corpora.subset, 51
 pimlico.modules.corpora.vocab_builder, 53
 pimlico.modules.corpora.vocab_counter, 54
 pimlico.modules.corpora.vocab_mapper, 55
 pimlico.modules.embeddings, 56
 pimlico.modules.embeddings.dependencies, 56
 pimlico.modules.embeddings.store_embeddings, 58
 pimlico.modules.embeddings.store_tsv, 58
 pimlico.modules.embeddings.store_word2vec, 59
 pimlico.modules.embeddings.word2vec, 60
 pimlico.modules.features, 61
 pimlico.modules.features.term_feature_compiler, 61
 pimlico.modules.features.term_feature_mapper, 63
 pimlico.modules.features.vocab_builder, 63
 pimlico.modules.features.vocab_mapper, 65
 pimlico.modules.gensim, 65
 pimlico.modules.gensim.lda, 66
 pimlico.modules.gensim.lda_doc_topics, 68
 pimlico.modules.input, 69
 pimlico.modules.input.embeddings, 69
 pimlico.modules.input.embeddings.fasttext, 69
 pimlico.modules.input.embeddings.fasttext.glove, 70
 pimlico.modules.input.embeddings.glove, 71
 pimlico.modules.input.embeddings.word2vec, 72
 pimlico.modules.input.text, 73
 pimlico.modules.input.text.raw_text_file, 73
 pimlico.modules.input.text.annotations, 74
 pimlico.modules.malt, 75
 pimlico.modules.malt.conll_parser_input, 75
 pimlico.modules.malt.parse, 75
 pimlico.modules.nltk, 77
 pimlico.modules.nltk.nist_tokenize, 77
 pimlico.modules.opennlp, 78
 pimlico.modules.opennlp.coreference, 78
 pimlico.modules.opennlp.coreference_pipeline, 79
 pimlico.modules.opennlp.ner, 81
 pimlico.modules.opennlp.parse, 82
 pimlico.modules.opennlp.pos, 83
 pimlico.modules.opennlp.tokenize, 84
 pimlico.modules.r, 85
 pimlico.modules.r.script, 85
 pimlico.modules.regex, 86
 pimlico.modules.regex.annotated_text, 86
 pimlico.modules.sklearn, 87
 pimlico.modules.sklearn.logistic_regression, 88
 pimlico.modules.sklearn.matrix_factorization, 89
 pimlico.modules.text, 90
 pimlico.modules.text.char_tokenize, 90
 pimlico.modules.text.normalize, 91
 pimlico.modules.text.simple_tokenize, 92
 pimlico.modules.text.text_normalize, 93
 pimlico.modules.text.untokenize, 94
 pimlico.modules.utility, 95
 pimlico.modules.utility.alias, 95
 pimlico.modules.utility.collect_files, 97
 pimlico.modules.utility.copy_file, 98
 pimlico.modules.visualization, 99
 pimlico.modules.visualization.bar_chart, 99
 pimlico.modules.visualization.embeddings_plot, 100

p
 pimlico.pipeline, 88

t
 pimlico.test, 179
 pimlico.test.pipeline, 178
 pimlico.test.suite, 179

u
 pimlico.utils, 187
 pimlico.utils.communicate, 179
 pimlico.utils.core, 180
 pimlico.utils.docs, 179

pimlico.utils.docs.rest, 179
pimlico.utils.email, 181
pimlico.utils.filesystem, 182
pimlico.utils.format, 182
pimlico.utils.linguistic, 182
pimlico.utils.logging, 183
pimlico.utils.network, 183
pimlico.utils.pipes, 183
pimlico.utils.pos, 183
pimlico.utils.probability, 184
pimlico.utils.progress, 185
pimlico.utils.strings, 186
pimlico.utils.system, 186
pimlico.utils.timeout, 186
pimlico.utils.urwid, 187
pimlico.utils.web, 187

A

- abs_path_or_model_dir_path() (in module pimlico.core.paths), 157
- add_arguments() (BrowseCmd method), 120
- add_arguments() (DepsCmd method), 117
- add_arguments() (DumpCmd method), 118
- add_arguments() (InputsCmd method), 118
- add_arguments() (InstallCmd method), 117
- add_arguments() (LoadCmd method), 118
- add_arguments() (MoveStoresCmd method), 119
- add_arguments() (OutputCmd method), 119
- add_arguments() (PimlicoCLISubcommand method), 122
- add_arguments() (PythonShellCmd method), 121
- add_arguments() (ResetCmd method), 121
- add_arguments() (RunCmd method), 121
- add_arguments() (ShellCLICmd method), 116
- add_arguments() (StatusCmd method), 121
- add_arguments() (UnlockCmd method), 120
- add_arguments() (VariantsCmd method), 119
- add_arguments() (VisualizeCmd method), 120
- add_buttons() (DialogDisplay method), 187
- add_execution_history_record() (BaseModuleInfo method), 139
- AlignedGroupedCorpora (class in pimlico.datatypes.corpora.grouped), 163
- all_dependencies() (SoftwareDependency method), 124
- all_inputs_ready() (BaseModuleInfo method), 142
- append_module() (PipelineConfig method), 153
- archive_iter() (AlignedGroupedCorpora method), 163
- archive_iter() (DocumentMapOutputTypeWrapper method), 131
- archive_iter() (FilterModuleOutputReader method), 131
- archive_iter_decorator() (in module pimlico.cli.debug.stepper), 114
- ask() (in module pimlico.cli.newmodule), 120
- available() (SoftwareDependency method), 123

B

- BaseModuleExecutor (class in pimlico.core.modules.base), 145
- BaseModuleInfo (class in pimlico.core.modules.base), 138
- batched_randint() (in module pimlico.utils.probability), 184
- BeautifulSoupDependency (class in pimlico.core.dependencies.python), 128
- browse_cmd() (in module pimlico.cli.browser.tool), 114
- browse_data() (in module pimlico.cli.browser.tools.corpus), 111
- browse_file() (NamedFileCollection method), 175
- browse_file() (ScoredRealFeatureSets method), 174
- browse_files() (in module pimlico.cli.browser.tools.files), 112
- BrowseCmd (class in pimlico.cli.main), 120
- button_press() (DialogDisplay method), 187

C

- cached_property (class in pimlico.utils.core), 181
- call_java() (in module pimlico.core.external.java), 129
- CharacterTokenizedDocumentType (class in pimlico.datatypes.corpora.tokenized), 166
- check_and_execute_modules() (in module pimlico.core.modules.execute), 146
- check_and_install() (in module pimlico.core.dependencies.base), 124
- check_for_cycles() (in module pimlico.core.config), 156
- check_for_error() (InputQueueFeeder method), 137
- check_invalid() (InputQueueFeeder method), 136
- check_java() (in module pimlico.core.dependencies.java), 126
- check_java_dependency() (in module pimlico.core.dependencies.java), 126
- check_modules_ready() (in module pimlico.core.modules.execute), 146
- check_pipeline() (in module pimlico.core.config), 156
- check_ready_to_run() (BaseModuleInfo method), 144
- check_ready_to_run() (MultistageModuleInfo method), 150
- check_release() (in module pimlico.core.config), 156

