
Eskapade Documentation

**KPMG Advanced Analytics
Big Data team**

Dec 09, 2018

Contents

1	Documentation	3
2	Check it out	5
3	Quick run	7
4	Release notes	9
5	Contact and support	11
6	Contents	13
6.1	Introduction	13
6.2	Installation	14
6.3	Tutorials	15
6.4	Command Line Arguments	36
6.5	Package structure	40
6.6	Release notes	44
6.7	Developing and Contributing	47
6.8	References	48
6.9	API	49
6.10	Appendices	100
6.11	Indices and tables	103
	Python Module Index	105

- Version: 0.9.3
- Released: December 2018
- Web page: <http://eskapade.kave.io>
- Repository: <https://github.com/kaveio/eskapade>
- Docker: <https://github.com/kaveio/eskapade-environment>

Eskapade is a light-weight, python-based data analysis framework, meant for developing and modularizing all sorts of data analysis problems into reusable analysis components.

Eskapade uses a modular approach to analytics, meaning that you can use pre-made operations (called ‘links’) to build an analysis. This is implemented in a chain-link framework, where you define a ‘Chain’, consisting of a number of Links. These links are the fundamental building block of your analysis. For example, a data loading link and a data transformation link will frequently be found together in a pre-processing Chain.

Each chain has a specific purpose, for example: data quality checks of incoming data, pre-processing of data, booking and/or training of predictive algorithms, validation of these algorithms, and their evaluation. By using this work methodology, analysis links can be more easily reused in future data analysis projects.

Eskapade is analysis-library agnostic. It is used to set up typical data analysis problems from multiple packages, e.g.: scikit-learn, Spark MLlib, and ROOT. Likewise, Eskapade can use a manner of different data structures to handle data, such as: pandas DataFrames, numpy arrays, Spark DataFrames/RDDs, and more.

For example, Eskapade has been used as a self-learning framework for typical machine learning problems. Trained algorithms can predict real-time or batch data, these models can be evaluated over time, and Eskapade can bookkeep and retrain their algorithms.

CHAPTER 1

Documentation

The entire Eskapade documentation including tutorials can be found [here](#).

CHAPTER 2

Check it out

Eskapade requires Python 3 and is pip friendly. To get started, simply do:

```
$ pip install Eskapade
```

or check out the code from our github repository:

```
$ git clone https://github.com/KaveIO/Eskapade.git
$ pip install -e Eskapade/
```

where in this example the code is installed in edit mode (option -e).

You can now use Eskapade in Python with:

```
import eskapade
```

Congratulations, you are now ready to use Eskapade!

CHAPTER 3

Quick run

To see the available Eskapade example, do:

```
$ export TUTDIR=`pip show Eskapade | grep Location | awk '{ print $2"/eskapade/
↪tutorials" }'`
$ ls -l $TUTDIR/
```

E.g. you can now run:

```
$ eskapade_run $TUTDIR/tutorial_1.py
```

For all available Eskapade example macros, please see our [tutorials section](#).

CHAPTER 4

Release notes

The Eskapade patch release v0.9.0 and corresponding docker images have the following features:

- The core functionality of Eskapade, namely: the `Link`, `Chain`, `process_manager`, `DataStore`, `ConfigObject` and corresponding tutorials, have been split off from the growing (analysis) Eskapade repository, into the new package `Eskapade-Core`. `Eskapade-Core` is a very light-weight Python3 package.
- A new module `data_mimic` has been add to Eskapade, including tutorials, meant for resimulating existing datasets.
- We have added `feather` i/o functionality for reading and writeng dataframes.
- The logger has been fixed, it is now possible to set the log-level of loggers again.
- The Eskapade docker files have been taken out of the Eskapade repository to avoid version conflicts, into the new git repo `Eskapade-Environment`.
- The Eskapade docker image `eskapade-usr` contain the latest working versions of Eskapade, `Eskapade-Core`, `Eskapade-ROOT`, and `Eskapade-Spark`. Type:

```
$ docker pull kave/eskapade-usr:latest
```

to pull it in.

See [release notes](#) for previous versions of Eskapade.

CHAPTER 5

Contact and support

- Issues & Ideas: <https://github.com/kaveio/eskapade/issues>
- Q&A Support: contact us at: kave [at] kpmg [dot] com

Please note that the KPMG Eskapade group provides support only on a best-effort basis.

6.1 Introduction

Welcome to Eskapade! This is a short introduction of the package and why we built it. In the next sections we will go over the installation, a tutorial on how to run Eskapade properly, some examples use-cases, and how to develop analysis code and run it in Jupyter and PyCharm.

6.1.1 What is Eskapade?

Eskapade is an abbreviation for: ‘Enterprise Solution KPMG Advanced Predictive Analytics Decision Engine’. It is a light-weight Python framework to help make your data analytics modular. This results in faster roll-out of analytics solutions in your business and less overhead when taking multiple analyses in production. In particular, it is intended for building data processing pipe lines and using implemented machine learning models in a production environment.

6.1.2 Why did we build this?

We found that the implementation phase of a data analytics solution at clients - a careful process - is often completely different from building the solution itself - which proceeds through explorative iterations. Therefore we made this analysis framework that makes it easier to set up a data analysis, while simultaneously making it easier to put it into production.

Next to that, it makes analyses modular, which has a lot of advantages. It is easier to work with multiple people on the same project, because the contributions are divided in clear blocks. Re-use of code becomes more straightforward, as old code is already put in a block with a clear purpose. Lastly, it gives you a universal basis for all your analyses, which can both be used across a company, or for different clients.

More about the purpose can be read at the general [readme](#).

6.2 Installation

Let's get Eskapade up and running! In order to make this as easy as possible, we provide both a Docker image and a virtual machine where everything you need is installed and working properly. Alternatively, you can download the repository and run it on your own machine.

- See *Eskapade on your own machine* for the local installation requirements.
- See *Eskapade with Docker* to get started with Docker.
- See *Eskapade on a virtual machine* to get started with Vagrant.

This manual is written for Linux systems, but Eskapade also runs fine on macOS systems.

6.2.1 Eskapade on your own machine

Eskapade can be installed as any other Python package with `easy_install` or `pip`. To get started, simply do:

```
$ pip install Eskapade
```

We have verified that this works on Linux 16.04 and MacOS based machines.

Or check out the code from our github repository:

```
$ git clone https://github.com/KaveIO/Eskapade.git
$ pip install -e Eskapade/
```

where the code is installed in editable mode (option `-e`).

You can now use Eskapade in Python with:

```
import eskapade
```

Congratulations, you are now ready to use Eskapade!

See the other parts of the documentation for specific usage.

Requirements

Eskapade requires Python 3 and some libraries, which can be found in `setup.py` at the root of the repository.

Eskapade can be installed as any other Python package with `easy_install` or `pip`:

```
$ pip install /path/to/eskapade
```

If `anaconda` is already installed in your machine, consider creating a `conda` virtual environment with Python 3.6 to install all the requirements and Eskapade itself to avoid collisions:

```
$ conda create -n eskapade_env36 python=3.6 anaconda
```

Then you can activate it as follows:

```
$ source activate eskapade_env36
```

More information about `conda` virtual environments can be found [here](#)

6.2.2 Eskapade with Docker

Type:

```
$ docker pull kave/eskapade-usr:latest
```

to pull in the Eskapade image from dockerhub.

For more details see the Eskapade repo with the [docker configurations](#).

6.2.3 Eskapade on a virtual machine

For detailed instruction on how to set up a vagrant box with Eskapade, go to the Eskapade repo with the [vagrant box](#).

6.2.4 Installing Eskapade on macOS

To install eskapade on macOS, see our [macOS appendix](#).

6.3 Tutorials

This section contains materials on how to use Eskapade. There are additional side notes on how certain aspects work and where to find parts of the code. For more in depth explanations on the functionality of the code-base, try the [API docs](#).

All command examples below can be run from any directory with write access.

6.3.1 Running your first macro

After successfully [installing](#) Eskapade, it is now time to run your very first macro, the classic code example: Hello World!

For ease of use, let's make shortcuts to the directories containing the Eskapade tutorials:

```
$ export TUTDIRC=`pip show Eskapade-Core | grep Location | awk '{ print $2"/escore/
↪tutorials" }'`
$ export TUTDIR=`pip show Eskapade | grep Location | awk '{ print $2"/eskapade/
↪tutorials" }'`
$ ls -l $TUTDIRC/ $TUTDIR/
```

Hello World!

If you just want to run it plain and simple, go to the root of the repository and run the following:

```
$ eskapade_run $TUTDIRC/esk101_helloworld.py
```

This will run the macro that prints out Hello World. There is a lot of output, but try to find back these lines (or similar):

```
2017-11-13T12:37:07.473512+00:00 [eskapade.core_ops.links.hello_world.HelloWorld
↪#INFO] Hello World
2017-11-13T12:37:07.473512+00:00 [eskapade.core_ops.links.hello_world.HelloWorld
↪#INFO] Hello World
```

Congratulations, you have just successfully run Eskapade!

Internal workings

To see what is actually happening under the hood, go ahead and open up `tutorials/esk101_helloworld.py`. The macro is like a recipe and it contains all of your analysis. It has all the ‘high level’ operations that are to be executed by Eskapade.

When we go into this macro we find the following piece of code:

```
hello = Chain(name='Hello')
link = core_ops.HelloWorld(name='HelloWorld')
link.logger.log_level = LogLevel.DEBUG
link.repeat = settings['n_repeat']
hello.add(link)
```

Which is the code that does the actual analysis (in this case, print out the statement). In this case `link` is an instance of the class `HelloWorld`, which itself is a `Link`. The `Link` class is the fundamental building block in Eskapade that contains our analysis steps. The code for `HelloWorld` can be found at:

```
$ less python/eskapade/core_ops/links/hello_world.py
```

Looking into this class in particular, in the code we find in the `execute()` function:

```
self.logger.info('Hello {hello}', hello=self.hello)
```

where `self.hello` is a parameter set in the `__init__` of the class. This setting can be overwritten as can be seen below. For example, we can make another link, `link2` and change the default `self.hello` into something else.

```
link2 = core_ops.HelloWorld(name='Hello2')
link2.hello = 'Lionel Richie'
ch.add(link2)
```

Rerunning results in us greeting the famous singer/songwriter.

There are many ways to run your macro and control the flow of your analysis. You can read more on this in the *Short introduction to the Framework* subsection below.

6.3.2 Tutorial 1: transforming data

Now that we know the basics of Eskapade we can go on to more advanced macros, containing an actual analysis.

Before we get started, we have to fetch some data, on your command line, type:

```
$ wget https://s3-eu-west-1.amazonaws.com/kpmg-eskapade-share/data/LAozone.data
```

To run the macro type on your CLI:

```
$ eskapade_run $TUTDIR/tutorial_1.py
```

If you want to add command line arguments, for example to change the output logging level, read the page on [command line arguments](#).

When looking at the output in the terminal we read something like the following:

```

2017-11-13T13:37:07.473512+00:00 [eskapade.core.execution#INFO] *
↳Welcome to Eskapade!
...
2017-11-13T13:37:08.085577+00:00 [eskapade.core.process_manager.ProcessManager#INFO]
↳Number of registered chains: 2
...
2017-11-13T13:37:11.316414+00:00 [eskapade.core.execution#INFO] *
↳Leaving Eskapade. Bye!

```

There is a lot more output than these lines (tens or hundred of lines depending on the log level). Eskapade has run the code from each link, and at the top of the output in your terminal you can see a summary.

When you look at the output in the terminal you can see that the macro contains two chains and a few Link are contained in these chains. Note that chain 2 is empty at this moment. In the code of the macro we see that in the first chain that data is loaded first and then a transformation is applied to this data.

Before we are going to change the code in the macro, there will be a short introduction to the framework.

Short introduction to the Framework

At this point we will not go into the underlying structure of the code that is underneath the macro, but later in this tutorial we will. For now we will take a look in the macro. So open `$TUTDIR/tutorial_1.py` in your favorite editor. We notice the structure: first imports, then defining all the settings, and finally the actual analysis: Chains and Links.

A chain is instantiated as follows:

```
data = Chain('Data')
```

and registered automatically with the ProcessManager. The ProcessManager is the main event processing loop and is responsible for processing the Chains and Links.

Next a Pandas data frame converter Link is initialized and its properties are set, and finally added to the data chain:

```
reader = analysis.ReadToDf(name='Read_LA_ozone', path='LAozone.data', reader=pd.read_
↳csv, key='data')
data.add(reader)
```

This means the Link is added to the chain and when Eskapade runs, it will execute the code in the Link.

Now that we know how the framework runs the code on a higher level, we will continue with the macro.

In the macro notice that under the second chain some code has been commented out. Uncomment the code and run the macro again with:

```
$ eskapade_run $TUTDIR/tutorial_1.py
```

And notice that it takes a bit longer to run, and the output is longer, since it now executes the Link in chain 2. This Link takes the data from chain 1 and makes plots of the data in the data set and saves it to your disk. Go to this path and open one of the pdfs found there:

```
$ results/Tutorial_1/data/v0/report/
```

The pdfs give an overview of all numerical variables in the data in histogram form. The binning, plotting and saving of this data is all done by the chain we just uncommented. If you want to take a look at how the Link works, it can be found in:

```
$ python/eskapade/visualization/links/df_summary.py
```

But for now, we will skip the underlying functionality of the links.

Let's do an exercise. Going back to the first link, we notice that the transformations that are executed are defined in `conv_funcs` passed to the link. We want to include in the plot the wind speed in km/h. There is already a part of the code available in the `conv_funcs` and the functions `comp_date` and `mi_to_km`. Use these functions as examples to write a function that converts the wind speed.

Add this to the transformation by adding your own code. Once this works you can also try to add the temperature in degrees Celsius.

Making a Link

Now we are going to add a new link that we create! To make a new link type the following:

```
$ eskapade_generate_link --dir python/eskapade/analysis/links YourLink
```

The command will make a link object named `YourLink` in the path specified in the first argument. The link we wish to add will do some textual transformation, so name it accordingly. And be sure to follow the instructions given by the command!

The command creates the skeleton file:

```
$ python/eskapade/analysis/links/yourlink.py
```

This skeleton file can be modified with your custom editor and then be imported and called inside a macro with `analysis>YourLink()`. Notice that the name of the class is CamelCase and that the name of the file is lowercase to conform to coding guidelines.