- check_type() (DynamicInputDatatypeRequirement method), 171
- check_type() (FilesInput method), 176
- check_type() (in module pimlico.core.modules.base), 145
- check_type() (IterableCorpus method), 158
- check_type() (NamedFileCollection method), 175
- check_type() (PimlicoDatatype method), 170
- choose_from_list() (in module pimlico.core.modules.options), 152
- chunk_list() (in module pimlico.utils.core), 181
- CleanCmd (class in pimlico.cli.clean), 117
- clear_output_queues() (Py4JInterface method), 129
- clear_storage_dir() (in module pimlico.test.pipeline), 178
- close() (DummyFileDescriptor method), 185
- cmdloop() (DataShell method), 115
- collect_runnable_modules() (in module pimlico.core.modules.base), 145
- collect_unexecuted_dependencies() (in module pimlico.core.modules.base), 145
- comma_separated_list() (in module pimlico.core.modules.options), 152
- comma_separated_strings() (in module pimlico.core.modules.options), 152
- command_desc (CleanCmd attribute), 117
- command_desc (DumpCmd attribute), 118
- command_desc (InputsCmd attribute), 118
- command_desc (ListStoresCmd attribute), 119
- command_desc (LoadCmd attribute), 118
- command_desc (MoveStoresCmd attribute), 119
- command_desc (NewModuleCmd attribute), 120
- command_desc (PimlicoCLISubcommand attribute), 122
- command_desc (UnlockCmd attribute), 120
- command_desc (VisualizeCmd attribute), 120
- command_help (BrowseCmd attribute), 120
- command_help (CleanCmd attribute), 117
- command_help (DepsCmd attribute), 117
- command_help (DumpCmd attribute), 118
- command_help (EmailCmd attribute), 122
- command_help (InputsCmd attribute), 118
- command_help (InstallCmd attribute), 117
- command_help (ListStoresCmd attribute), 119
- command_help (LoadCmd attribute), 118
- command_help (MoveStoresCmd attribute), 119
- command_help (NewModuleCmd attribute), 120
- command_help (OutputCmd attribute), 119
- command_help (PimlicoCLISubcommand attribute), 122
- command_help (PythonShellCmd attribute), 121
- command_help (ResetCmd attribute), 121
- command_help (RunCmd attribute), 121
- command_help (ShellCLICmd attribute), 116
- command_help (StatusCmd attribute), 121
- command_help (UnlockCmd attribute), 120
- command_help (VariantsCmd attribute), 119
- command_help (VisualizeCmd attribute), 120
- command_name (BrowseCmd attribute), 120
- command_name (CleanCmd attribute), 117
- command_name (DepsCmd attribute), 117
- command_name (DumpCmd attribute), 118
- command_name (EmailCmd attribute), 122
- command_name (InputsCmd attribute), 118
- command_name (InstallCmd attribute), 117
- command_name (ListStoresCmd attribute), 119
- command_name (LoadCmd attribute), 118
- command_name (MoveStoresCmd attribute), 119
- command_name (NewModuleCmd attribute), 120
- command_name (OutputCmd attribute), 119
- command_name (PimlicoCLISubcommand attribute), 122
- command_name (PythonShellCmd attribute), 121
- command_name (ResetCmd attribute), 121
- command_name (RunCmd attribute), 121
- command_name (ShellCLICmd attribute), 116
- command_name (StatusCmd attribute), 121
- command_name (UnlockCmd attribute), 120
- command_name (VariantsCmd attribute), 119
- command_name (VisualizeCmd attribute), 120
- commands (CountInvalidCmd attribute), 157
- commands (MetadataCmd attribute), 116
- commands (PythonCmd attribute), 116
- commands (ShellCommand attribute), 115
- compare_dotted_versions() (in module pimlico.core.dependencies.versions), 129
- copy_dir_with_progress() (in module pimlico.utils.filesystem), 182
- CORE_PIMLICO_DEPENDENCIES (in module pimlico.core.dependencies.core), 125
- CorpusAlignmentError, 164
- CorpusState (class in pimlico.cli.browser.tools.corpus), 111
- CorpusWithTypeFromInput (class in pimlico.datatypes.corpora.grouped), 164
- CountInvalidCmd (class in pimlico.datatypes.corpora.base), 157
- create_pool() (DocumentMapModuleExecutor method), 136
- create_pool() (MultiprocessingMapModuleExecutor method), 132
- create_pool() (SingleThreadMapModuleExecutor method), 133
- create_pool() (ThreadingMapModuleExecutor method), 134
- create_pop_up() (InputPopupLauncher method), 112
- create_pop_up() (MessagePopupLauncher method), 112
- create_queue() (DocumentProcessorPool static method), 137
- create_queue() (MultiprocessingMapPool static method), 132
- create_queue() (ThreadingMapPool static method), 134

D

- data_point_type_opt() (in module pimlico.datatypes.corpora.base), 157
- data_ready() (DocumentMapOutputTypeWrapper method), 131
- DataPointError, 161
- DataPointType (class in pimlico.datatypes.corpora.data_points), 158
- DataPointType.Document (class in pimlico.datatypes.corpora.data_points), 159
- DataShell (class in pimlico.cli.shell.base), 115
- DATATYPE (DefaultFormatter attribute), 113
- DATATYPE (DocumentBrowserFormatter attribute), 113
- DATATYPE (FloatListsFormatter attribute), 162
- DATATYPE (VectorFormatter attribute), 163
- datatype_doc_info (DynamicInputDatatypeRequirement attribute), 171
- datatype_doc_info (FilesInput attribute), 176
- datatype_full_class_name() (pimlico.datatypes.base.PimlicoDatatype class method), 170
- datatype_name (CorpusWithTypeFromInput attribute), 164
- datatype_name (Dict attribute), 172
- datatype_name (Dictionary attribute), 172
- datatype_name (DynamicOutputDatatype attribute), 171
- datatype_name (Embeddings attribute), 173
- datatype_name (GensimLdaModel attribute), 176
- datatype_name (GroupedCorpus attribute), 163
- datatype_name (GroupedCorpusWithTypeFromInput attribute), 163
- datatype_name (IterableCorpus attribute), 158
- datatype_name (NamedFile attribute), 176
- datatype_name (NamedFileCollection attribute), 175
- datatype_name (NumpyArray attribute), 167
- datatype_name (PimlicoDatatype attribute), 169
- datatype_name (ScipySparseMatrix attribute), 167
- datatype_name (ScoredRealFeatureSets attribute), 174
- datatype_name (SklernModel attribute), 177
- datatype_name (StringList attribute), 172
- datatype_name (TextFile attribute), 176
- datatype_name (TSVVecFiles attribute), 174
- datatype_options (IterableCorpus attribute), 158
- datatype_options (NamedFile attribute), 176
- datatype_options (NamedFileCollection attribute), 175
- datatype_options (PimlicoDatatype attribute), 169
- datatype_options (TextFile attribute), 176
- DatatypeLoadError, 172
- DatatypeWriteError, 172
- default() (DataShell method), 115
- DefaultFormatter (class in pimlico.cli.browser.tools.formatter), 113
- dependencies (BaseModuleInfo attribute), 143
- dependencies() (NLTKResource method), 128
- dependencies() (Py4JSoftwareDependency method), 126
- dependencies() (SoftwareDependency method), 123
- DependencyCheckerError, 130
- DependencyError, 146
- DepsCmd (class in pimlico.cli.check), 117
- DialogDisplay (class in pimlico.utils.urwid), 187
- DialogExit, 187
- Dict (class in pimlico.datatypes.core), 172
- Dictionary (class in pimlico.datatypes.dictionary), 172
- dirsize() (in module pimlico.utils.filesystem), 182
- do_EOF() (DataShell method), 115
- Document (CharacterTokenizedDocumentType attribute), 166
- Document (FloatListDocumentType attribute), 162
- Document (FloatListsDocumentType attribute), 162
- Document (IntegerListDocumentType attribute), 165
- Document (IntegerListsDocumentType attribute), 165
- Document (IntegerTableDocumentType attribute), 166
- Document (InvalidDocument attribute), 160
- Document (RawDocumentType attribute), 161
- Document (RawTextDocumentType attribute), 161
- Document (SegmentedLinesDocumentType attribute), 166
- Document (TextDocumentType attribute), 161
- Document (TokenizedDocumentType attribute), 166
- Document (VectorDocumentType attribute), 163
- document() (DocumentMapModuleInfo method), 135
- document_preprocessors (GroupedCorpus attribute), 163
- DocumentBrowserFormatter (class in pimlico.cli.browser.tools.formatter), 113
- DocumentMapModuleExecutor (class in pimlico.core.modules.map), 135
- DocumentMapModuleInfo (class in pimlico.core.modules.map), 135
- DocumentMapOutputTypeWrapper (class in pimlico.core.modules.map.filter), 130
- DocumentMapProcessMixin (class in pimlico.core.modules.map), 137
- DocumentProcessorPool (class in pimlico.core.modules.map), 137
- download_file() (in module pimlico.utils.web), 187
- DummyFileDescriptor (class in pimlico.utils.progress), 185
- DumpCmd (class in pimlico.cli.loaddump), 118
- DynamicInputDatatypeRequirement (class in pimlico.datatypes.base), 171
- DynamicOutputDatatype (class in pimlico.datatypes.base), 171

E

- EmailCmd (class in pimlico.cli.testemail), 122
- EmailConfig (class in pimlico.utils.email), 181
- EmailError, 181

- Embeddings (class in pimlico.datatypes.embeddings), 173
 - empty() (PipelineConfig static method), 154
 - empty_all_queues() (DocumentProcessorPool method), 137
 - empty_all_queues() (MultiprocessingMapPool method), 132
 - emptyline() (DataShell method), 115
 - enable_step() (PipelineConfig method), 155
 - enable_step_for_pipeline() (in module pimlico.cli.debug.stepper), 114
 - encode() (StreamCommunicationPacket method), 180
 - execute() (BaseModuleExecutor method), 145
 - execute() (CountInvalidCmd method), 157
 - execute() (DocumentMapModuleExecutor method), 136
 - execute() (MetadataCmd method), 116
 - execute() (PythonCmd method), 116
 - execute() (ShellCommand method), 115
 - execute_modules() (in module pimlico.core.modules.execute), 147
 - execution_history (BaseModuleInfo attribute), 139
 - execution_history_path (BaseModuleInfo attribute), 139
 - extract_archive() (in module pimlico.utils.filesystem), 182
 - extract_file() (FilterModuleOutputReader method), 131
 - extract_from_archive() (in module pimlico.utils.filesystem), 182
 - extract_input_options() (pimlico.core.modules.base.BaseModuleInfo class method), 139
- ## F
- FileInput (in module pimlico.datatypes.files), 176
 - FilesInput (class in pimlico.datatypes.files), 176
 - filter_document() (DocumentBrowserFormatter method), 113
 - filter_document() (InvalidDocumentFormatter method), 113
 - FilterModuleOutputReader (class in pimlico.core.modules.map.filter), 131
 - find_data() (PipelineConfig method), 155
 - find_data_path() (PipelineConfig method), 154
 - find_data_store() (PipelineConfig method), 155
 - finish() (LittleOutputtingProgressBar method), 186
 - FloatListDocumentType (class in pimlico.datatypes.corpora.floats), 162
 - FloatListsDocumentType (class in pimlico.datatypes.corpora.floats), 161
 - FloatListsFormatter (class in pimlico.datatypes.corpora.floats), 162
 - fmt_frame_info() (in module pimlico.cli.debug), 115
 - format_document() (DefaultFormatter method), 113
 - format_document() (DocumentBrowserFormatter method), 113
 - format_document() (FloatListsFormatter method), 162
 - format_document() (InvalidDocumentFormatter method), 113
 - format_document() (VectorFormatter method), 163
 - format_execution_dependency_tree() (in module pimlico.core.modules.execute), 147
 - format_execution_error() (in module pimlico.cli.util), 123
 - format_file_size() (in module pimlico.utils.filesystem), 182
 - format_option_type() (in module pimlico.core.modules.options), 152
 - formatters (DataPointType attribute), 159
 - formatters (VectorDocumentType attribute), 162
 - from_local_config() (pimlico.utils.email.EmailConfig class method), 181
 - full_class_name() (pimlico.datatypes.corpora.data_points.DataPointType class method), 159
 - full_datatype_name() (IterableCorpus method), 158
 - full_datatype_name() (PimlicoDatatype method), 170
- ## G
- gateway_client_to_running_server() (in module pimlico.core.external.java), 129
 - GensimLdaModel (class in pimlico.datatypes.gensim), 176
 - get() (OutputQueue method), 183
 - get_absolute_output_dir() (BaseModuleInfo method), 140
 - get_all_executed_modules() (BaseModuleInfo method), 144
 - get_available() (OutputQueue method), 183
 - get_base_datatype_class() (DynamicOutputDatatype method), 171
 - get_base_datatype_class() (GroupedCorpusWithTypeFromInput method), 163
 - get_classpath() (in module pimlico.core.dependencies.java), 126
 - get_classpath_components() (JavaDependency method), 125
 - get_console_logger() (in module pimlico.utils.logging), 183
 - get_data_search_paths() (PipelineConfig method), 155
 - get_datatype() (CorpusWithTypeFromInput method), 164
 - get_datatype() (DynamicOutputDatatype method), 171
 - get_datatype() (GroupedCorpusWithTypeFromInput method), 164
 - get_dependencies() (in module pimlico.core.config), 156
 - get_dependent_modules() (PipelineConfig method), 153
 - get_detailed_status() (BaseModuleInfo method), 144
 - get_detailed_status() (DocumentMapModuleInfo method), 135