Now open up the link in your editor. In the `execute` function of the Link, we see that a `DataStore` is called. This is the central in-memory object in which all data is saved. `DataStore` inherits from a dict, so by calling the right key we can get objects. Call:

```
df = ds['data']
```

to get the `DataFrame` that includes the latest transformations.

Now we are going to make a completely different transformation in the Link and apply it to the object in the `DataStore`. We want to add a column to the data that states how humid it is. When column 'humidity' is less than 50 it is 'dry', otherwise it is 'humid'. You will have to use some pandas functionality or perhaps something else if you prefer. Save the new column back into the `DataFrame` and then put the `DataFrame` in the `DataStore` under the key 'data_new'.

We are going to let our plot functionality loose on this `DataFrame` once more, to see what happens to our generated textual data. It can not be plotted. In the future this functionality will be available for most data types.

Now run the entire macro with the new code and compile the output `.tex` file. This can be done on the command line with

```
$ cd results/Tutorial_1/data/v0/report/  
$ pdflatex report.tex
```

If you have `pdflatex` installed on your machine.

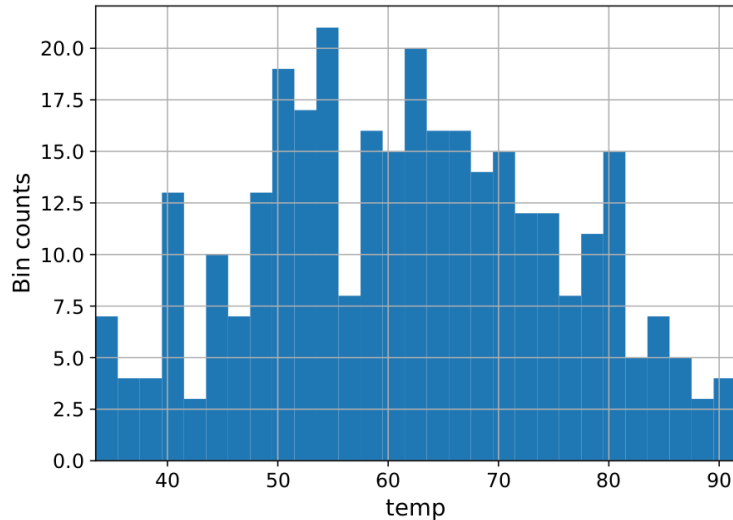
Note: If you don't have `pdflatex` installed on your machine you can install it by executing the following command: ..
code-block:: bash

```
$ yum install texlive-latex-recommended
```

Now take a look at the output pdf. The final output should look something like this:

temp

count	330
filled	330
distinct	63
mean	61.7545
std	14.4587
min	25
max	93
p01	31.29
p05	39
p16	47.64
p50	62
p84	78
p95	85.55
p99	92



Your plot should be quite similar (though it might be different in its make up.)

In summary, the work method of Eskapade is to run chains of custom code chunks (links). Each chain should have a specific purpose, for example pre-processing incoming data, booking and/or training predictive algorithms, validating these predictive algorithms, evaluating the algorithms.

By using this work method, links can be easily reused in future projects. Some links are provided by default. For example, links used to load data from a json file, book predictive algorithms, predict the training and test data set and produce evaluation metrics. If you want to use your own predictive model just go ahead and add your own links!

6.3.3 Tutorial 2: macro from basic links

In this tutorial we are going to build a macro using existing Links. We start by using templates to make a new macro. The command

```
$ eskapade_generate_macro --dir python/eskapade/tutorials tutorial_2
```

makes a new macro from a template macro. When we open the macro we find a lot of options that we can use. For now we will actually not use them, but if you want to learn more about them, read the [Examples](#) section below.

First we will add new chains to the macro. These are the higher level building blocks that can be controlled when starting a run of the macro. At the bottom of the macro we find a commented out Link, the classic Hello World link. You can uncomment it and run the macro if you like, but for now we are going to use the code to make a few chains.

So use the code and add 3 chains with different names:

```
ch = Chain('CHAINNAME')
```

When naming chains, remember that the output of Eskapade will print per chain-link combination the logs that are defined in the Links. So name the chains appropriately, so when you run the macro the logging actually makes sense.

This tutorial will be quite straight-forward, it has 3 short steps, which is why we made 3 chains.

1. In the first chain: Read a data file of your choosing into Eskapade using the pandas links in the analysis subpackage.
2. In the second chain: Copy the DataFrame you created in the DataStore using the core_ops subpackage.
3. In the third chain: Delete the entire DataStore using a Link in the core_ops subpackage.

To find the right Links you have to go through the Eskapade documentation (or code!), and to find within its subpackages the proper Links you have to understand the package structure. Every package is specific for a certain task, such as analysis, core tasks (like the `ProcessManager`), or data quality. Every subpackage contains links in its `links/` subdirectory. See for example the subpackages `core_ops`, `analysis` or `visualization`.

In *All available examples* we give some tips to find the right Links your analysis, and how to configure them properly.

6.3.4 Tutorial 3: Jupyter notebook

This section contains materials on how to use Eskapade in Jupyter Notebooks. There are additional side notes on how certain aspects work and where to find parts of the code. For more in depth explanations, try the [API-docs](#).

Next we will demonstrate how Eskapade can be run and debugged interactively from within a Jupyter notebook.

An Eskapade notebook

To run Eskapade use the `eskapade_generate_notebook` command to create a template notebook. For example:

```
$ eskapade_generate_notebook --dir ./ notebook_name
```

The minimal code you need to run a notebook is the following:

```
from eskapade import process_manager, resources, ConfigObject, DataStore
from eskapade.core import execution, persistence
from eskapade.logger import LogLevel

# --- basic config
settings = process_manager.service(ConfigObject)
settings['macro'] = resources.tutorial('tutorial_1.py')
settings['version'] = 0
settings['logLevel'] = LogLevel.DEBUG

# --- optional running parameters
# settings['beginWithChain'] = 'startChain'
# settings['endWithChain'] = 'endChain'
# settings['resultsDir'] = 'resultsdir'
settings['storeResultsEachChain'] = True

# --- other global flags (just some examples)
# settings['set_mongo'] = False
# settings['set_training'] = False

# --- run eskapade!
execution.eskapade_run(settings)

# --- To rerun eskapade, clear the memory state first!
# execution.reset_eskapade()
```

Make sure to fill out all the necessary parameters for it to run. The macro has to be set obviously, but not all settings in this example are needed to be set to a value. The function `execution.eskapade_run(settings)` runs Eskapade with the settings you specified.

To inspect the state of the Eskapade objects (datastore and configurations) after the various chains see the command line examples below. .. note:

```
Inspecting intermediate states requires Eskapade to be run with the option_
↳storeResultsEachChain
(command line: ``-w``) on.
```

```
# --- example inspecting the data store after the preprocessing chain
ds = DataStore.import_from_file('./results/Tutorial_1/proc_service_data/v0/_Summary/
↳eskapade.core.process_services.DataStore.pkl')
ds.keys()
ds.Print()
ds['data'].head()

# --- example showing Eskapade settings
settings = ConfigObject.import_from_file('./results/Tutorial_1/proc_service_data/v0/_
↳Summary/eskapade.core.process_services.ConfigObject.pkl')
settings.Print()
```

The `import_from_file` function imports a pickle file that was written out by Eskapade, containing the `DataStore`. This can be used to start from an intermediate state of your Eskapade. For example, you do some operations on your `DataStore` and then save it. At a later time you load this saved `DataStore` and continue from there.

Running in a notebook

In this tutorial we will make a notebook and run the macro from [tutorial 1](#). This macro shows the basics of Eskapade. Once we have Eskapade running in a terminal, we can run it also in Jupyter. Make sure you have properly [installed Jupyter](#).

We start by making a notebook:

```
$ eskapade_generate_notebook tutorial_3_notebook
```

This will create a notebook in the current directory with the name `tutorial_3_notebook` running macro `tutorial_1.py`. You can set a destination directory by specifying the command argument `--dir`. Now open Jupyter and take a look at the notebook.

```
$ jupyter notebook
```

Try to run the notebook. You might get an error if the notebook can not find the data for the data reader. Unless you luckily are in the right folder. By default, `tutorial_1.py` looks for the data file `LAozone.data` in the working directory. Use:

```
!pwd
```

In Jupyter to find which path you are working on, and put the data to the path. Or change the load path in the macro to the proper one. But in the end it depends on your setup.

Intermezzo: you can run bash commands in Jupyter by prepending the command with a !

Now run the cells in the notebook and check if the macro runs properly. The output be something like:

```
2017-02-14 14:04:55,506 DEBUG [link/execute_link]: Now executing link 'LA ozone data'
2017-02-14 14:04:55,506 DEBUG [readtodf/execute]: reading datasets from files ["../
↳data/LAozone.data"]
2017-02-14 14:04:55,507 DEBUG [readtodf/pandasReader]: using Pandas reader "<function_
↳_make_parser_function.<locals>.parser_f at 0x7faaac7f4d08>"
```

(continues on next page)

(continued from previous page)

```

2017-02-14 14:04:55,509 DEBUG [link/execute_link]: Done executing link 'LA ozone data'
2017-02-14 14:04:55,510 DEBUG [link/execute_link]: Now executing link 'Transform'
2017-02-14 14:04:55,511 DEBUG [applyfunctodataframe/execute]: Applying function
↳<function <lambda> at 0x7faa8ba2e158>
2017-02-14 14:04:55,512 DEBUG [applyfunctodataframe/execute]: Applying function
↳<function <lambda> at 0x7faa8ba95f28>
2017-02-14 14:04:55,515 DEBUG [link/execute_link]: Done executing link 'Transform'
2017-02-14 14:04:55,516 DEBUG [chain/execute]: Done executing chain 'Data'
2017-02-14 14:04:55,516 DEBUG [chain/finalize]: Now finalizing chain 'Data'
2017-02-14 14:04:55,517 DEBUG [link/finalize_link]: Now finalizing link 'LA ozone data
↳'
2017-02-14 14:04:55,518 DEBUG [link/finalize_link]: Done finalizing link 'LA ozone_
↳data'
2017-02-14 14:04:55,518 DEBUG [link/finalize_link]: Now finalizing link 'Transform'
2017-02-14 14:04:55,519 DEBUG [link/finalize_link]: Done finalizing link 'Transform'
2017-02-14 14:04:55,519 DEBUG [chain/finalize]: Done finalizing chain 'Data'

```

with a lot more text surrounding this output. Now try to run the macro again. The run should fail, and you get the following error:

```

KeyError: Processor "<Chain name=Data parent=<... ProcessManager ...> id=...>"
↳already exists!'

```

This is because the `ProcessManager` is a singleton. This means there is only one of this in memory allowed, and since the Jupyter python kernel was still running the object still existed and running the macro gave an error. The macro tried to add a chain, but it already exists in the `ProcessManager`. Therefore the final line in the notebook template has to be ran every time you want to rerun Eskapade. So run this line:

```
execution.reset_eskapade()
```

And try to rerun the notebook without restarting the kernel.

```
execution.eskapade_run(settings)
```

If one wants to call the objects used in the run, `execute` contains them. For example calling

```
ds = process_manager.service(DataStore)
```

is the `DataStore`, and similarly the other ‘master’ objects can be called. Resetting will clear the process manager singleton from memory, and now the macro can be rerun without any errors.

Note: restarting the Jupyter kernel also works, but might take more time because you have to re-execute all of the necessary code.

Reading data from a pickle

Continuing with the notebook we are going to load a pickle file that is automatically written away when the engine runs. First we must locate the folder where it is saved. By default this is in:

```
./results/$MACRO/proc_service_data/v$VERSION/latest/eskapade.core.process_services.
↳DataStore.pkl'
```

Where `$MACRO` is the macro name you specified in the settings, `$VERSION` is the version you specified and `latest` refers to the last chain you wrote to disk. By default, the version is 0 and the name is `v0` and the chain is the last chain of your macro.

You can write a specific chain with the [command line arguments](#), otherwise use the default, the last chain of the macro.

Now we are going to load the pickle from `tutorial_1`.

So make a new cell in Jupyter and add:

```
from eskapade import DataStore
```

to import the `DataStore` module. Now to import the actual pickle and convert it back to the `DataStore` do:

```
ds = DataStore.import_from_file('./results/Tutorial_1/proc_service_data/v0/latest/  
↳eskapade.core.process_services.DataStore.pkl')
```

to open the saved `DataStore` into variable `ds`. Now we can call the keys of the `DataStore` with

```
ds.Print()
```

We see there are two keys: `data` and `transformed_data`. Call one of them and see what is in there. You will find of course the pandas `DataFrames` that we used in the tutorial. Now you can use them in the notebook environment and directly interact with the objects without running the entirety of `Eskapade`.

Similarly you can open old `ConfigObject` objects if they are available. By importing and calling:

```
from eskapade import ConfigObject  
settings = ConfigObject.import_from_file('./results/Tutorial_1/proc_service_data/v0/  
↳latest/eskapade.core.process_services.ConfigObject.pkl')
```

one can import the saved singleton at the path. The singleton can be any of the above mentioned stores/objects. Finally, by default there are soft-links in the results directory at `results/$MACRO/proc_service_data/$VERSION/latest/` that point to the pickles of the associated objects from the last chain in the macro.

6.3.5 Tutorial 4: creating a new analysis project

Now that we have covered how to make a link, macro, and notebook we can create a new analysis project. To generate a new project type the following:

```
$ eskapade_bootstrap --project_root_dir ./yourproject -m yourmacro -l YourLink --n_  
↳yournotebook yourpackage
```

The script will create a Python package called `yourpackage` in the path specified in the `--project_root_dir` argument. The arguments `-m`, `-l`, and `-n` are optional, if not specified the default values are used.

The generated project has the following structure:

```
|yourproject  
  |yourpackage  
    |-links  
      |-__init__.py  
      |-yourlink.py  
    |-__init__.py  
    |-yourmacro.py  
    |yournotebook.ipynb  
  |-setup.py
```

The project contains a link module called `yourlink.py` under `links` directory, a macro `yourmacro.py`, and a Jupyter notebook `yournotebook.ipynb` with required dependencies. To add more of each to the project you can use the commands `generate_link`, `generate_macro`, and `generate_notebook` like it was done before.