- [get_detailed_status\(\)](#) (MultistageModuleInfo method), 150
[get_execution_dependency_tree\(\)](#) (BaseModuleInfo method), 144
[get_extra_outputs_from_options\(\)](#) (BaseModuleInfo static method), 140
[get_input\(\)](#) (BaseModuleInfo method), 142
[get_input_datatype\(\)](#) (BaseModuleInfo method), 142
[get_input_decorator\(\)](#) (in module pimlico.cli.debug.stepper), 114
[get_input_module_connection\(\)](#) (BaseModuleInfo method), 141
[get_input_reader_setup\(\)](#) (BaseModuleInfo method), 142
[get_input_software_dependencies\(\)](#) (BaseModuleInfo method), 144
[get_input_software_dependencies\(\)](#) (MultistageModuleInfo method), 150
[get_installed_version\(\)](#) (PythonPackageDependency method), 127
[get_installed_version\(\)](#) (PythonPackageOnPip method), 128
[get_installed_version\(\)](#) (SoftwareDependency method), 124
[get_key_info_table\(\)](#) (pimlico.core.modules.base.BaseModuleInfo class method), 139
[get_key_info_table\(\)](#) (pimlico.core.modules.multistage.MultistageModuleInfo class method), 150
[get_log_file\(\)](#) (in module pimlico.core.logs), 157
[get_metadata\(\)](#) (BaseModuleInfo method), 139
[get_module_classpath\(\)](#) (in module pimlico.core.dependencies.java), 126
[get_module_output_dir\(\)](#) (BaseModuleInfo method), 140
[get_module_schedule\(\)](#) (PipelineConfig method), 153
[get_names\(\)](#) (DataShell method), 115
[get_new_log_filename\(\)](#) (BaseModuleInfo method), 145
[get_next_output_document\(\)](#) (InputQueueFeeder method), 136
[get_next_stage\(\)](#) (MultistageModuleInfo method), 150
[get_nowait\(\)](#) (OutputQueue method), 183
[get_output\(\)](#) (BaseModuleInfo method), 141
[get_output_datatype\(\)](#) (BaseModuleInfo method), 140
[get_output_dir\(\)](#) (BaseModuleInfo method), 140
[get_output_reader_setup\(\)](#) (BaseModuleInfo method), 141
[get_output_software_dependencies\(\)](#) (BaseModuleInfo method), 144
[get_output_writer\(\)](#) (BaseModuleInfo method), 141
[get_pipeline\(\)](#) (in module pimlico.cli.pyshell), 121
[get_pop_up_parameters\(\)](#) (InputPopupLauncher method), 112
[get_pop_up_parameters\(\)](#) (MessagePopupLauncher method), 112
[get_progress_bar\(\)](#) (in module pimlico.utils.progress), 185
[get_redirect_func\(\)](#) (in module pimlico.core.external.java), 129
[get_software_dependencies\(\)](#) (BaseModuleInfo method), 143
[get_software_dependencies\(\)](#) (Embeddings method), 173
[get_software_dependencies\(\)](#) (GensimLdaModel method), 176
[get_software_dependencies\(\)](#) (MultistageModuleInfo method), 149
[get_software_dependencies\(\)](#) (NumpyArray method), 167
[get_software_dependencies\(\)](#) (PimlicoDatatype method), 169
[get_software_dependencies\(\)](#) (ScipySparseMatrix method), 167
[get_software_dependencies\(\)](#) (SklearnModel method), 177
[get_struct\(\)](#) (in module pimlico.datatypes.corpora.table), 165
[get_transitive_dependencies\(\)](#) (BaseModuleInfo method), 143
[get_uninstalled_dependencies\(\)](#) (TestPipeline method), 178
[get_unused_local_port\(\)](#) (in module pimlico.utils.network), 183
[get_unused_local_ports\(\)](#) (in module pimlico.utils.network), 183
[get_writer\(\)](#) (PimlicoDatatype method), 170
[get_writer_software_dependencies\(\)](#) (Embeddings method), 173
[get_writer_software_dependencies\(\)](#) (PimlicoDatatype method), 169
[get_writers\(\)](#) (DocumentMapModuleInfo method), 135
[GroupedCorpus](#) (class in pimlico.datatypes.corpora.grouped), 163
[GroupedCorpusIterationError](#), 164
[GroupedCorpusWithTypeFromInput](#) (class in pimlico.datatypes.corpora.grouped), 163
- ## H
- [help_text](#) (CountInvalidCmd attribute), 157
[help_text](#) (MetadataCmd attribute), 116
[help_text](#) (PythonCmd attribute), 116
[help_text](#) (ShellCommand attribute), 115
- ## I
- [import_member\(\)](#) (in module pimlico.utils.core), 180
[import_package\(\)](#) (BeautifulSoupDependency method), 128
[import_package\(\)](#) (PythonPackageDependency method), 127
[increment\(\)](#) (SafeProgressBar method), 185

- infinite_cycle() (in module pimlico.utils.core), 180
- input_corpora (DocumentMapModuleInfo attribute), 135
- input_module_factory() (in module pimlico.core.modules.inputs), 147
- input_names (BaseModuleInfo attribute), 139
- input_ready() (BaseModuleInfo method), 142
- InputDialog (class in pimlico.cli.browser.tools.corpus), 111
- InputModuleInfo (class in pimlico.core.modules.inputs), 147
- InputPopupLauncher (class in pimlico.cli.browser.tools.corpus), 112
- InputQueueFeeder (class in pimlico.core.modules.map), 136
- InputsCmd (class in pimlico.cli.locations), 118
- install() (in module pimlico.core.dependencies.base), 124
- install() (JavaJarsDependency method), 126
- install() (NLTKResource method), 128
- install() (Py4JSoftwareDependency method), 126
- install() (PythonPackageOnPip method), 128
- install() (SoftwareDependency method), 124
- install_core_dependencies() (in module pimlico), 188
- install_dependencies() (in module pimlico.core.dependencies.base), 124
- installable() (JavaDependency method), 125
- installable() (JavaJarsDependency method), 125
- installable() (NLTKResource method), 128
- installable() (Py4JSoftwareDependency method), 126
- installable() (PythonPackageOnPip method), 127
- installable() (PythonPackageSystemwideInstall method), 127
- installable() (SoftwareDependency method), 123
- installable() (SystemCommandDependency method), 124
- installation_instructions() (PythonPackageSystemwideInstall method), 127
- installation_instructions() (SoftwareDependency method), 123
- InstallationError, 124
- InstallCmd (class in pimlico.cli.check), 117
- instantiate_from_options() (pimlico.datatypes.base.PimlicoDatatype class method), 170
- instantiate_output_datatype() (BaseModuleInfo method), 141
- instantiate_output_datatype_decorator() (in module pimlico.cli.debug.stepper), 114
- instantiate_output_reader() (BaseModuleInfo method), 141
- instantiate_output_reader_setup() (BaseModuleInfo method), 141
- IntegerListDocumentType (class in pimlico.datatypes.corpora.ints), 165
- IntegerListsDocumentType (class in pimlico.datatypes.corpora.ints), 164
- IntegerTableDocumentType (class in pimlico.datatypes.corpora.table), 165
- internal_available() (DataPointType.Document method), 160
- internal_data (DataPointType.Document attribute), 160
- internal_to_raw() (DataPointType.Document method), 160
- InternalModuleConnection (class in pimlico.core.modules.multistage), 151
- invalid_doc_on_error() (in module pimlico.core.modules.map), 136
- invalid_docs_on_error() (in module pimlico.core.modules.map), 136
- InvalidDocument (class in pimlico.datatypes.corpora.data_points), 160
- InvalidDocumentFormatter (class in pimlico.cli.browser.tools.formatter), 113
- is_binary_file() (in module pimlico.cli.browser.tools.files), 112
- is_binary_string() (in module pimlico.cli.browser.tools.files), 112
- is_filter() (pimlico.core.modules.base.BaseModuleInfo class method), 142
- is_identifier() (in module pimlico.utils.core), 180
- is_input() (pimlico.core.modules.base.BaseModuleInfo class method), 143
- is_locked() (BaseModuleInfo method), 145
- is_locked() (MultistageModuleInfo method), 150
- is_multiple_input() (BaseModuleInfo method), 141
- is_type_for_doc() (DataPointType method), 159
- iterable_input_reader() (in module pimlico.core.modules.inputs), 148
- IterableCorpus (class in pimlico.datatypes.corpora.base), 157
- ## J
- jars (Py4JSoftwareDependency attribute), 126
- JavaDependency (class in pimlico.core.dependencies.java), 125
- JavaJarsDependency (class in pimlico.core.dependencies.java), 125
- JavaProcessError, 130
- json_dict() (in module pimlico.core.modules.options), 152
- json_string() (in module pimlico.core.modules.options), 152
- ## K
- keypress() (InputDialog method), 111
- keys (DataPointType.Document attribute), 160
- ## L
- launch_gateway() (in module pimlico.core.external.java), 129