The script also generates `setup.py` file and the package can be installed as any other pip package.

Let's try to debug the project interactively within a Jupyter notebook. First, go to your project directory and install the package in an editable mode:

```
$ cd yourproject
$ pip install -e .
```

As you can see in the output, installation checks if `eskapade` and its requirements are installed.

If the installation was successful, run Jupyter and open `yournotebook.ipynb` in `yourpackage` directory:

```
$ jupyter notebook
```

As you see in the code the notebook runs `yourmacro.py`:

```
settings['macro'] = '<...>/yourproject/yourpackage/yourmacro.py'
```

Now run the cells in the notebook and check if the macro runs properly.

6.3.6 Tutorial 5: Data Mimic

This section contains materials on how to use Eskapade in to re-simulate data. We will explain how to use the links present in the submodule and what choice were made.

Running the macro

To run the tutorial macro enter the command in a shell:

```
$ eskapade_run $TUTDIR/esk701_mimic_data.py
```

The tutorial macro illustrates how to resample a dataset using kernel density estimation (KDE). The macro can handle continuous data, and both ordered and unordered categorical data. The macro is build up in the following way:

- A dataset is simulated containing mixed data types, representing general input data.
- Some cleaning steps are performed on the dataset
- KDE is applied to the dataset
- Using the estimated bandwidths of the KDE, the data is resampled
- An evaluation is done on the resulting resimulated dataset

We'll go through each of the links and explain the workings, choices made and available options so to help facilitate the use after a client engagement.

To run the macro for a client engagement you need to change out the first link, which simulates fake data, to the ReadToDf link in order to read in the data:

```
settings['input_path'] = 'diamonds.csv'
settings['reader'] = 'csv'
settings['separator'] = ','
settings['maps'] = maps

np.random.seed(42)

ch = Chain('do_all')
```

(continues on next page)

(continued from previous page)

```

reader = analysis.ReadToDf(name='reader',
                           path=settings['input_path'],
                           sep=settings['separator'],
                           key='df',
                           reader=settings['reader'])

ch.add(reader)

```

For ordered variables that are strings, it is important that you provide the `settings['maps']` variable. It contains a dictionary mapping the string values of the variable to a numeric value that matches the order from low to high. For example, a variable ordered containing 'low', 'medium' and 'high' values should be mapped to 0, 1 and 3. If the mapping is not included the macros will assign numeric values but in order of first appearance, thus not guaranteeing the implied order of the strings. (aka: the macro doesn't know 'low' is the lowest, or 'V1' is the lowest value in your category.)

Warning: When providing the map for ordered categorical variables, you also need to add them for the unordered categorical variables which are strings. The macro will either create maps for all string variables, or will use the provided maps input.

Additionally, you need to provide the macro with the data type contained in each column, ordered-, and unordered categorical, or continuous. Finally also provide the columns which contain strings, so the macro will transfer them to numerical values. In case of ordered categorical values this will happen using the `maps` you mentioned above. If it is not provided, the macro will automatically create a map and add it to the datastore. The automatically created map does not take any implicit order of categories into account.

```

settings['unordered_categorical_columns'] = ['a', 'b']
settings['ordered_categorical_columns'] = ['d', 'c']
settings['continuous_columns'] = ['e', 'f', 'g']
settings['string_columns'] = ['a', 'c']

```

The rest of the macro can be run as is as far as functionality goes. There are a few parameters that can be tweaked to optimize the results:

- `ds['bw']` in `Resampler` is a list of bandwidths estimated by the kde corresponding to the columns of the data as stored in `ds`.
- `bins` in `ResampleEvaluation` controls the bins used to bin the data for chi square evaluation. For more details on the impact of choosing these bins please refer to [Chi-square](#)

Mixed Variables Simulation

The first link exists to create some fake data to run the tutorial macro on. Naturally, if you want to run the `data_mimic` to resimulate a specific dataset, this link is not necessary.

The link takes some parameters as input that will determine how the data is generated, such as the number of observations, the probabilities associated with each category per dimension, and the mean stds for the continuous data.

KDE Preparation

In order to do Kernel Density Estimation we need to prepare the data. This means:

- Converting strings to dummy integer variables and saving the mapping

- Remove any NaN values present
- Find peaks in distributions
- Transforms continuous variables to copula space (transform to a normal distribution)
- Performs a PCA transformation (optional)
- Hash column names and variables if required

Each of these transformations are saved to the datastore including their properties needed to transform the data back.

Kernel Density Estimation

This link performs the actual KDE. The link takes the normalized data in case of the continuous variables, and the data without NaN values for categorical variables. It then applies the `KDEMultivariate` implementation from `statsmodels`, using the 'normal-rule-of-thumb' for quick calculation.

Note: There is a faster method available if only unordered variables are present. This method is selected automatically.

The result is a bandwidth bw for each variable, which are saved in a list in the datastore.

Resampler

Currently the resampler loops over each datapoint and variable j and resamples by:

- Resamples a new point from a normal distribution centered at the original datapoint, with `std=bw[j]`, for continuous variables.
- Resamples randomly from all categories != current category if `bw[j] > np.random.rand()` for unordered categorical variables.
- Resamples using a Wang-Ryzin kernel defined at the datapoint using bandwidth `bw[j]` for ordered categorical variables.

ResampleEvaluation

Evaluates the distribution similarities based on Chi-square, Kolmogorov-Smirnov and the Cosine distance. The formulas and applications of these metrics to the dataset are explained below.

Chi-square

When applying the two sample chi-square test we are testing whether two datasets come from a common distribution.

- H_0 : The two sets come from a common distribution
- H_1 : H_0 is false, the sets come from different distributions

The Chi-square test we use is calculated using the formula:

$$\chi^2 = \sum_{i=1}^k \frac{(K_R * R - K_E * E)^2}{R}$$

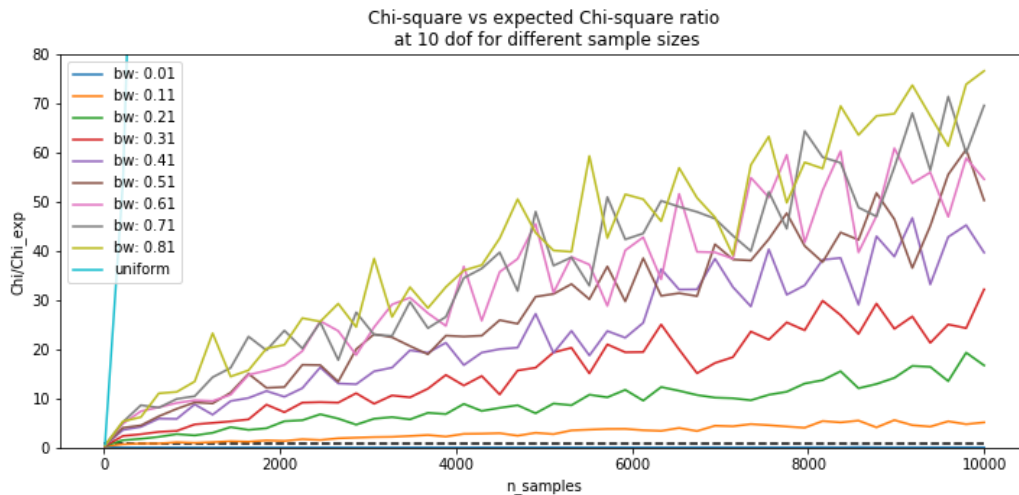
where R is the resampled dataset and E the expected values, or in our context, the original dataset.

$$K_R = \sqrt{\frac{\sum(E_i)}{\sum(R_i)}}$$

$$K_E = \sqrt{\frac{\sum(R_i)}{\sum(E_i)}}$$

In case the datasets are not of equal sample size, we can still use the Chi-square test using the scaling constants. If the sets of equal sample size, the constants will go to 1, and we are left with the ‘traditional’ chi-square formula:

$$\chi^2 = \sum_{i=1}^k \frac{(R - E)^2}{E}$$



Kolmogorov-Smirnov

The Kolmogorov–Smirnov test may also be used to test whether two underlying one-dimensional probability distributions differ. In this case, we apply the KS test to each variable and save the results.

Note: The KS test assumes the distribution is continuous. Although the test is run for all variables, we should keep this in mind when looking at the results for the categorical variables.

If the K-S statistic is small or the p-value is high, then we cannot reject the null-hypothesis that the distributions of the two samples are the same. Aka: They are sufficiently similar to say they are from the same distribution.

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I_{[-\text{inf}, x]}(X_i)$$

$$D_n = \sup_n |F_n(x) - F(x)|$$

Cosine-distance

We also tried to define a distance from the original dataset to the resampled one. We employ the cosine distance applied to each point and its resampled point, represented as a vector. The distance will be 0 when on top of each other, and the max distance is.

We define a vector as the combination of all variables for one datapoint (or row in your dataset).

- All continuous values are represented as is
- All ordered categorical values are mapped to numerical values going from 0 to # of categories available, where 0 corresponds to the lowest ranking category.
- All unordered categorical are ignored for now since we have not yet defined a sufficient distance measure for these.

Then, the cosine distance is calculated for each point and it's corresponding resample.

$$\cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

Mimic Report

The mimic report link will create a standard eskapade style pdf report. The report includes per variable:

- A stacked histogram plot showing before and after the resampling
- A stacked histogram plot of the data per variable in the copula space and a normal distribution
- A correlation matrix of numeric values before and after resimulation

Each variable page also contains the chi-square values comparing before and after the resampling (also see *Chi-square*). For each variable, there is a table containing several values. The values correspond the chisquare calculation done on a 1D histogram of the variable itself, and done on 2D histograms of two variables as listed in the table.

Example: On the page of variable *d*

	Chi2	p-value	dof
d	1.22018	0.269325	3
e	1034.82	0	3
f	317.124	0	3
g	1118.11	0	3
a	7.92157	0.0476607	3
b	1.4137	0.84181	3
c	1.43721	0.696837	3

The value 1.22 corresponds to the calculation variable *d* before and after the resampling. The value of 1034.82 corresponds to the calculations done on a 2D histogram of variables *d* and *e*, before and after the resampling.

Finally, two other metrics, the Kolmogorov-Smirnov and the cosine distance, are also calculated and shown in the report. You can find these on the last page.

6.3.7 All available examples

To see the available Eskapade examples, do:

```
$ export TUTDIR=`pip show Eskapade-Core | grep Location | awk '{ print $2"/escore/
↪tutorials" }'`
$ export TUTDIR=`pip show Eskapade | grep Location | awk '{ print $2"/eskapade/
↪tutorials" }'`
$ ls -l $TUTDIR/ $TUTDIR/
```


Many Eskapade example macros exist in the tutorials directory. The numbering of the example macros follows the package structure:

- esk100+: basic macros describing the chains, links, and datastore functionality of Eskapade.
- esk200+: macros describing links to do basic processing of pandas dataframes.
- esk300+: visualization macros for making histograms, plots and reports of datasets.
- esk500+: macros for doing data quality assessment and cleaning.
- esk700+: macros for doing data simulation.

The Eskapade macros are briefly described below. They explain the basic architecture of Eskapade, i.e. how the chains, links, datastore, and process manager interact.

Hopefully you now have enough knowledge to run and explore Eskapade by yourself. You are encouraged to run all examples to see what Eskapade can do for you!

Example esk101: Hello World!

Macro 101 runs the Hello World Link. It runs the Link twice using a repeat kwarg, showing how to use kwargs in Links.

```
$ eskapade_run $TUTDIRC/esk101_helloworld.py
```

Example esk102: Multiple chains

Macro 102 uses multiple chains to print different kinds of output from one Link. This link is initialized multiple times with different kwargs and names. There are if-statements in the macro to control the usage of the chains.

```
$ eskapade_run $TUTDIRC/esk102_multiple_chains.py
```

Example esk103: Print the DataStore

Macro 103 has some objects in the DataStore. The contents of the DataStore are printed in the standard output.

```
$ eskapade_run $TUTDIRC/esk103_printdatastore.py
```

Example esk104: Basic DataStore operations

Macro 104 adds some objects from a dictionary to the DataStore and then moves or deletes some of the items. Next it adds more items and prints some of the objects.

```
$ eskapade_run $TUTDIRC/esk104_basic_datastore_operations.py
```

Example esk105: DataStore Pickling

Macro 105 has 3 versions: A, B and C. These are built on top of the basic macro esk105. Each of these 3 macro's does something slightly different:

- A does not store any output pickles,
- B stores all output pickles,

- C starts at the 3rd chain of the macro.

Using these examples one can see how the way macro's are run can be controlled and what it saves to disk.

```
$ eskapade_run $TUTDIRC/esk105_A_dont_store_results.py
$ eskapade_run $TUTDIRC/esk105_B_store_each_chain.py
$ eskapade_run $TUTDIRC/esk105_C_begin_at_chain3.py
```

Example esk106: Command line arguments

Macro 106 shows us how command line arguments can be used to control the chains in a macro. By adding the arguments from the message inside of the macro we can see that the chains are not run.

```
$ eskapade_run $TUTDIRC/esk106_cmdline_options.py
```

Example esk107: Chain loop

Example 107 adds a chain to the macro and using a repeater Link it repeats the chain 10 times in a row.

```
$ eskapade_run $TUTDIRC/esk107_chain_looper.py
```

Example esk108: Event loop

Example 108 processes a textual data set, to loop through every word and do a Map and Reduce operation on the data set. Finally a line printer prints out the result.

```
$ source $TUTDIRC/esk108_eventlooper.sh
```

Example esk109: Debugging tips

This macro illustrates basic debugging features of Eskapade. The macro shows how to start a python session while running through the chains, and also how to break out of a chain.

```
$ eskapade_run $TUTDIRC/esk109_debugging_tips.py
```

Example esk110: Code profiling

This macro demonstrates how to run Eskapade with code profiling turned on.

```
$ eskapade_run $TUTDIRC/esk110_code_profiling.py
```

Example esk111: Loading a datastore from file

Macro illustrates how to load an external datastore from file.

```
$ eskapade_run $TUTDIRC/esk111_load_datastore_from_file.py
```

Example esk201: Read data

Macro 201 illustrates how to open files as pandas datasets. It reads a file into the DataStore. The first chain reads one csv into the DataStore, the second chain reads multiple files (actually the same file multiple times) into the DataStore. (Looping over data is shown in example esk209.)

```
$ eskapade_run $TUTDIR/esk201_readdata.py
```

Example esk202: Write data

Macro 202 illustrate writing pandas dataframes to file. It reads a DataFrame into the data store and then writes the DataFrame to csv format on the disk.

```
$ eskapade_run $TUTDIR/esk202_writedata.py
```

Example esk203: apply func to pandas df

Illustrates the link that calls basic apply() to columns of a pandas dataframes. See for more information pandas documentation:

<http://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.apply.html>

```
$ eskapade_run $TUTDIR/esk203_apply_func_to_pandas_df.py
```

Example esk204: apply query to pandas df

Illustrates the link that applies basic queries to pandas dataframe. See for more information pandas documentation:

<http://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.query.html>

```
$ eskapade_run $TUTDIR/esk204_apply_query_to_pandas_df.py
```

Example esk205: concatenate pandas dfs

Illustrates the link that calls basic concat() of pandas dataframes. See for more information pandas documentation:

<http://pandas.pydata.org/pandas-docs/stable/merging.html>

```
$ eskapade_run $TUTDIR/esk205_concatenate_pandas_dfs.py
```

Example esk206: merge pandas dfs

Illustrate link that calls basic merge() of pandas dataframes. For more information see pandas documentation:

<http://pandas.pydata.org/pandas-docs/stable/merging.html>

```
$ eskapade_run $TUTDIR/esk206_merge_pandas_dfs.py
```

Example esk207: record vectorizer

This macro performs the vectorization of an input column of an input dataframe. E.g. a column x with values 1, 2 is transformed into columns x_1 and x_2, with values True or False assigned per record.

```
$ eskapade_run $TUTDIR/esk207_record_vectorizer.py
```

Example esk208: record factorizer

This macro performs the factorization of an input column of an input dataframe. E.g. a column x with values 'apple', 'tree', 'pear', 'apple', 'pear' is transformed into columns x with values 0, 1, 2, 0, 2.

```
$ eskapade_run $TUTDIR/esk208_record_factorizer.py
```

Example esk209: read big data itr

Macro to that illustrates how to loop over multiple (possibly large!) datasets in chunks.

```
$ eskapade_run $TUTDIR/esk209_read_big_data_itr.py
```

Example esk210: dataframe restoration

Macro to illustrate writing pandas dataframes to file and reading them back in whilst retaining the datatypes and index using numpy and feather file formats.

```
$ eskapade_run $TUTDIR/esk210_dataframe_restoration.py
```

Example esk301: dfsummary plotter

Macro shows how to plot the content of a dataframe in a nice summary pdf file. (Example similar to \$TUTDIR/tutorial_1.py.)

```
$ eskapade_run $TUTDIR/esk301_dfsummary_plotter.py
```

Example esk302: histogram_filler_plotter

Macro that illustrates how to loop over multiple (possibly large!) datasets in chunks, in each loop fill a (common) histogram, and plot the final histogram.

```
$ eskapade_run $TUTDIR/esk302_histogram_filler_plotter.py
```

Example esk303: histogrammar filler plotter

Macro that illustrates how to loop over multiple (possibly large!) datasets in chunks, in each loop fill a histogrammar histograms, and plot the final histograms.

```
$ eskapade_run $TUTDIR/esk302_histogram_filler_plotter.py
```

Example esk304: df boxplot

Macro shows how to boxplot the content of a dataframe in a nice summary pdf file.

```
$ eskapade_run $TUTDIR/esk304_df_boxplot.py
```

Example esk305: correlation summary

Macro to demonstrate generating nice correlation heatmaps using various types of correlation coefficients.

```
$ eskapade_run $TUTDIR/esk305_correlation_summary.py
```

Example esk306: concatenate reports

This macro illustrates how to concatenate the reports of several visualization links into one big report.

```
$ eskapade_run $TUTDIR/esk306_concatenate_reports.py
```

Example esk501: fix pandas dataframe

Macro illustrates how to call FixPandasDataFrame link that gives columns consistent names and datatypes. Default settings perform the following clean-up steps on an input dataframe:

- Fix all column names. Eg. remove punctuation and strange characters, and convert spaces to underscores.
- Check for various possible nans in the dataset, then make all nans consistent by turning them into numpy.nan (= float)
- Per column, assess dynamically the most consistent datatype (ignoring all nans in that column). Eg. bool, int, float, datetime64, string.
- Per column, make the data types of all rows consistent, by using the identified (or imposed) data type (by default ignoring all nans)

```
$ eskapade_run $TUTDIR/esk501_fix_pandas_dataframe.py
```

Example esk701: Mimic dataset

Macro that illustrates how to resample a dataset using kernel density estimation (KDE). The macro can handle continuous data, and both ordered and unordered categorical data. The macro is build up in the following way:

- A dataset is simulated containing mixed data types, representing general input data.
- Some cleaning steps are performed on the dataset
- KDE is applied to the dataset
- Using the estimated bandwidths of the KDE, the data is resampled
- An evaluation is done on the resulting resimulated dataset

```
$ eskapade_run $TUTDIR/esk701_mimic_data.py
```

Example esk702: Mimic data only unordered

This macro illustrates how to resample an existing data set, containing only unordered categorical data, using kernel density estimation (KDE) and a direct resampling technique.

```
$ eskapade_run $TUTDIR/esk702_mimic_data_only_unordered.py
```

6.3.8 Tips on coding

This section contains a general description on how to use Eskapade in combination with other tools, in particular for the purpose of developing code.

Eskapade in PyCharm

PyCharm is a very handy IDE for debugging Python source code. It can be used to run Eskapade stand-alone (i.e. like from the command line) and with an API.

Stand-alone

- Install PyCharm on your machine.
- Open project and point to the Eskapade source code
- **Configuration, in ‘Preferences’, check the following desired values:**
 - Under ‘Project: eskapade’ / ‘Project Interpreter’:
 - * The correct Python version (Python 3)
 - Under ‘Build, Execution & Deployment’ / ‘Console’ / ‘Python Console’:
 - * The correct Python version (Python 3)
- Install Eskapade in editable mode
- **Run/Debug Configuration:**
 - Under ‘Python’ add new configuration
 - Script: path to the console script `eskapade_run` (located in the same directory as the interpreter specified above in ‘Project Interpreter’)
 - Script parameters: path to a macro to be debugged, e.g. `$ESKAPADE/python/eskapade/tutorials/tutorial_1.py`, and `eskapade_run` command line arguments, e.g. `--begin-with=Summary`
 - Python interpreter: check if it is the correct Python version (Python 3)

You should now be able to press the ‘play button’ to run Eskapade with the specified parameters.

Writing a new Link using Jupyter and notebooks

Running the framework works best from the command line (in our experience), but running experiments and trying new ideas is better left to an interactive environment like Jupyter. How can we reconcile the difference in these work flows? How can we use them together to get the most out of it?

Well, when using the [data and config import functionality](#) of Eskapade together with Jupyter we can interactively work on our objects and when we are satisfied with the results integration into links is straight-forward. The steps to undertake this are *in general* the following:

1. Import the DataStore and/or ConfigObject. Once you have imported the ConfigObject, run it to generate the output you want to use.
2. Grab the data you want from the DataStore using `ds = DataStore.import_from_file` and `data = ds[key]`.
3. Now you can apply the operation you want to do on the data, experiment on it and work towards the end result you want to have.
4. Create a new link in the appropriate link folder using the `eskapade_generate_link` command.
5. Copy the operations (code) you want to do to the link.
6. Add assertions and checks to make sure the Link is safe to run.
7. Add the Link to your macro and run it!

These steps are very general and we will now go into steps 5, 6 and 7. Steps 1, 2, 3 and 4 have already been covered by various parts of the documentation.

So assuming you wrote some new functionality that you want to add to a Link called YourLink and you have created a new Link from the template we are going to describe how you can get your code into the Link and test it.

Developing Links in notebooks

This subsection starts with a short summary of the workflow for developing Links:

1. Make your code in a notebook
2. Make a new Link
3. Port the code into the Link
4. Import the Link into your notebook
5. Test if the Link has the desired effect.
6. Finish the Link code
7. Write a unit test (optional but advised if you want to contribute)

We continue with a longer description of the steps above.

When adding the new code to a new link the following conventions are used:

In the `__init__` you specify the key word arguments of the Link and their default values, if you want to get an object from the DataStore or you want to write an object back into it, use the name `read_key` and `store_key`. Other keywords are free to use as you see fit.

In the `initialize` function in the Link you define and initialize functions that you want to call when executing the code on your objects. If you want to import something, you can do this at the root of the Link, as per PEP8.

In the `execute` function you put the actual code in this format:

```
settings = process_manager.service(ConfigObject)
ds = process_manager.service(DataStore)

## --- your code follows here
```

Now you can call the objects that contain all the settings and data of the macro in your Link, and in the code below you can add your analysis code that calls from the objects and writes back in case that this is necessary. Another possibility is writing a file to the disk, for example writing out a plot you made.

If you quickly want to test the Link without running the entire Eskapade framework, you can import it into your notebook sessions:

```
import eskapade.analysis.links.yourlink
from yourlink import YourLink
l = YourLink()
l.execute()
```

should run your link. You can also call the other functions. However, `execute()` is supposed to contain the bulk of your operations, so running that should give you your result. Now you can change the code in your link if it is not how you want it to run. The notebook kernel however keeps everything in memory, so you either have to restart the kernel, or use

```
import importlib
importlib.reload(eskapade.analysis.links.yourlink)
from yourlink import YourLink
l = YourLink()
l.execute()
```

to reload the link you changed. This is equivalent to the Python2 function `reload(eskapade)`.

Combined with the importing of the other objects it becomes clear that you can run every piece of the framework from a notebook. However working like this is only recommended for development purposes, running an entire analysis should be done from the command line.

Finally after finishing all the steps you use the function `finalize()` to clean up all objects you do not want to save.

After testing whether the Link gives the desired result you have to add the proper assertions and other types of checks into your Link code to make sure that it does not have use-cases that are improperly defined. It is advised that you also write a unit test for the Link, but unless you want it merged into the master, it will not be enforced.

Now you can run Eskapade with your macro from your command line, using the new functionality that you first created in a notebook and then ported into a stand-alone Link.

6.4 Command Line Arguments

6.4.1 Overview

We start this section with a short overview of a few often used arguments of the Eskapade command `eskapade_run`. The only required argument is a configuration file, which can be a Python script (Eskapade macro) or a pickled Eskapade configuration object. This section gives an overview of the optional arguments of the run command.

At the end of running the Eskapade program, by default the DataStore and configuration object are pickled and written out to:

```
$ ls -l results/Tutorial_1/proc_service_data/v0/latest/
```

When you are working on a macro, once you are done tweaking it, you can also store the results of each chain in pickle files:

```
$ eskapade_run --store-all python/eskapade/tutorials/tutorial_1.py
```

Eskapade uses these pickle files to load the trained models and uses them to predict new samples real-time, but also to pick up running at a later stage in the chain setup.

For example, if running Eskapade takes a long time, you can run one chain as well:


```
$ eskapade_run --single-chain=Data python/eskapade/tutorials/tutorial_1.py
```

This command uses as input the stored pickle files from the previous chain. This might come in handy when, for example, data pre-processing of your training set takes a long time. In that case, you can run the pre-processing chain over night, store the results in a pickle file and start with the training chain the next day.

Start running Eskapade from a specified chain:

```
$ eskapade_run --begin-with=Summary python/eskapade/tutorials/tutorial_1.py
```

Stop running after a specified chain:

```
$ eskapade_run --end-with=Data python/eskapade/tutorials/tutorial_1.py
```

Below the most important command-line options are explained in detail.

6.4.2 Table of all arguments

The following table summarizes the available command-line options. Most of these options set variables in the Eskapade configuration object and can be overwritten by settings in the configuration macro.

Option	Short option	Argument	Description
-help	-h		show help message and exit
-analysis-name	-n	NAME	set name of analysis in run
-analysis-version	-v	VERSION	set version of analysis version in run
-batch-mode			run in batch mode (no X Windows)
-interactive	-i		start Python shell after run
-log-level	-L	LEVEL	set logging level
-log-format		FORMAT	set log-message format
-unpickle-config			interpret first CONFIG_FILE as path to pickled settings
-profile			run profiler for Python code
-conf-var	-c	KEY=VALUE	set configuration variable
-begin-with	-b	CHAIN_NAME	begin execution with chain CHAIN_NAME
-end-with	-e	CHAIN_NAME	end execution with chain CHAIN_NAME
-single-chain	-s	CHAIN_NAME	only execute chain CHAIN_NAME
-store-all			store run-process services after every chain
-store-one		CHAIN_NAME	store run-process services after chain CHAIN_NAME
-store-none			do not store run-process services
-results-dir		RESULTS_DIR	set directory path for results output
-data-dir		DATA_DIR	set directory path for data
-macros-dir		MACROS_DIR	set directory path for macros
-templates-dir		TEMPLATES_DIR	set directory path for template files
-spark-cfg-file		SPARK_CONFIG_FILE	set path of Spark configuration file
-seed		SEED	set seed for random-number generation

6.4.3 Description and examples

This section contains the most used options with a longer description of what it does and how it works combined with examples.

Set log level

The log level is controlled with the `--log-level` option. For example, to set the log level to “debug”, add:

```
--log-level=DEBUG
```

to the command line:

```
$ eskapade_run -L DEBUG python/eskapade/tutorials/tutorial_1.py
```

The available log levels are:

```
NOTSET,  
DEBUG,  
INFO,  
WARNING,  
ERROR,  
FATAL
```

They correspond to the appropriate POSIX levels.

When writing your own Link, these levels can be accessed with the logger module:

```
self.logger.debug('Text to be printed when logging at DEBUG level')
```

All output is done in this manner, never with the python print function, since this yields us more control over the process.

Help

Help can be called by running the following:

```
$ eskapade_run --help
```

Interactive python mode

To keep the results in memory at end of session and access them in an interactive session, run Eskapade in interactive mode. This is controlled with `--interactive`:

```
$ eskapade_run -i python/eskapade/tutorials/tutorial_1.py
```

At the end of the session a Python console is started from which e.g. the data store can be accessed.