- launch_shell() (in module pimlico.cli.shell.runner), 116
 length (StreamCommunicationPacket attribute), 180
 length_struct (IntegerListsDocumentType attribute), 165
 limited_shuffle() (in module pimlico.utils.probability), 184
 limited_shuffle_numpy() (in module pimlico.utils.probability), 184
 list_archive_iter() (FilterModuleOutputReader method), 131
 ListDialogDisplay (class in pimlico.utils.urwid), 187
 ListStoresCmd (class in pimlico.cli.locations), 119
 LittleOutputtingProgressBar (class in pimlico.utils.progress), 185
 load() (PipelineConfig static method), 154
 load_datatype() (in module pimlico.datatypes), 177
 load_executor() (BaseModuleInfo method), 138
 load_formatter() (in module pimlico.cli.browser.tools.formatter), 113
 load_local_config() (PipelineConfig static method), 154
 load_module_executor() (in module pimlico.core.modules.base), 146
 load_module_info() (in module pimlico.core.modules.base), 146
 load_pipeline() (TestPipeline static method), 178
 LoadCmd (class in pimlico.cli.loaddump), 118
 lock() (BaseModuleInfo method), 144
 lock_path (BaseModuleInfo attribute), 144
 long_term_store (PipelineConfig attribute), 153
- ## M
- main_module (BaseModuleInfo attribute), 138
 make_py4j_errors_safe() (in module pimlico.core.external.java), 130
 make_table() (in module pimlico.utils.docs.rest), 179
 MessageDialog (class in pimlico.cli.browser.tools.corpus), 111
 MessagePopupLauncher (class in pimlico.cli.browser.tools.corpus), 112
 metadata (FilterModuleOutputReader attribute), 131
 metadata_defaults (DataPointType attribute), 159
 metadata_defaults (FloatListsDocumentType attribute), 161
 metadata_defaults (IntegerListDocumentType attribute), 165
 metadata_defaults (IntegerListsDocumentType attribute), 164
 metadata_defaults (IntegerTableDocumentType attribute), 165
 metadata_filename (BaseModuleInfo attribute), 139
 MetadataCmd (class in pimlico.cli.shell.commands), 116
 missing_data() (BaseModuleInfo method), 143
 missing_module_data() (BaseModuleInfo method), 142
 module_dependencies (PipelineConfig attribute), 153
 module_dependents (PipelineConfig attribute), 153
 module_executable (BaseModuleInfo attribute), 138
 module_executable (InputModuleInfo attribute), 147
 module_executable (MultistageModuleInfo attribute), 149
 module_executor_override (BaseModuleInfo attribute), 138
 module_inputs (BaseModuleInfo attribute), 138
 module_number_to_name() (in module pimlico.cli.util), 122
 module_numbers_to_names() (in module pimlico.cli.util), 122
 module_optional_inputs (BaseModuleInfo attribute), 138
 module_optional_outputs (BaseModuleInfo attribute), 138
 module_options (BaseModuleInfo attribute), 138
 module_outputs (BaseModuleInfo attribute), 138
 module_outputs (DocumentMapModuleInfo attribute), 135
 module_package_name() (pimlico.core.modules.base.BaseModuleInfo class method), 144
 module_readable_name (BaseModuleInfo attribute), 138
 module_status() (in module pimlico.cli.status), 122
 module_status_color() (in module pimlico.cli.status), 121
 module_type_name (BaseModuleInfo attribute), 138
 module_type_name (InputModuleInfo attribute), 147
 ModuleAlreadyCompletedError, 147
 ModuleConnection (class in pimlico.core.modules.multistage), 151
 ModuleExecutionError, 147
 ModuleExecutorLoadError, 145
 ModuleInfoLoadError, 145
 ModuleInputConnection (class in pimlico.core.modules.multistage), 151
 ModuleNotReadyError, 147
 ModuleOptionParseError, 152
 ModuleOutputConnection (class in pimlico.core.modules.multistage), 151
 modules (PipelineConfig attribute), 153
 ModuleStage (class in pimlico.core.modules.multistage), 150
 ModuleTypeError, 145
 move_dir_with_progress() (in module pimlico.utils.filesystem), 182
 MoveStoresCmd (class in pimlico.cli.locations), 119
 msgbox() (in module pimlico.utils.urwid), 187
 multiline_tablate() (in module pimlico.utils.format), 182
 MultipleInputs (class in pimlico.datatypes.base), 171
 multiprocessing_executor_factory() (in module pimlico.core.modules.map.multiproc), 132
 MultiprocessingMapModuleExecutor (class in pimlico.core.modules.map.multiproc), 132
 MultiprocessingMapPool (class in pimlico.core.modules.map.multiproc), 132

MultiprocessingMapProcess (class in pimlico.core.modules.map.multiproc), 132
 multistage_module() (in module pimlico.core.modules.multistage), 150
 MultistageModuleInfo (class in pimlico.core.modules.multistage), 149
 MultistageModulePreparationError, 151
 multiwith() (in module pimlico.utils.core), 180

N

name (DataPointType attribute), 159
 named_storage_locations (PipelineConfig attribute), 154
 NamedFile (class in pimlico.datatypes.files), 176
 NamedFileCollection (class in pimlico.datatypes.files), 175
 new_client() (Py4JInterface method), 129
 new_filename() (in module pimlico.utils.filesystem), 182
 NewModuleCmd (class in pimlico.cli.newmodule), 120
 next_document() (CorpusState method), 111
 NLTKResource (class in pimlico.core.dependencies.python), 128
 no_retry_gateway() (in module pimlico.core.external.java), 129
 non_filter_datatype (DocumentMapOutputTypeWrapper attribute), 131
 NonOutputtingProgressBar (class in pimlico.utils.progress), 185
 NonPTBTagError, 184
 normalize_cell() (in module pimlico.utils.docs.rest), 179
 notify_no_more_inputs() (DocumentMapProcessMixin method), 137
 notify_no_more_inputs() (DocumentProcessorPool method), 137
 notify_no_more_inputs() (MultiprocessingMapPool method), 132
 notify_no_more_inputs() (MultiprocessingMapProcess method), 132
 notify_no_more_inputs() (ThreadingMapThread method), 134
 NumpyArray (class in pimlico.datatypes.arrays), 167

O

on_exit() (DialogDisplay method), 187
 on_exit() (ListDialogDisplay method), 187
 opt_type_example() (in module pimlico.core.modules.options), 152
 opt_type_help() (in module pimlico.core.modules.options), 152
 option_message() (in module pimlico.cli.debug.stepper), 114
 options_dialog() (in module pimlico.utils.urwid), 187
 output_name (DocumentMapOutputTypeWrapper attribute), 131
 output_names (BaseModuleInfo attribute), 139

output_p4j_error_info() (in module pimlico.core.external.java), 130
 output_path (PipelineConfig attribute), 154
 output_ready() (BaseModuleInfo method), 141
 output_stack_trace() (in module pimlico.cli.debug), 115
 OutputCmd (class in pimlico.cli.locations), 118
 OutputConsumer (class in pimlico.core.external.java), 130
 OutputQueue (class in pimlico.utils.pipes), 183