Saving states

To write out the intermediate results from every chain, add the command line argument `--store-all`. This will write pickles in `results/NAME/proc_service_data/VERSION/`, containing the state of Eskapade at the end of the chain:

```
$ eskapade_run --store-all python/eskapade/tutorials/tutorial_1.py
```

To write out the state after a particular chain, use option `--store-one`:

```
$ eskapade_run --store-one=Data python/eskapade/tutorials/tutorial_1.py
```

To not store any pickle files, run with the option `--store-none`:

```
$ eskapade_run --store-none python/eskapade/tutorials/tutorial_1.py
```

Single Chain

To run a single chain, use the option `--single-chain`. This picks up the data stored by the previous chain in the macro. It is, therefore, necessary to have run the previous chain, otherwise the engine can not start:

```
$ eskapade_run -s Summary python/eskapade/tutorials/tutorial_1.py
```

Start from a Chain

To start from a chain use the command line argument `--begin-with`. This picks up the data stored by the previous chain in the macro.

```
$ eskapade_run -b Summary python/eskapade/tutorials/tutorial_1.py
```

Stop at a Chain

To end the running of the engine at a chain use, the command line argument `--end-with`:

```
$ eskapade_run -e Data python/eskapade/tutorials/tutorial_1.py
```

Changing analysis version

A version number is assigned to each analysis, which by default is 0. It can be upgraded by using the option `--analysis-version`. When working on an analysis, it is recommended to update this number regularly for bookkeeping purposes. The command line always has higher priority over the macro. If the macro is version 0 and the command line uses version 1, the command line will overrule the macro.

```
$ eskapade_run -v 1 python/eskapade/tutorials/tutorial_1.py
```

Notice that the output of this analysis is now stored in the directory:

```
$ ls -l results/Tutorial_1/data/v1/report/
```

Notice as well that, for bookkeeping purposes, a copy of the (evolving) configuration macro is always stored as well, under:

```
$ ls -l results/Tutorial_1/config/v1/tutorial_1.py
```

Running an old configuration (macro)

Settings for the Eskapade run are stored in a configuration object, which is accessed as a run-process service. This run-time service can be persisted as a file, which is normally done at the end of the run.

Persisted settings can be used in a following run by providing the file path of the `ConfigObject` pickle file as the configuration file argument. The option `--unpickle-config` is required to indicate that this file contains persisted settings:

```
$ eskapade_run --unpickle-config results/Tutorial_1/proc_service_data/v0/latest/  
↳eskapade.core.process_services.ConfigObject.pkl
```

In this way, rolling back to a previous point is straight-forward.

For lookup purposes a copy of the configuration macro is always stored as well, under:

```
$ ls -l results/Tutorial_1/config/v0/tutorial_1.py
```

Profiling your code

You can profile the execution of your analysis functions with the option `--profile`:

```
$ eskapade_run --profile=cumulative python/eskapade/tutorials/tutorial_1.py
```

After running this prints out a long list of all functions called, including the time it took to run each of them, where the functions are sorted based on cumulative time.

To get the the list of sorting options for the profiling, run:

```
$ eskapade_run --help
```

Combining arguments

Of course you can add multiple arguments to the command line, the result would be for example an interactive session in debug mode that writes out intermediate results from each chain:

```
$ eskapade_run -i --store-all -L DEBUG -c do_chain0=False -c mydict="{ 'f': 'y=pi', 'pi  
↳': 3.14}" python/eskapade/tutorials/esk106_cmdline_options.py
```

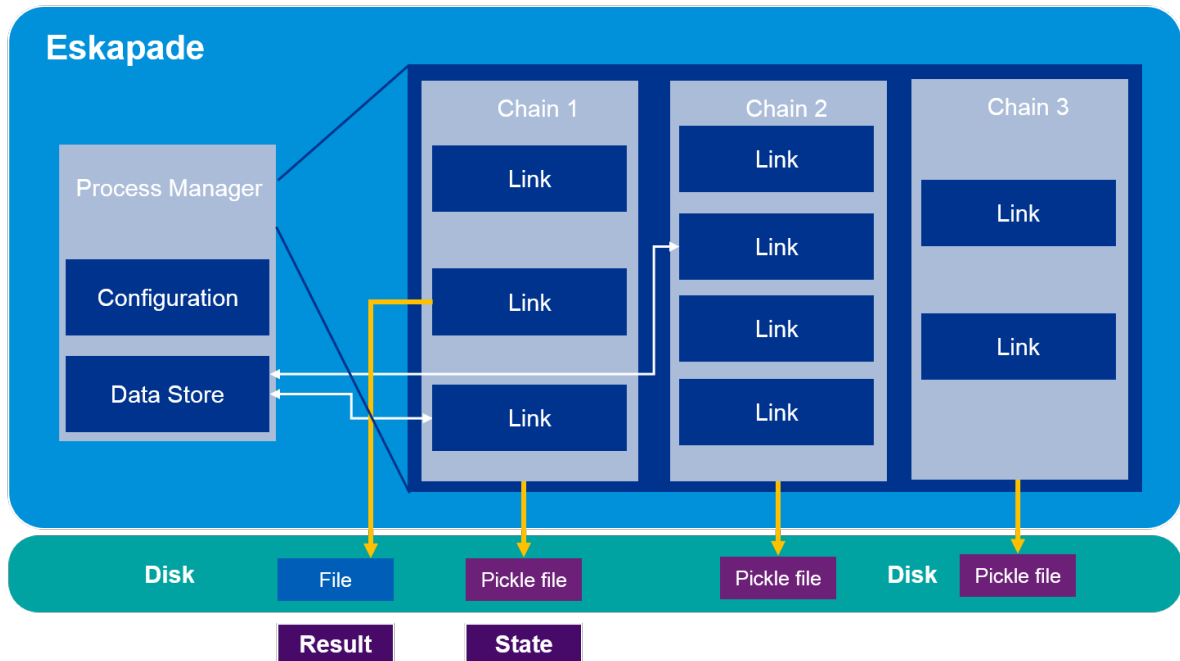
6.5 Package structure

Eskapade contains many tools, and to find and use them most efficiently it is necessary to understand how the repository is build up. This section discusses the structure of the code and how the framework handles subpackages.

6.5.1 Architecture

The architecture of Eskapade can be summarized in this picture:

Architecture of Eskapade



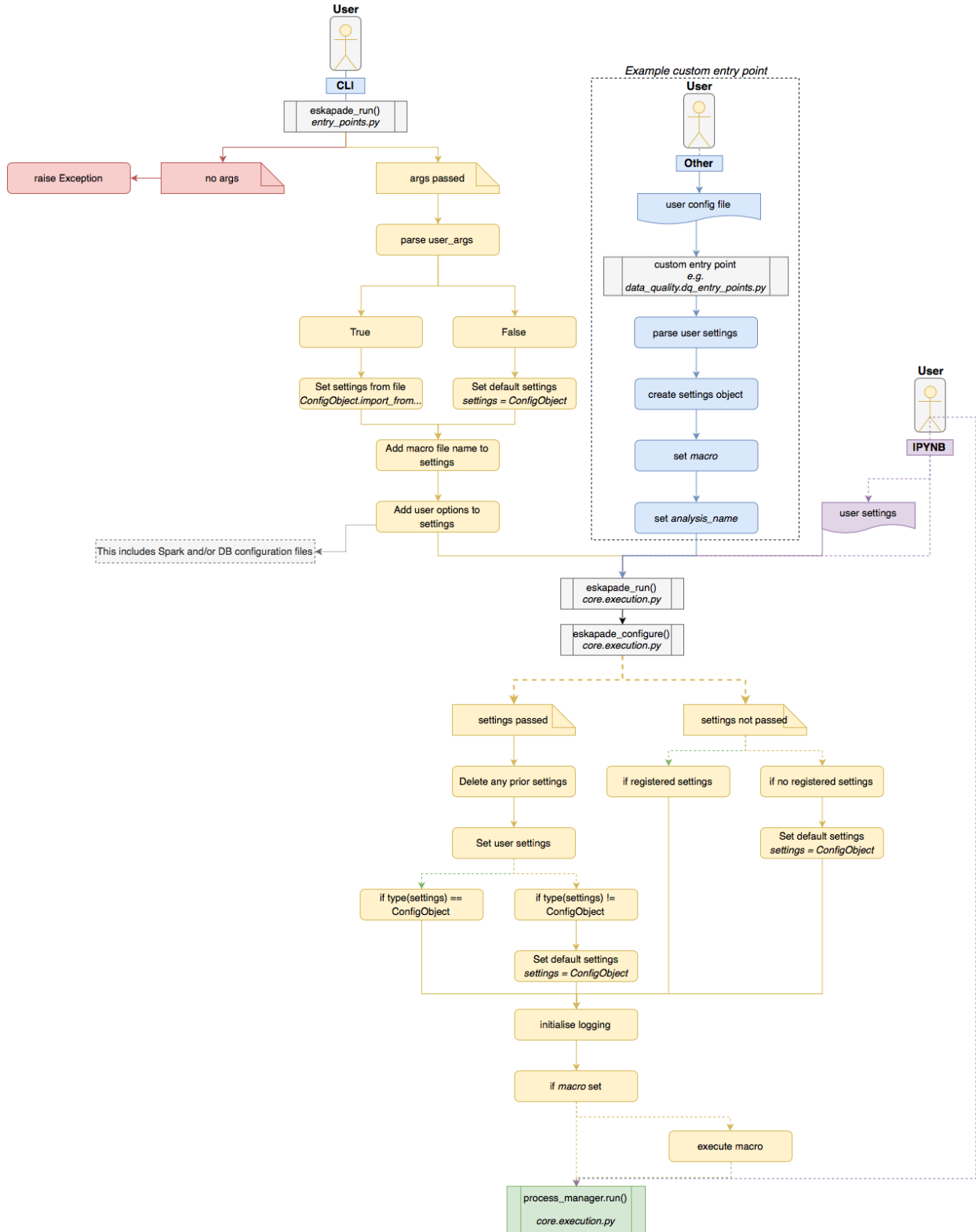
The example we just discussed generally shows how the framework works. The steps it takes are the following:

- `eskapade_run` configures Eskapade based on user settings and/or the macro file, and executes all chains and links;
- Macros (python file) contain Chains;
- Chains (python object) contains Links;
- Links (python class) contain analysis code.

The chains are run in the order of 'registering' them in the `ProcessManager`.

The `ProcessManager` is the ultimate object that executes all the code in your macro. It also keeps track of the configuration of Eskapade, and of the objects in the `data store` that are passable between links.

The settings flow of Eskapade is shown in the following picture:



The components of the architecture of Eskapade are explained in further detail in the [Tutorials section](#).

6.5.2 Structure

When using Eskapade it is important to understand where all components are located. The components can be for example links or utilities that you want to use.

The Eskapade framework is contained in the Python package `eskapade`, which lives in the `python` directory. Every specific subject has its subpackage in `eskapade`, containing the utilities it needs, the links that are defined for the subject.

The core of the framework is implemented in the `core` subpackage. This subpackage contains the low-level machinery for running analysis algorithms in chains of links. The `core_ops` subpackage contains basic links for operating this framework.

An example of a typical subpackage is `eskapade.analysis`, which contains basic analysis tools. Its structure is common to all Eskapade subpackages:

```
|-eskapade
  |-analysis
    |-links
```

The subpackage contains several modules, which contain classes and functions to be applied in links. The `eskapade.analysis.statistics` module, for example, contains code to generate an overview of the statistical properties of variables in given input data.

Eskapade links are located in the `links` directory. There is a separate module for each link, defining the link class instance. By convention, the names of the module and class are both the link name, the former in snake case and the latter in camel case. For example, the module `read_to_df` defines the link class `ReadToDf`.

The tests are contained separately in the Python package `eskapade_python` under `tests` directory. Ideally, there is a test module for each (link) module in the Eskapade package. Optionally, integration tests are implemented in `integration`. For the `eskapade.analysis` package, there is the module `test_tutorial_macros` with integration tests that run the tutorial macros corresponding to this subpackage:

```
|-eskapade_python
  |-analysis
    |-integration
      |-test_tutorial_macros.py
```

6.5.3 Subpackages

Eskapade contains the following list of subpackages:

- `core` is the package that contains the core framework of Eskapade.
- `core_ops` contains links pertaining to the core functionality of Eskapade.
- `analysis` contains pandas links and code.
- `visualization` contains visualization code and plotter links.
- `data_quality` contains links and code for fixing messy data.
- `data_mimic` contains links and code for re-simulating data with mixed data-types

6.5.4 Imports

Main elements of the Eskapade framework are imported directly from the `eskapade` package. For example, the run-configuration object and the run-process manager are part of the `core` subpackage, but are imported by

```
from eskapade import process_manager, ConfigObject
```

Links are imported directly from their subpackage:

```
from eskapade.analysis import ReadToDf
```

In a macro, you can now instantiate and configure the `ReadToDf` link and add it to a chain in the process manager.

6.5.5 Results

Results of a macro are written out by default in the `results` directory. The analysis run is persisted in the results directory by the `analysis_name` given in the macro. This directory has the following structure:

- `config`: the configuration macro
- `proc_service_data`: persisted states of run-process services
- `data`: analysis results, such as graphs or a trained model

The data for each of these elements are stored by the analysis version, e.g. `v0`, `v1`, `v2`, etc. For example, the report produced by the tutorial `esk301_dfsummary_plotter` is saved in the directory `results/esk301_dfsummary_plotter/data/v0/report`.

6.5.6 Debugging

When building new Links or other functionality you will want to debug at some point. There are multiple ways to do this, because there are multiple ways of running the framework. A few ways are:

- Running in the terminal. In this scenario you have to work in a virtual environment (or adjust your own until it has all dependencies) and debug using the terminal output.
- Running in a notebook. This way the code is run in a notebook and you can gather the output from the browser.
- Running in a docker. The code is run in the docker and the repository is mounted into the container. The docker (terminal) returns output.
- Running in a VM. In this case you run the code in the VM and mount the code into the VM. The output can be gathered in the VM and processed in the VM.

In the first three options you want to use an IDE or text-editor in a ‘normal’ environment to debug your code and in the last option you can use an editor in the VM or outside of it.

6.5.7 Troubleshooting

The least error prone ways are docker and VMs, because they automatically have the dependencies set.

6.6 Release notes

6.6.1 Version 0.9.0

Eskapade release v0.9.0 has the following useful upgrades:

- The core functionality of Eskapade, namely: the `Link`, `Chain`, `process_manager`, `DataStore`, `ConfigObject` and corresponding tutorials, have been split off from the growing (ever more analysis related) Eskapade repository, into the new package `Eskapade-Core`. `Eskapade-Core` is a very light-weight Python3 package.
- A new module `data_mimic` has been add to Eskapade, including tutorials, meant for resimulating existing datasets.
- We have added `feather` i/o functionality for reading and writeng dataframes.
- The logger has been fixed, it is now possible to set the log-level of loggers again.
- The Eskapade docker files have been taken out of the Eskapade repository to avoid version conflicts, into the new git repo `Eskapade-Environment`.
- The Eskapade docker image `eskapade-usr` contain the latest working versions of Eskapade, `Eskapade-Core`, `Eskapade-ROOT`, and `Eskapade-Spark`. Type:

```
$ docker pull kave/eskapade-usr:latest
```

to pull it in.

6.6.2 Version 0.8.2

The Eskapade patch release v0.8.2 and corresponding docker containers fix two issues:

- The `matplotlib` backend is no longer set to batchmode when running Eskapade in a jupyter notebook. By default, batch mode is only turned on when no `DISPLAY` environment variable is set, and when not running in a notebook; the batch-mode flag can also be controlled with the command line option `-batch-mode`.
- The Eskapade docker containers contain working version of Eskapade, `Eskapade-ROOT`, and `Eskapade-Spark`. Type:

```
$ docker pull kave/eskapade-usr:0.8.2
```

to pull it in.

6.6.3 Version 0.8

In version 0.8 of Eskapade (August 2018) the modules `root-analysis` and `spark-analysis` have been split off into separate packages called `Eskapade-ROOT` and `Eskapade-Spark` .

So the (core) Eskapade package no longer depends on `ROOT` and `Spark`, just on plain python packages. This make it much easier for people to try out the core functionality of Eskapade.

- To install `Eskapade-ROOT` and `Eskapade-Spark`, do:

```
$ pip install Eskapade-ROOT
$ pip install Eskapade-Spark
```

or check out the code from out github repository:

```
$ git clone https://github.com/KaveIO/Eskapade-ROOT.git eskapade-root
$ pip install -e eskapade-root/
$ git clone https://github.com/KaveIO/Eskapade-Spark.git eskapade-spark
$ pip install -e eskapade-spark/
```

where in this example the code is installed in edit mode (option `-e`).

You can now use these in Python with:

```
import eskapadespark
import esroofit
```

6.6.4 Version 0.7

Version 0.7 of Eskapade (February 2018) contains several major updates:

- The Eskapade code has been made pip friendly. One can now simply do:

```
$ pip install Eskapade
```

or check out the code from our github repository:

```
$ git clone https://github.com/KaveIO/Eskapade.git
$ pip install -e Eskapade/
```

where in this example the code is installed in edit mode (option `-e`).

You can now use Eskapade in Python with:

```
import eskapade
```

This change has resulted in some restructuring of the python directories, making the overall structure more transparent: all python code, including the tutorials, now fall under the (single) `python/` directory. Additionally, thanks to the pip convention, our prior dependence on environment variables (`$ESKAPADE`) has now been fully stripped out of the code.

- There has been a cleanup of the core code, removing obsolete code and making it better maintainable. This has resulted in a (small) change in the api of the process manager, adding chains, and using the logger. All tutorials and example macro files have been updated accordingly. See the [migration section](#) for detailed tips on migrating existing Eskapade code to version 0.7.
- All eskapade commands now start with the prefix `eskapade_`. All tutorials have been updated accordingly. We have the commands:
 - `eskapade_bootstrap`, for creating a new Eskapade analysis project. See this new [tutorial](#) for all the details.
 - `eskapade_run`, for running the Eskapade macros.
 - `eskapade_trail`, for running the Eskapade unit and integration tests.
 - `eskapade_generate_link`, `eskapade_generate_macro`, `eskapade_generate_notebook`, for generating a new link, macro, or Jupyter notebook respectively.

6.6.5 Version 0.6

The primary feature of version 0.6 (August 2017) is the inclusion of Spark, but this version also includes several other new features and analyses.

We include multiple Spark links and 10 Spark examples on:

- The configuration of Spark, reading, writing and converting Spark dataframes, applying functions and queries to dataframes, filling histograms and (very useful!) applying arbitrary functions (e.g. pandas) to groupby calls.

In addition we have added:

- A ROOT analysis for studying and quantifying between sets of (non-)categorical and observables. This is useful for finding outliers in arbitrary datasets (e.g. surveys), and we include a tutorial of how to do this.
- A ROOT analysis on predictive maintenance that decomposes a distribution of time difference between malfunctions by fitting this multiple Weibull distributions.
- New flexible features to create and chain analysis reports with several analysis and visualization links.

6.6.6 Version 0.5

Our 0.5 release (May 2017) contains multiple new features, in particular:

- Support for ROOT, including multiple examples on using data analysis, fitting and simulation examples using RooFit.
- Histogram conversion and filling support, using ROOT, numpy, Histogrammar and Eskapade-internal histograms.
- Automated data-quality fixes for buggy columns datasets, including data type fixing and NaN conversion.
- New visualization utilities, e.g. plotting multiple types of (non-linear) correlation matrices and dendograms.
- And most importantly, many new and interesting example macros illustrating the new features above!

6.6.7 Version 0.4

In our 0.4 release (Feb 2017) we are releasing the core code to run the framework. It is written in python 3. Anyone can use this to learn Eskapade, build data analyses with the link-chain methodology, and start experiencing its advantages.

The focus of the provided documentation is on constructing a data analysis setup in Eskapade. Machine learning interfaces will be included in an upcoming release.

6.7 Developing and Contributing

6.7.1 Preliminaries

6.7.2 Working on Eskapade

You have some cool feature and/or algorithm you want to add to Eskapade. How do you go about it?

First clone Eskapade.

```
git clone https://github.com/KaveIO/Eskapade.git eskapade
```

then

```
pip install -e eskapade
```

this will install Eskapade in editable mode, which will allow you to edit the code and run it as you would with a normal installation of eskapade.

To make sure that everything works try executing eskapade without any arguments, e.g.

```
eskapade_run
```

or you could just execute the tests using either the eskapade test runner, e.g.

```
cd eskapade
eskapade_trial .
```

or

```
cd eskapade
python setup.py test
```

That's it.

6.7.3 Contributing

When contributing to this repository, please first discuss the change you wish to make via issue, email, or any other method with the owners of this repository before making a change. You can find the contact information on the [index](#) page.

Note that when contributing that all tests should succeed.

6.7.4 Tips and Tricks

- Enable auto reload in ipython:

```
%load_ext autoreload
```

this will reload modules before executing any user code.

6.8 References

- Web page: <https://eskapade.readthedocs.io>
- Repository: <https://github.com/kaveio/eskapade>
- Docker: <https://github.com/kaveio/eskapade-environment>
- Issues & Ideas: <https://github.com/kaveio/eskapade/issues>
- Eskapade: <http://eskapade.kave.io>
- Contact us at: kave [at] kpmg [dot] com

6.9 API

6.9.1 API Documentation

Eskapade

eskapade package

Subpackages

eskapade.analysis package

Subpackages

eskapade.analysis.links package

Submodules

eskapade.analysis.links.apply_func_to_df module

Project: Eskapade - A python-based package for data analysis.

Class: ApplyFuncToDf

Created: 2016/11/08

Description: Algorithm to apply one or more functions to a (grouped) dataframe column or to an entire dataframe.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

```
class eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf (**kwargs)
    Bases: escore.core.element.Link
```

Apply functions to data-frame.

Applies one or more functions to a (grouped) dataframe column or an entire dataframe. In the latter case, this can be done row wise or column wise. The input dataframe will be overwritten.

```
__init__ (**kwargs)
    Initialize link instance.
```

Parameters

- **read_key** (*str*) – data-store input key
- **store_key** (*str*) – data-store output key
- **apply_funcs** (*list*) – functions to apply (list of dicts) - ‘func’: function to apply - ‘colout’ (string): output column - ‘colin’ (string, optional): input column - ‘entire’ (boolean, optional): apply to the entire dataframe? - ‘args’ (tuple, optional): args for ‘func’ - ‘kwargs’ (dict, optional): kwargs for ‘func’ - ‘groupby’ (list, optional): column names to group by - ‘groupbyColout’ (string) output column after the split-apply-combine combination

- **add_columns** (*dict*) – columns to add to output (name, column)

add_apply_func (*func*, *out_column*, *in_column*=”, *args, **kwargs)
Add function to be applied to dataframe.

execute ()
Execute the link.

groupbyapply (*df*, *groupby_columns*, *applyfunc*, *args, **kwargs)
Apply groupby to dataframe.

initialize ()
Initialize the link.

eskapade.analysis.links.apply_selection_to_df module

Project: Eskapade - A python-based package for data analysis.

Class: ApplySelectionToDf

Created: 2016/11/08

Description: Algorithm to apply queries to input dataframe

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.links.apply_selection_to_df.ApplySelectionToDf` (**kwargs)
Bases: `escore.core.element.Link`

Applies queries with sub-selections to a pandas dataframe.

__init__ (**kwargs)
Initialize link instance.

Input dataframe is not overwritten, unless instructed to do so in kwargs.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of data to read from data store
- **store_key** (*str*) – key of data to store in data store. If not set read_key is overwritten.
- **query_set** (*list*) – list of strings, query expressions to evaluate in the same order, see pandas documentation
- **select_columns** (*list*) – column names to select after querying
- **continue_if_failure** (*bool*) – if True continues with next query after failure (optional)
- **kwargs** – all other key word arguments are passed on to the pandas queries.

execute ()
Execute the link.

Applies queries or column selection to a pandas DataFrame. Input dataframe is not overwritten, unless told to do so in kwargs.

1. Apply queries, in order of provided query list.

2. Select columns (if provided).

```
initialize ()
    Initialize the link.

    Perform checks on provided attributes.
```

eskapade.analysis.links.basic_generator module

Project: Eskapade - A python-based package for data analysis.

Class: BasicGenerator

Created: 2017/02/26

Description: Link to generate random data with basic distributions

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

```
class eskapade.analysis.links.basic_generator.BasicGenerator (**kwargs)
    Bases: escore.core.element.Link
```

Generate data with basic distributions.

```
__init__ (**kwargs)
    Initialize link instance.
```

Parameters

- **key** (*str*) – key of output data in data store
- **columns** (*list*) – output column names
- **size** (*int*) – number of variable values
- **gen_config** (*dict*) – generator configuration for each variable
- **gen_seed** (*int*) – generator random seed

```
execute ()
    Execute the link.
```

```
initialize ()
    Initialize the link.
```

eskapade.analysis.links.df_concatenator module

Project: Eskapade - A python-based package for data analysis.

Class: DataFrameColumnRenamer

Created: 2016/11/08

Description: Algorithm to concatenate multiple pandas datadrames

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.links.df_concatenator.DfConcatenator` (**kwargs)

Bases: `escore.core.element.Link`

Concatenates multiple pandas datadrames.

`__init__` (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **store_key** (*str*) – key of data to store in data store
- **read_keys** (*list*) – keys of pandas dataframes in the data store
- **ignore_missing_input** (*bool*) – Skip missing input datasets. If all missing, store empty dataset. Default is false.
- **kwargs** – all other key word arguments are passed on to pandas concat function.

`execute` ()

Execute the link.

Perform concatenation of multiple pandas datadrames.

`initialize` ()

Initialize the link.

eskapade.analysis.links.df_merger module

Project: Eskapade - A python-based package for data analysis.

Class: DfMerger

Created: 2016/11/08

Description: Algorithm to Merges two pandas DataFrames

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.links.df_merger.DfMerger` (**kwargs)

Bases: `escore.core.element.Link`

Merges two pandas dataframes.

`__init__` (**kwargs)

Initialize link instance.

Store the configuration of the link.

Parameters

- **name** (*str*) – name of link
- **input_collection1** (*str*) – datastore key of the first pandas.DataFrame to merge
- **input_collection2** (*str*) – datastore key of the second pandas.DataFrame to merge
- **output_collection** (*str*) – datastore key of the merged output pandas.DataFrame
- **how** (*str*) – merge modus. See pandas documentation.

- **on** (*list*) – column names. See pandas documentation.
- **columns1** (*list*) – column names of the first pandas.DataFrame. Only these columns are included in the merge. If not set, use all columns.
- **columns2** (*list*) – column names of the second pandas.DataFrame. Only these columns are included in the merge. If not set, use all columns.
- **remove_duplicate_cols2** (*bool*) – if True duplicate columns will be taken out before the merge (default=True)
- **kwargs** – all other key word arguments are passed on to the pandas merge function.

execute ()

Perform merging of input dataframes.

initialize ()

Perform basic checks on provided attributes.

eskapade.analysis.links.histogrammar_filler module

Project: Eskapade - A python-based package for data analysis.

Class: HistogrammarFiller

Created: 2017/03/21

Description: Algorithm to fill histogrammar sparse-bin histograms. It is possible to do cleaning of these histograms by rejecting certain keys or removing inconsistent data types. Timestamp columns are converted to nanoseconds before the binning is applied.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.links.histogrammar_filler.HistogrammarFiller` (**kwargs)

Bases: `eskapade.analysis.histogram_filling.HistogramFillerBase`

Fill histogrammar sparse-bin histograms.

Algorithm to fill histogrammar style sparse-bin and category histograms.

It is possible to do after-filling cleaning of these histograms by rejecting certain keys or removing inconsistent data types. Timestamp columns are converted to nanoseconds before the binning is applied. Final histograms are stored in the datastore.

Example is available in: `tutorials/esk303_hgr_filler_plotter.py`

__init__ (**kwargs)

Initialize link instance.