P

palette (DialogDisplay attribute), 187
 path_relative_to_config() (PipelineConfig method), 153
 pimlico (module), 188
 pimlico.cfg (module), 187
 pimlico.cli (module), 123
 pimlico.cli.browser (module), 114
 pimlico.cli.browser.tool (module), 114
 pimlico.cli.browser.tools (module), 114
 pimlico.cli.browser.tools.corpus (module), 111
 pimlico.cli.browser.tools.files (module), 112
 pimlico.cli.browser.tools.formatter (module), 112
 pimlico.cli.check (module), 117
 pimlico.cli.clean (module), 117
 pimlico.cli.debug (module), 115
 pimlico.cli.debug.stepper (module), 114
 pimlico.cli.loaddump (module), 118
 pimlico.cli.locations (module), 118
 pimlico.cli.main (module), 119
 pimlico.cli.newmodule (module), 120
 pimlico.cli.pysshell (module), 120
 pimlico.cli.reset (module), 121
 pimlico.cli.run (module), 121
 pimlico.cli.shell (module), 117
 pimlico.cli.shell.base (module), 115
 pimlico.cli.shell.commands (module), 116
 pimlico.cli.shell.runner (module), 116
 pimlico.cli.status (module), 121
 pimlico.cli.subcommands (module), 122
 pimlico.cli.testemail (module), 122
 pimlico.cli.util (module), 122
 pimlico.core (module), 157
 pimlico.core.config (module), 153
 pimlico.core.dependencies (module), 129
 pimlico.core.dependencies.base (module), 123
 pimlico.core.dependencies.core (module), 125
 pimlico.core.dependencies.java (module), 125
 pimlico.core.dependencies.python (module), 127
 pimlico.core.dependencies.versions (module), 129
 pimlico.core.external (module), 130
 pimlico.core.external.java (module), 129
 pimlico.core.logs (module), 157
 pimlico.core.modules (module), 152
 pimlico.core.modules.base (module), 138

- pimlico.core.modules.execute (module), 146
- pimlico.core.modules.inputs (module), 147
- pimlico.core.modules.map (module), 135
- pimlico.core.modules.map.filter (module), 130
- pimlico.core.modules.map.multiproc (module), 132
- pimlico.core.modules.map.singleproc (module), 133
- pimlico.core.modules.map.threaded (module), 134
- pimlico.core.modules.multistage (module), 149
- pimlico.core.modules.options (module), 152
- pimlico.core.paths (module), 157
- pimlico.datatypes (module), 177
- pimlico.datatypes.arrays (module), 167
- pimlico.datatypes.base (module), 168
- pimlico.datatypes.core (module), 172
- pimlico.datatypes.corpora (module), 167
- pimlico.datatypes.corpora.base (module), 157
- pimlico.datatypes.corpora.data_points (module), 158
- pimlico.datatypes.corpora.floats (module), 161
- pimlico.datatypes.corpora.grouped (module), 163
- pimlico.datatypes.corpora.ints (module), 164
- pimlico.datatypes.corpora.table (module), 165
- pimlico.datatypes.corpora.tokenized (module), 166
- pimlico.datatypes.dictionary (module), 172
- pimlico.datatypes.embeddings (module), 173
- pimlico.datatypes.features (module), 174
- pimlico.datatypes.files (module), 175
- pimlico.datatypes.gensim (module), 176
- pimlico.datatypes.sklearn (module), 177
- pimlico.modules (module), 40
- pimlico.modules.candc (module), 40
- pimlico.modules.corenlp (module), 41
- pimlico.modules.corpora (module), 43
- pimlico.modules.corpora.concat (module), 43
- pimlico.modules.corpora.corpus_stats (module), 44
- pimlico.modules.corpora.format (module), 45
- pimlico.modules.corpora.group (module), 46
- pimlico.modules.corpora.interleave (module), 47
- pimlico.modules.corpora.list_filter (module), 49
- pimlico.modules.corpora.split (module), 49
- pimlico.modules.corpora.store (module), 51
- pimlico.modules.corpora.subset (module), 51
- pimlico.modules.corpora.vocab_builder (module), 53
- pimlico.modules.corpora.vocab_counter (module), 54
- pimlico.modules.corpora.vocab_mapper (module), 55
- pimlico.modules.embeddings (module), 56
- pimlico.modules.embeddings.dependencies (module), 56
- pimlico.modules.embeddings.store_embeddings (module), 58
- pimlico.modules.embeddings.store_tsv (module), 58
- pimlico.modules.embeddings.store_word2vec (module), 59
- pimlico.modules.embeddings.word2vec (module), 60
- pimlico.modules.features (module), 61
- pimlico.modules.features.term_feature_compiler (module), 61
- pimlico.modules.features.term_feature_matrix_builder (module), 63
- pimlico.modules.features.vocab_builder (module), 63
- pimlico.modules.features.vocab_mapper (module), 65
- pimlico.modules.gensim (module), 65
- pimlico.modules.gensim.lda (module), 66
- pimlico.modules.gensim.lda_doc_topics (module), 68
- pimlico.modules.input (module), 69
- pimlico.modules.input.embeddings (module), 69
- pimlico.modules.input.embeddings.fasttext (module), 69
- pimlico.modules.input.embeddings.fasttext_gensim (module), 70
- pimlico.modules.input.embeddings.glove (module), 71
- pimlico.modules.input.embeddings.word2vec (module), 72
- pimlico.modules.input.text (module), 73
- pimlico.modules.input.text.raw_text_files (module), 73
- pimlico.modules.input.text.annotations (module), 74
- pimlico.modules.malt (module), 75
- pimlico.modules.malt.conll_parser_input (module), 75
- pimlico.modules.malt.parse (module), 75
- pimlico.modules.nltk (module), 77
- pimlico.modules.nltk.nist_tokenize (module), 77
- pimlico.modules.opennlp (module), 78
- pimlico.modules.opennlp.coreference (module), 78
- pimlico.modules.opennlp.coreference_pipeline (module), 79
- pimlico.modules.opennlp.ner (module), 81
- pimlico.modules.opennlp.parse (module), 82
- pimlico.modules.opennlp.pos (module), 83
- pimlico.modules.opennlp.tokenize (module), 84
- pimlico.modules.r (module), 85
- pimlico.modules.r.script (module), 85
- pimlico.modules.regex (module), 86
- pimlico.modules.regex.annotated_text (module), 86
- pimlico.modules.sklearn (module), 87
- pimlico.modules.sklearn.logistic_regression (module), 88
- pimlico.modules.sklearn.matrix_factorization (module), 89
- pimlico.modules.text (module), 90
- pimlico.modules.text.char_tokenize (module), 90
- pimlico.modules.text.normalize (module), 91
- pimlico.modules.text.simple_tokenize (module), 92
- pimlico.modules.text.text_normalize (module), 93
- pimlico.modules.text.untokenize (module), 94
- pimlico.modules.utility (module), 95
- pimlico.modules.utility.alias (module), 95
- pimlico.modules.utility.collect_files (module), 97
- pimlico.modules.utility.copy_file (module), 98
- pimlico.modules.visualization (module), 99
- pimlico.modules.visualization.bar_chart (module), 99