Store and do basic check on the attributes of link HistogrammarFiller.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **store_key** (*str*) – key of output data to store histograms in data store
- **columns** (*list*) – columns to pick up from input data. (default is all columns)
- **bin_specs** (*dict*) – dictionaries used for rebinning numeric or timestamp columns

Example `bin_specs` dictionary is:

```
>>> bin_specs = {'x': {'bin_width': 1, 'bin_offset': 0},
                 'y': {'bin_edges': [0, 2, 3, 4, 5, 7, 8]}}
```

Parameters

- **var_dtype** (*dict*) – dict of datatypes of the columns to study from dataframe. If not provided, try to determine datatypes directly from dataframe.
- **quantity** (*dict*) – dictionary of lambda functions of how to parse certain columns

Example `quantity` dictionary is:

```
>>> quantity = {'y': lambda x: x}
```

Parameters

- **store_at_finalize** (*bool*) – Store histograms in datastore at `finalize()`, not at `execute()`. Useful when looping over datasets. Default is `False`.
- **dict** (*drop_keys*) – dictionary used for dropping specific keys from bins dictionaries of histograms

Example `drop_keys` dictionary is:

```
>>> drop_keys = {'x': [1, 4, 8, 19],
                 'y': ['apple', 'pear', 'tomato'],
                 'x:y': [(1, 'apple'), (19, 'tomato')]}
```

construct_empty_hist (*columns*)

Create an (empty) histogram of right type.

Create a multi-dim histogram by iterating through the columns in reverse order and passing a single-dim hist as input to the next column.

Parameters `columns` (*list*) – histogram columns

Returns created histogram

Return type `histogrammar.Count`

fill_histogram (*idf, columns*)

Fill input histogram with column(s) of input dataframe.

Parameters

- **idf** – input data frame used for filling histogram
- **columns** (*list*) – histogram column(s)

eskapade.analysis.links.random_sample_splitter module

Project: Eskapade - A python-based package for data analysis.

Class: `RandomSampleSplitter`

Created: 2016/11/08

Description: Algorithm to randomly assign records to a number of classes

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.links.random_sample_splitter.RandomSampleSplitter` (***kwargs*)
 Bases: `escore.core.element.Link`

Link that randomly assigns records of an input dataframe to a number of classes.

After assigning classes does one of the following:

- splits the input dataframe into sub dataframes according classes and stores the sub dataframes into the datastore;
- add a new column with assigned classes to the dataframe.

Records are assigned randomly.

`__init__` (***kwargs*)
 Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of data to read from datastore
- **store_key** (*list*) – keys of datasets to store in datastore. Number of sub samples equals length of store_key list (optional instead of ‘column’ and ‘nclasses’).
- **column** (*str*) – name of new column that specifies the randomly assigned class. Default is randomclass (optional instead of ‘store_key’).
- **nclasses** (*int*) – number of random classes. Needs to be set (optional instead of ‘store_key’).
- **fractions** (*list*) – list of fractions ($0 < \text{fraction} < 1$) of records assigned to the sub samples. Can be one less than n classes. Sum can be less than 1. Needs to be set.
- **nevents** (*list*) – list of number of random records assigned to the sub samples Can be one less than n classes (optional instead of ‘fractions’).

`execute` ()
 Execute the link.

`initialize` ()
 Check and initialize attributes of the link.

`eskapade.analysis.links.read_to_df` module

Project: Eskapade - A python-based package for data analysis.

Class: ReadToDf

Created: 2016/11/08

Description: Algorithm to write pandas dataframes picked up from the datastore.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.links.read_to_df.ReadToDf` (**kwargs)

Bases: `escore.core.element.Link`

Reads input file(s) to a pandas dataframe.

You give the link a path where your file is located and some kwargs that go into a pandas DataFrame. The kwargs are passed into the file reader.

__init__ (**kwargs)

Initialize link instance.

Store the configuration of link ReadToDf.

Parameters

- **name** (*str*) – Name given to the link
- **path** (*str*) – path of your file to read into pandas DataFrame .
- **key** (*str*) – storage key for the DataStore.
- **reader** – reader is determined automatically. But can be set by hand, e.g. csv, xlsx. To use the numpy reader one of the following should be true:
 - reader is { 'numpy', 'np', 'npz', 'npz' }
 - path contains extensions { 'npz', 'npz' }
 - param *file_type* is { 'npz', 'npz' }

To use the feather reader one of the following should be true:

- reader is { 'feather', 'ft' }
- path contains extensions 'ft'

When to use feather or which numpy type see the `esk210_dataframe_restoration` tutorial :param bool `restore_index`: whether to store the index in the metadata. Default is False when the index is numeric, True otherwise. :param str `file_type`: { 'npz', 'npz' } when using the numpy reader Optional, see reader for details. :param bool `itr_over_files`: Iterate over individual files, default is false. If false, are files are collected in one dataframe. NB chunksize takes priority! :param int `chunksize`: Default is none. If positive integer then will always iterate. chunksize requires `pd.read_csv` or `pd.read_table`. :param kwargs: all other key word arguments are passed on to the pandas reader.

execute ()

Execute the link.

Reads the input file(s) and puts the dataframe in the datastore.

initialize ()

Initialize the link.

is_finished () → bool

Try to assess if looper is done iterating over files.

Assess if looper is done or if a next dataset is still coming up.

latest_data_length ()

Return length of current dataset.

set_chunk_size (*size*)

Set chunksize setting.

Parameters *size* – chunk size

sum_data_length()

Return sum length of all datasets processed sofar.

`eskapade.analysis.links.read_to_df.feather_reader(path, restore_index)`

Read from feather file from disk to DataFrame, restoring the metadata

Parameters

- **path** (*str*) – target file location
- **restore_index** (*bool*) – store index in DataFrame Default is True

Returns df the DF read from disk

Return type `pd.DataFrame`

`eskapade.analysis.links.read_to_df.numpy_reader(path, restore_index, file_type)`

Read from numpy file from disk to DataFrame, restoring the metadata

Parameters

- **path** (*str*) – target file location
- **restore_index** (*bool*) – store index in DataFrame Default is True
- **file_type** (*str*) – the file type used { 'npz', 'npz' }

Raises

- **AmbiguousFileType** – when we can't determine whether the file type is npy or npz
- **UnhandledFileType** – generic catch for when the type logic fails to exclude case

Returns df the DF read from disk

Return type `pd.DataFrame`

`eskapade.analysis.links.read_to_df.set_reader(path, reader, *args, **kwargs)`

Pick the correct reader.

Based on provided reader setting, or based on file extension.

eskapade.analysis.links.record_factorizer module

Project: Eskapade - A python-based package for data analysis.

Class: RecordFactorizer

Created: 2016/11/08

Description: Algorithm to perform the factorization of an input column of an input dataframe. E.g. a column x with values 'apple', 'tree', 'pear', 'apple', 'pear' is tranformed into columns x with values 0, 1, 2, 0, 2, etc.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.links.record_factorizer.RecordFactorizer(**kwargs)`

Bases: `escore.core.element.Link`

Factorize data-frame columns.

Perform factorization of input column of an input dataframe. E.g. a column x with values 'apple', 'tree', 'pear', 'apple', 'pear' is tranformed into columns x with values 0, 1, 2, 0, 2, etc. Resulting dataset stored as new dataset. Alternatively, map transformed columns back to orginal format.

`__init__` (**kwargs)

Initialize link instance.

Store and do basic check on the attributes of link RecordFactorizer

Parameters

- **read_key** (*str*) – key to read dataframe from the data store. Dataframe of records that is to be transformed.
- **columns** (*list*) – list of columns that are to be factorized
- **inplace** (*bool*) – replace original columns. Default is False. Overwrites store_key to read_key.
- **convert_all_categories** (*bool*) – if true, convert all category observables. Default is false.
- **convert_all_booleans** (*bool*) – if true, convert all boolean observables. Default is false.
- **map_to_original** (*dict*) – dictionary or key to dictionary to map back factorized columns to original. map_to_original is a dict of dicts, one dict for each column.
- **store_key** (*str*) – store key of output dataframe. Default is read_key + ‘_fact’. (optional)
- **sk_map_to_original** (*str*) – store key of dictionary to map factorized columns to original. Default is ‘key’ + ‘_’ + store_key + ‘_to_original’. (optional)
- **sk_map_to_factorized** (*str*) – store key of dictionary to map original to factorized columns. Default is ‘key’ + ‘_’ + read_key + ‘_to_factorized’. (optional)

execute ()

Execute the link.

Perform factorization input columns ‘columns’ of input dataframe. Resulting dataset stored as new dataset. Alternatively, map transformed columns back to original format.

initialize ()

Initialize the link.

Initialize and (further) check the assigned attributes of the RecordFactorizer

eskapade.analysis.links.record_vectorizer module

Project: Eskapade - A python-based package for data analysis.

Class: RecordVectorizer

Created: 2016/11/08

Description: Algorithm to perform the vectorization of an input column of an input dataframe.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.links.record_vectorizer.RecordVectorizer` (**kwargs)

Bases: `escore.core.element.Link`

Vectorize data-frame columns.

Perform vectorization of input column of an input dataframe. E.g. a column `x` with values 1, 2 is transformed into columns `x_1` and `x_2`, with values True or False assigned per record.

`__init__` (**kwargs)

Initialize link instance.

Store and do basic check on the attributes of link RecordVectorizer.

Parameters

- **read_key** (*str*) – key to read dataframe from the data store. Dataframe of records that is to be transformed.
- **columns** (*list*) – list of columns that are to be vectorized
- **store_key** (*str*) – store key of output dataframe. Default is `read_key + '_vectorized'`. (optional)
- **column_compare_with** (*dict*) – dict of unique items per column with which column values are compared. If not given, this is derived automatically from the column. (optional)
- **astype** (*type*) – store answer of comparison of column with value as certain type. Default is `bool`. (optional)

`execute` ()

Execute the link.

Perform vectorization input column 'column' of input dataframe. Resulting dataset stored as new dataset.

`initialize` ()

Initialize the link.

Initialize and (further) check the assigned attributes of RecordVectorizer.

```
eskapade.analysis.links.record_vectorizer.record_vectorizer (df, column_to_vectorize,
                                                             col-
                                                             umn_compare_set,
                                                             astype=<class
                                                             'bool'>)
```

Vectorize data-frame column.

Takes the new record that is already transformed and vectorizes the given columns.

Parameters

- **df** – dataframe of the new record to vectorize
- **column_to_vectorize** (*str*) – string, column in the new record to vectorize.
- **column_compare_set** (*list*) – list of values to compare the column with.

Returns dataframe of the new records.

eskapade.analysis.links.value_counter module

Project: Eskapade - A python-based package for data analysis.

Class: ValueCounter

Created: 2017/03/02

Description: Algorithm to do `value_counts()` on single columns of a pandas dataframe, or `groupby().size()` on multiple columns, both returned as dictionaries. It is possible to do cleaning of these dicts by rejecting certain keys or removing inconsistent data types. Numeric and timestamp columns are converted to bin indices before the binning is applied. Results are stored as 1D Histograms or as ValueCounts objects.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.links.value_counter.ValueCounter` (***kwargs*)
 Bases: `eskapade.analysis.histogram_filling.HistogramFillerBase`

Count values in Pandas data frame.

ValueCounter does `value_counts()` on single columns of a pandas dataframe, or `groupby().size()` on multiple columns. Results of both are returned as same-style dictionaries.

Numeric and timestamp columns are converted to bin indices before the binning is applied. The binning can be provided as input.

It is possible to do cleaning of these dicts by rejecting certain keys or removing inconsistent data types. Results are stored as 1D Histograms or as ValueCounts objects.

Example is available in: `tutorials/esk302_histogram_filling_plotting.py`

`__init__` (***kwargs*)
 Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **store_key_counts** (*str*) – key of output data to store ValueCounts objects in data store
- **store_key_hists** (*str*) – key of output data to store histograms in data store
- **columns** (*list*) – columns to pick up from input data (default is all columns)
- **bin_specs** (*dict*) – dictionaries used for rebinning numeric or timestamp columns

Example `bin_specs` dictionary is:

```

>>> bin_specs = {'x': {'bin_width': 1, 'bin_offset': 0},
>>>                'y': {'bin_edges': [0, 2, 3, 4, 5, 7, 8]},
>>>                'date': {'bin_width': np.timedelta64(30, 'D'),
>>>                          'bin_offset': np.datetime64('2010-01-04')}}
    
```

Parameters

- **var_dtype** (*dict*) – dict of datatypes of the columns to study from dataframe. If not provided, try to determine datatypes directly from dataframe.
- **store_at_finalize** (*bool*) – Store histograms and/or ValueCount object in datatore at `finalize()`, not at `execute()`. Useful when looping over datasets. Default is `False`.
- **drop_inconsistent_key_types** (*bool*) – cleanup histograms and/or ValueCount objects by removing alls bins/keys with inconsistent datatypes. By default compare with data types in `var_dtype` dictionary.

- **dict** (*drop_keys*) – dictionary used for dropping specific keys from created value_counts dictionaries

Example drop_keys dictionary is:

```
>>> drop_keys = {'x': [1, 4, 8, 19],
>>>               'y': ['apple', 'pear', 'tomato'],
>>>               'x:y': [(1, 'apple'), (19, 'tomato')]}
```

drop_inconsistent_keys (*columns, obj*)

Drop inconsistent keys.

Drop inconsistent keys from a ValueCounts or Histogram object.

Parameters

- **columns** (*list*) – columns key to retrieve desired datatypes
- **obj** (*object*) – ValueCounts or Histogram object to drop inconsistent keys from

fill_histogram (*idf, columns*)

Fill input histogram with column(s) of input dataframe.

Parameters

- **idf** – input data frame used for filling histogram
- **columns** (*list*) – histogram column(s)

finalize ()

Finalize ValueCounter.

initialize ()

Initialize the link.

process_and_store ()

Make, clean, and store ValueCount objects.

process_columns (*df*)

Process columns before histogram filling.

Specifically, convert timestamp columns to integers and numeric variables are converted to indices

Parameters **df** – input (pandas) data frame

Returns output (pandas) data frame with converted timestamp columns

Return type pandas DataFrame

eskapade.analysis.links.write_from_df module

Project: Eskapade - A python-based package for data analysis.

Class: WriteFromDf

Created: 2016/11/08

Description: Algorithm to write a DataFrame from the DataStore to disk

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.links.write_from_df.WriteFromDf` (**kwargs)

Bases: `escore.core.element.Link`

Write a DataFrame from the DataStore to disk.

`__init__` (**kwargs)

Store the configuration of the link.

Parameters

- **name** (*str*) – Name given to the link
- **key** (*str*) – the DataStore key
- **path** (*str*) – path where to save the DataFrame
- **writer** – file extension that can be written by a pandas writer function from `pd.DataFrame`, or the numpy- feather writers. For example: ‘csv’ will trigger the `DataFrame.to_csv`. To use numpy_writer specify one of the following:

{‘numpy’, ‘np’, ‘npy’, ‘npz’, }

To use feather specify: {‘feather’, ‘ft’} If writer is not passed the path must contain a known file extension. Valid numpy extensions {‘npy’, ‘npz’} or feather {‘ft’}

Note the numpy and feather writers will preserve the metadata such as dtypes for each column and the index if non numeric.