- pimlico.modules.visualization.embeddings_plot (module), 100
- pimlico.test (module), 179
- pimlico.test.pipeline (module), 178
- pimlico.test.suite (module), 179
- pimlico.utils (module), 187
- pimlico.utils.communicate (module), 179
- pimlico.utils.core (module), 180
- pimlico.utils.docs (module), 179
- pimlico.utils.docs.rest (module), 179
- pimlico.utils.email (module), 181
- pimlico.utils.filesystem (module), 182
- pimlico.utils.format (module), 182
- pimlico.utils.linguistic (module), 182
- pimlico.utils.logging (module), 183
- pimlico.utils.network (module), 183
- pimlico.utils.pipes (module), 183
- pimlico.utils.pos (module), 183
- pimlico.utils.probability (module), 184
- pimlico.utils.progress (module), 185
- pimlico.utils.strings (module), 186
- pimlico.utils.system (module), 186
- pimlico.utils.timeout (module), 186
- pimlico.utils.urwid (module), 187
- pimlico.utils.web (module), 187
- PimlicoCLISubcommand (class in pimlico.cli.subcommands), 122
- PimlicoDatatype (class in pimlico.datatypes.base), 168
- PimlicoJavaLibrary (class in pimlico.core.dependencies.java), 126
- PimlicoPythonShellContext (class in pimlico.cli.pyshell), 120
- PipelineCheckError, 156
- PipelineConfig (class in pimlico.core.config), 153
- PipelineConfigParseError, 155
- PipelineStructureError, 155
- POOL_TYPE (MultiprocessingMapModuleExecutor attribute), 132
- POOL_TYPE (ThreadingMapModuleExecutor attribute), 134
- pos_tag_to_ptb() (in module pimlico.utils.pos), 183
- pos_tags_to_ptb() (in module pimlico.utils.pos), 183
- postloop() (DataShell method), 115
- postprocess() (DocumentMapModuleExecutor method), 136
- postprocess() (MultiprocessingMapModuleExecutor method), 132
- postprocess() (ThreadingMapModuleExecutor method), 134
- preloop() (DataShell method), 115
- preprocess() (DocumentMapModuleExecutor method), 136
- preprocess_config_file() (in module pimlico.core.config), 156
- print_dependency_leaf_problems() (in module pimlico.core.config), 156
- print_execution_error() (in module pimlico.cli.util), 123
- print_missing_dependencies() (in module pimlico.core.config), 156
- problems() (JavaDependency method), 125
- problems() (NLTKResource method), 128
- problems() (PythonPackageDependency method), 127
- problems() (SoftwareDependency method), 123
- problems() (SystemCommandDependency method), 124
- process_document() (DocumentMapProcessMixin method), 137
- process_documents() (DocumentMapProcessMixin method), 137
- process_module_options() (in module pimlico.core.modules.options), 152
- process_module_options() (pimlico.core.modules.base.BaseModuleInfo class method), 139
- process_setup() (FilterModuleOutputReader method), 131
- PROCESS_TYPE (MultiprocessingMapPool attribute), 132
- ProcessOutput (class in pimlico.core.modules.map), 136
- ProgressBarIter (class in pimlico.utils.progress), 186
- prompt (DataShell attribute), 115
- provide_further_outputs() (BaseModuleInfo method), 140
- Py4JInterface (class in pimlico.core.external.java), 129
- Py4JSafeJavaError, 130
- Py4JSoftwareDependency (class in pimlico.core.dependencies.java), 126
- PythonCmd (class in pimlico.cli.shell.commands), 116
- PythonPackageDependency (class in pimlico.core.dependencies.python), 127
- PythonPackageOnPip (class in pimlico.core.dependencies.python), 127
- PythonPackageSystemwideInstall (class in pimlico.core.dependencies.python), 127
- PythonShellCmd (class in pimlico.cli.pyshell), 120
- ## Q
- qget() (in module pimlico.utils.pipes), 183
- ## R
- raw_available() (DataPointType.Document method), 160
- raw_data (DataPointType.Document attribute), 160
- raw_to_internal() (DataPointType.Document method), 160
- RawDocumentType (class in pimlico.datatypes.corpora.data_points), 161
- RawTextDocumentType (class in pimlico.datatypes.corpora.data_points), 161
- read() (DummyFileDescriptor method), 185

- read() (StreamCommunicationPacket static method), 180
- Reader (Dict attribute), 172
- Reader (Dictionary attribute), 172
- Reader (Embeddings attribute), 174
- Reader (GensimLdaModel attribute), 177
- Reader (GroupedCorpus attribute), 163
- Reader (IterableCorpus attribute), 158
- Reader (NamedFile attribute), 176
- Reader (NamedFileCollection attribute), 175
- Reader (NumpyArray attribute), 167
- Reader (PimlicoDatatype attribute), 170
- Reader (ScipySparseMatrix attribute), 168
- Reader (ScoredRealFeatureSets attribute), 174
- Reader (SklearnModel attribute), 177
- Reader (StringList attribute), 172
- Reader (TextFile attribute), 176
- Reader (TSVVecFiles attribute), 174
- reader_init() (DataPointType method), 159
- reader_init() (FloatListDocumentType method), 162
- reader_init() (FloatListsDocumentType method), 161
- reader_init() (IntegerListDocumentType method), 165
- reader_init() (IntegerListsDocumentType method), 164
- reader_init() (IntegerTableDocumentType method), 165
- reader_init() (VectorDocumentType method), 162
- readLine() (DummyFileDescriptor method), 185
- recursive_deps() (in module pimlico.core.dependencies.base), 124
- remove_duplicates() (in module pimlico.utils.core), 180
- remove_temporary_redirects() (OutputConsumer method), 130
- reset_all_modules() (PipelineConfig method), 153
- reset_execution() (BaseModuleInfo method), 144
- reset_execution() (MultistageModuleInfo method), 150
- ResetCmd (class in pimlico.cli.reset), 121
- retrieve_processing_status() (DocumentMapModuleExecutor method), 136
- retry_open() (in module pimlico.utils.filesystem), 182
- run() (InputQueueFeeder method), 136
- run() (MultiprocessingMapProcess method), 132
- run() (OutputConsumer method), 130
- run() (ThreadingMapThread method), 134
- run_browser() (IterableCorpus method), 158
- run_browser() (NamedFileCollection method), 175
- run_browser() (PimlicoDatatype method), 170
- run_command() (BrowseCmd method), 120
- run_command() (CleanCmd method), 117
- run_command() (DepsCmd method), 117
- run_command() (DumpCmd method), 118
- run_command() (EmailCmd method), 122
- run_command() (InputsCmd method), 118
- run_command() (InstallCmd method), 117
- run_command() (ListStoresCmd method), 119
- run_command() (LoadCmd method), 118
- run_command() (MoveStoresCmd method), 119
- run_command() (NewModuleCmd method), 120
- run_command() (OutputCmd method), 119
- run_command() (PimlicoCLISubcommand method), 122
- run_command() (PythonShellCmd method), 121
- run_command() (ResetCmd method), 121
- run_command() (RunCmd method), 121
- run_command() (ShellCLICmd method), 116
- run_command() (StatusCmd method), 121
- run_command() (UnlockCmd method), 120
- run_command() (VariantsCmd method), 119
- run_command() (VisualizeCmd method), 120
- run_test_pipeline() (in module pimlico.test.pipeline), 178
- run_test_suite() (in module pimlico.test.pipeline), 178
- RunCmd (class in pimlico.cli.run), 121
- ## S
- safe_import_bs4() (in module pimlico.core.dependencies.python), 128
- SafeProgressBar (class in pimlico.utils.progress), 185
- satisfies_typecheck() (in module pimlico.core.modules.base), 145
- save_popup_launcher() (in module pimlico.cli.browser.tools.corpus), 112
- ScipySparseMatrix (class in pimlico.datatypes.arrays), 167
- ScoredRealFeatureSets (class in pimlico.datatypes.features), 174
- SegmentedLinesDocumentType (class in pimlico.datatypes.corpora.tokenized), 166
- send_final_report_email() (in module pimlico.core.modules.execute), 147
- send_module_report_email() (in module pimlico.core.modules.execute), 147
- send_pimlico_email() (in module pimlico.utils.email), 181
- send_text_email() (in module pimlico.utils.email), 181
- sequential_document_sample() (in module pimlico.utils.probability), 184
- sequential_sample() (in module pimlico.utils.probability), 184
- set_metadata_value() (BaseModuleInfo method), 139
- set_metadata_values() (BaseModuleInfo method), 139
- set_proc_title() (in module pimlico.utils.system), 186
- set_up() (DocumentMapProcessMixin method), 137
- Setup (FilterModuleOutputReader attribute), 131
- shell_commands (IterableCorpus attribute), 158
- shell_commands (PimlicoDatatype attribute), 169
- ShellCLICmd (class in pimlico.cli.shell.runner), 116
- ShellCommand (class in pimlico.cli.shell.base), 115
- ShellContextError, 121
- ShellError, 116
- short_term_store (PipelineConfig attribute), 153
- shutdown() (DocumentProcessorPool method), 137
- shutdown() (InputQueueFeeder method), 137

- shutdown() (MultiprocessingMapPool method), 132
- shutdown() (ThreadingMapPool method), 134
- shutdown() (ThreadingMapThread method), 134
- signals (InputDialog attribute), 111
- similarities() (in module pimlico.utils.strings), 186
- single_process_executor_factory() (in module pimlico.core.modules.map.singleproc), 133
- SINGLE_PROCESS_TYPE (MultiprocessingMapPool attribute), 132
- SingleThreadMapModuleExecutor (class in pimlico.core.modules.map.singleproc), 133
- skip() (CorpusState method), 111
- skip_invalid() (in module pimlico.core.modules.map), 136
- skip_invalids() (in module pimlico.core.modules.map), 136
- skip_popup_launcher() (in module pimlico.cli.browser.tools.corpus), 112
- SklearnModel (class in pimlico.datatypes.sklearn), 177
- slice_progress() (in module pimlico.utils.progress), 186
- SoftwareDependency (class in pimlico.core.dependencies.base), 123
- SoftwareVersion (class in pimlico.core.dependencies.versions), 129
- sorted_by_similarity() (in module pimlico.utils.strings), 186
- split_seq() (in module pimlico.utils.core), 180
- split_seq_after() (in module pimlico.utils.core), 180
- stages (MultistageModuleInfo attribute), 149
- start() (LittleOutputtingProgressBar method), 185
- start() (Py4JInterface method), 129
- start_java_process() (in module pimlico.core.external.java), 129
- start_worker() (MultiprocessingMapPool method), 132
- start_worker() (ThreadingMapPool method), 134
- status (BaseModuleInfo attribute), 139
- status (MultistageModuleInfo attribute), 150
- status_colored() (in module pimlico.cli.status), 121
- StatusCmd (class in pimlico.cli.status), 121
- step (PipelineConfig attribute), 155
- Stepper (class in pimlico.cli.debug.stepper), 114
- stop() (Py4JInterface method), 129
- StopProcessing, 147
- store_names (PipelineConfig attribute), 154
- str_to_bool() (in module pimlico.core.modules.options), 152
- StreamCommunicationError, 180
- StreamCommunicationPacket (class in pimlico.utils.communicate), 179
- StringList (class in pimlico.datatypes.core), 172
- strip_punctuation() (in module pimlico.utils.linguistic), 182
- struct (IntegerListDocumentType attribute), 165
- struct (IntegerListsDocumentType attribute), 165
- subsample() (in module pimlico.utils.probability), 185
- SystemCommandDependency (class in pimlico.core.dependencies.base), 124
- ## T
- table_div() (in module pimlico.utils.docs.rest), 179
- tear_down() (DocumentMapProcessMixin method), 137
- terminate_process() (in module pimlico.utils.communicate), 179
- test_all_modules() (TestPipeline method), 178
- test_input_module() (TestPipeline method), 178
- test_module_execution() (TestPipeline method), 178
- TestPipeline (class in pimlico.test.pipeline), 178
- TestPipelineRunError, 178
- TextDocumentType (class in pimlico.datatypes.corpora.data_points), 161
- TextFile (class in pimlico.datatypes.files), 176
- THREAD_TYPE (ThreadingMapPool attribute), 134
- threading_executor_factory() (in module pimlico.core.modules.map.threaded), 134
- ThreadingMapModuleExecutor (class in pimlico.core.modules.map.threaded), 134
- ThreadingMapPool (class in pimlico.core.modules.map.threaded), 134
- ThreadingMapThread (class in pimlico.core.modules.map.threaded), 134
- timeout() (in module pimlico.utils.timeout), 186
- timeout_process() (in module pimlico.utils.communicate), 179
- title_box() (in module pimlico.utils.format), 182
- TokenizedDocumentType (class in pimlico.datatypes.corpora.tokenized), 166
- trim_docstring() (in module pimlico.utils.docs), 179
- truncate() (in module pimlico.utils.strings), 186
- TSVVecFiles (class in pimlico.datatypes.embeddings), 174
- type_checking_name() (DynamicInputDatatypeRequirement method), 171
- type_checking_name() (IterableCorpus method), 158
- type_checking_name() (PimlicoDatatype method), 170
- typecheck_formatter() (in module pimlico.cli.browser.tools.formatter), 113
- typecheck_input() (BaseModuleInfo method), 143
- typecheck_inputs() (BaseModuleInfo method), 143
- typecheck_inputs() (MultistageModuleInfo method), 149
- TypeCheckError, 145
- ## U
- unhandled_key() (ListDialogDisplay method), 187
- unlock() (BaseModuleInfo method), 145
- UnlockCmd (class in pimlico.cli.main), 119
- update() (SafeProgressBar method), 185
- update_processing_status() (DocumentMapModuleExecutor method), 136

V

VariantsCmd (class in pimlico.cli.main), 119
 VectorDocumentType (class in pimlico.datatypes.corpora.floats), 162
 VectorFormatter (class in pimlico.datatypes.corpora.floats), 163
 VisualizeCmd (class in pimlico.cli.main), 120

W

WorkerShutdownError, 137
 WorkerStartupError, 137
 wrap_module_info_as_filter() (in module pimlico.core.modules.map.filter), 131
 wrap_tarred_corpus() (in module pimlico.cli.debug.stepper), 114
 wrapped_module_info (DocumentMapOutputTypeWrapper attribute), 131
 write() (DummyFileDescriptor method), 185
 Writer (Dict attribute), 172
 Writer (Dictionary attribute), 173
 Writer (Embeddings attribute), 174
 Writer (GensimLdaModel attribute), 177
 Writer (GroupedCorpus attribute), 163
 Writer (IterableCorpus attribute), 158
 Writer (NamedFile attribute), 176
 Writer (NamedFileCollection attribute), 175
 Writer (NumpyArray attribute), 167
 Writer (PimlicoDatatype attribute), 171
 Writer (ScipySparseMatrix attribute), 168
 Writer (ScoredRealFeatureSets attribute), 174
 Writer (SklearnModel attribute), 177
 Writer (StringList attribute), 172
 Writer (TextFile attribute), 176
 Writer (TSVVecFiles attribute), 174
 writer_init() (DataPointType method), 159
 writer_init() (FloatListDocumentType method), 162
 writer_init() (FloatListsDocumentType method), 162
 writer_init() (IntegerListDocumentType method), 165
 writer_init() (IntegerListsDocumentType method), 164
 writer_init() (IntegerTableDocumentType method), 166
 writer_init() (VectorDocumentType method), 163

Y

yesno_dialog() (in module pimlico.utils.urwid), 187