Parameters

- **dictionary** (*dict*) – keys (as in the arg above) and paths (as in the arg above) it will write out all the keys to the associated paths.
- **add_counter_to_name** (*bool*) – if true, add an index to the output file name. Useful when running in loops. Default is false.
- **store_index** (*bool*) – whether the index should be stored as metadata. Default is False unless the index is non-numeric
- **kwargs** – all other key word arguments are passed on to the pandas writers.

`execute` ()

Execute the link.

Pick up the dataframe and write to disk.

`initialize` ()

Initialize the link.

`eskapade.analysis.links.write_from_df.feather_writer` (*df, path, store_index*)

Write df to disk in feather format; preserving the metadata

Parameters

- **df** (*DataFrame*) – pandas DataFrame to write out
- **path** (*str*) – target file location
- **store_index** (*bool*) – store index in DataFrame, default is True

`eskapade.analysis.links.write_from_df.get_writer` (*path, writer, *args, **kwargs*)

Pick the correct writer.

Based on provided writer setting, or based on file extension.

`eskapade.analysis.links.write_from_df.numpy_writer` (*df, path, store_index*)

Write df to disk in numpy format; preserving the metadata

Parameters

- **df** (*DataFrame*) – pandas Dataframe to write out
- **path** (*str*) – target file location
- **store_index** (*bool*) – store index in DataFrame

Module contents

class `eskapade.analysis.links.ApplyFuncToDf` (**kwargs)

Bases: `escore.core.element.Link`

Apply functions to data-frame.

Applies one or more functions to a (grouped) dataframe column or an entire dataframe. In the latter case, this can be done row wise or column wise. The input dataframe will be overwritten.

`__init__` (**kwargs)

Initialize link instance.

Parameters

- **read_key** (*str*) – data-store input key
- **store_key** (*str*) – data-store output key
- **apply_funcs** (*list*) – functions to apply (list of dicts) - ‘func’: function to apply - ‘colout’ (string): output column - ‘colin’ (string, optional): input column - ‘entire’ (boolean, optional): apply to the entire dataframe? - ‘args’ (tuple, optional): args for ‘func’ - ‘kwargs’ (dict, optional): kwargs for ‘func’ - ‘groupby’ (list, optional): column names to group by - ‘groupbyColout’ (string) output column after the split-apply-combine combination
- **add_columns** (*dict*) – columns to add to output (name, column)

add_apply_func (*func, out_column, in_column=*”, *args, **kwargs)

Add function to be applied to dataframe.

execute ()

Execute the link.

groupbyapply (*df, groupby_columns, applyfunc, *args, **kwargs*)

Apply groupby to dataframe.

initialize ()

Initialize the link.

class `eskapade.analysis.links.ApplySelectionToDf` (**kwargs)

Bases: `escore.core.element.Link`

Applies queries with sub-selections to a pandas dataframe.

`__init__` (**kwargs)

Initialize link instance.

Input dataframe is not overwritten, unless instructed to do so in kwargs.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of data to read from data store

- **store_key** (*str*) – key of data to store in data store. If not set read_key is overwritten.
- **query_set** (*list*) – list of strings, query expressions to evaluate in the same order, see pandas documentation
- **select_columns** (*list*) – column names to select after querying
- **continue_if_failure** (*bool*) – if True continues with next query after failure (optional)
- **kwargs** – all other key word arguments are passed on to the pandas queries.

execute ()

Execute the link.

Applies queries or column selection to a pandas DataFrame. Input dataframe is not overwritten, unless told to do so in kwargs.

1. Apply queries, in order of provided query list.
2. Select columns (if provided).

initialize ()

Initialize the link.

Perform checks on provided attributes.

class `eskapade.analysis.links.BasicGenerator` (**kwargs)

Bases: `escore.core.element.Link`

Generate data with basic distributions.

__init__ (**kwargs)

Initialize link instance.

Parameters

- **key** (*str*) – key of output data in data store
- **columns** (*list*) – output column names
- **size** (*int*) – number of variable values
- **gen_config** (*dict*) – generator configuration for each variable
- **gen_seed** (*int*) – generator random seed

execute ()

Execute the link.

initialize ()

Initialize the link.

class `eskapade.analysis.links.DfConcatenator` (**kwargs)

Bases: `escore.core.element.Link`

Concatenates multiple pandas datadrames.

__init__ (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **store_key** (*str*) – key of data to store in data store
- **read_keys** (*list*) – keys of pandas dataframes in the data store

- **ignore_missing_input** (*bool*) – Skip missing input datasets. If all missing, store empty dataset. Default is false.
- **kwargs** – all other key word arguments are passed on to pandas concat function.

execute()

Execute the link.

Perform concatenation of multiple pandas datadrames.

initialize()

Initialize the link.

class `eskapade.analysis.links.DfMerger` (**kwargs)

Bases: `escore.core.element.Link`

Merges two pandas dataframes.

__init__ (**kwargs)

Initialize link instance.

Store the configuration of the link.

Parameters

- **name** (*str*) – name of link
- **input_collection1** (*str*) – datastore key of the first pandas.DataFrame to merge
- **input_collection2** (*str*) – datastore key of the second pandas.DataFrame to merge
- **output_collection** (*str*) – datastore key of the merged output pandas.DataFrame
- **how** (*str*) – merge modus. See pandas documentation.
- **on** (*list*) – column names. See pandas documentation.
- **columns1** (*list*) – column names of the first pandas.DataFrame. Only these columns are included in the merge. If not set, use all columns.
- **columns2** (*list*) – column names of the second pandas.DataFrame. Only these columns are included in the merge. If not set, use all columns.
- **remove_duplicate_cols2** (*bool*) – if True duplicate columns will be taken out before the merge (default=True)
- **kwargs** – all other key word arguments are passed on to the pandas merge function.

execute()

Perform merging of input dataframes.

initialize()

Perform basic checks on provided attributes.

class `eskapade.analysis.links.HistogrammarFiller` (**kwargs)

Bases: `eskapade.analysis.histogram_filling.HistogramFillerBase`

Fill histogrammar sparse-bin histograms.

Algorithm to fill histogrammar style sparse-bin and category histograms.

It is possible to do after-filling cleaning of these histograms by rejecting certain keys or removing inconsistent data types. Timestamp columns are converted to nanoseconds before the binning is applied. Final histograms are stored in the datastore.

Example is available in: `tutorials/esk303_hgr_filler_plotter.py`

`__init__` (**kwargs)

Initialize link instance.

Store and do basic check on the attributes of link HistogrammarFiller.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **store_key** (*str*) – key of output data to store histograms in data store
- **columns** (*list*) – cols to pick up from input data. (default is all columns)
- **bin_specs** (*dict*) – dictionaries used for rebinning numeric or timestamp columns

Example bin_specs dictionary is:

```
>>> bin_specs = {'x': {'bin_width': 1, 'bin_offset': 0},
                 'y': {'bin_edges': [0, 2, 3, 4, 5, 7, 8]}}
```

Parameters

- **var_dtype** (*dict*) – dict of datatypes of the columns to study from dataframe If not provided, try to determine datatypes directly from dataframe.
- **quantity** (*dict*) – dictionary of lambda functions of how to pars certain columns

Example quantity dictionary is:

```
>>> quantity = {'y': lambda x: x}
```

Parameters

- **store_at_finalize** (*bool*) – Store histograms in datastore at finalize(), not at execute(). Useful when looping over datasets. Default is False.
- **dict** (*drop_keys*) – dictionary used for dropping specific keys from bins dictionaries of histograms

Example drop_keys dictionary is:

```
>>> drop_keys = {'x': [1, 4, 8, 19],
                 'y': ['apple', 'pear', 'tomato'],
                 'x:y': [(1, 'apple'), (19, 'tomato')]}
```

construct_empty_hist (*columns*)

Create an (empty) histogram of right type.

Create a multi-dim histogram by iterating through the columns in reverse order and passing a single-dim hist as input to the next column.

Parameters **columns** (*list*) – histogram columns

Returns created histogram

Return type histogrammar.Count

fill_histogram (*idf, columns*)

Fill input histogram with column(s) of input dataframe.

Parameters

- **idf** – input data frame used for filling histogram
- **columns** (*list*) – histogram column(s)

class `eskapade.analysis.links.RandomSampleSplitter` (**kwargs)

Bases: `escore.core.element.Link`

Link that randomly assigns records of an input dataframe to a number of classes.

After assigning classes does one of the following:

- splits the input dataframe into sub dataframes according classes and stores the sub dataframes into the datastore;
- add a new column with assigned classes to the dataframe.

Records are assigned randomly.

`__init__` (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of data to read from datastore
- **store_key** (*list*) – keys of datasets to store in datastore. Number of sub samples equals length of store_key list (optional instead of ‘column’ and ‘nclasses’).
- **column** (*str*) – name of new column that specifies the randomly assigned class. Default is `randomclass` (optional instead of ‘store_key’).
- **nclasses** (*int*) – number of random classes. Needs to be set (optional instead of ‘store_key’).
- **fractions** (*list*) – list of fractions ($0 < \text{fraction} < 1$) of records assigned to the sub samples. Can be one less than n classes. Sum can be less than 1. Needs to be set.
- **nevents** (*list*) – list of number of random records assigned to the sub samples Can be one less than n classes (optional instead of ‘fractions’).

`execute` ()

Execute the link.

`initialize` ()

Check and initialize attributes of the link.

class `eskapade.analysis.links.ReadToDf` (**kwargs)

Bases: `escore.core.element.Link`

Reads input file(s) to a pandas dataframe.

You give the link a path where your file is located and some kwargs that go into a pandas DataFrame. The kwargs are passed into the file reader.

`__init__` (**kwargs)

Initialize link instance.

Store the configuration of link `ReadToDf`.

Parameters

- **name** (*str*) – Name given to the link
- **path** (*str*) – path of your file to read into pandas DataFrame .

- **key** (*str*) – storage key for the DataStore.
- **reader** – reader is determined automatically. But can be set by hand, e.g. csv, xlsx. To use the numpy reader one of the following should be true:
 - reader is { 'numpy', 'np', 'npy', 'npz' }
 - path contains extensions { 'npz', 'npz' }
 - param *file_type* is { 'npz', 'npz' }

To use the feather reader one of the following should be true:

- reader is { 'feather', 'ft' }
- path contains extensions 'ft'

When to use feather or which numpy type see the esk210_dataframe_restoration tutorial :param bool restore_index: whether to store the index in the metadata. Default is False when the index is numeric, True otherwise. :param str file_type: { 'npz', 'npz' } when using the numpy reader Optional, see reader for details. :param bool itr_over_files: Iterate over individual files, default is false. If false, are files are collected in one dataframe. NB chunksize takes priority! :param int chunksize: Default is none. If positive integer then will always iterate. chunksize requires pd.read_csv or pd.read_table. :param kwargs: all other key word arguments are passed on to the pandas reader.

execute()

Execute the link.

Reads the input file(s) and puts the dataframe in the datastore.

initialize()

Initialize the link.

is_finished() → bool

Try to assess if looper is done iterating over files.

Assess if looper is done or if a next dataset is still coming up.

latest_data_length()

Return length of current dataset.

set_chunk_size(size)

Set chunksize setting.

Parameters size – chunk size

sum_data_length()

Return sum length of all datasets processed sofar.

class eskapade.analysis.links.**RecordFactorizer** (**kwargs)

Bases: escore.core.element.Link

Factorize data-frame columns.

Perform factorization of input column of an input dataframe. E.g. a column x with values 'apple', 'tree', 'pear', 'apple', 'pear' is tranformed into columns x with values 0, 1, 2, 0, 2, etc. Resulting dataset stored as new dataset. Alternatively, map transformed columns back to orginal format.

__init__ (**kwargs)

Initialize link instance.

Store and do basic check on the attributes of link RecordFactorizer

Parameters

- **read_key** (*str*) – key to read dataframe from the data store. Dataframe of records that is to be transformed.
- **columns** (*list*) – list of columns that are to be factorized
- **inplace** (*bool*) – replace original columns. Default is False. Overwrites store_key to read_key.
- **convert_all_categories** (*bool*) – if true, convert all category observables. Default is false.
- **convert_all_booleans** (*bool*) – if true, convert all boolean observables. Default is false.
- **map_to_original** (*dict*) – dictionary or key to dictionary to map back factorized columns to original. map_to_original is a dict of dicts, one dict for each column.
- **store_key** (*str*) – store key of output dataframe. Default is read_key + ‘_fact’. (optional)
- **sk_map_to_original** (*str*) – store key of dictionary to map factorized columns to original. Default is ‘key’ + ‘_’ + store_key + ‘_to_original’. (optional)
- **sk_map_to_factorized** (*str*) – store key of dictionary to map original to factorized columns. Default is ‘key’ + ‘_’ + read_key + ‘_to_factorized’. (optional)

execute()

Execute the link.

Perform factorization input columns ‘columns’ of input dataframe. Resulting dataset stored as new dataset. Alternatively, map transformed columns back to original format.

initialize()

Initialize the link.

Initialize and (further) check the assigned attributes of the RecordFactorizer

```
class eskapade.analysis.links.RecordVectorizer (**kwargs)
```

Bases: `escore.core.element.Link`

Vectorize data-frame columns.

Perform vectorization of input column of an input dataframe. E.g. a column x with values 1, 2 is transformed into columns x_1 and x_2, with values True or False assigned per record.

```
__init__ (**kwargs)
```

Initialize link instance.

Store and do basic check on the attributes of link RecordVectorizer.

Parameters

- **read_key** (*str*) – key to read dataframe from the data store. Dataframe of records that is to be transformed.
- **columns** (*list*) – list of columns that are to be vectorized
- **store_key** (*str*) – store key of output dataframe. Default is read_key + ‘_vectorized’. (optional)
- **column_compare_with** (*dict*) – dict of unique items per column with which column values are compared. If not given, this is derived automatically from the column. (optional)

- **astype** (*type*) – store answer of comparison of column with value as certain type. Default is bool. (optional)

execute ()

Execute the link.

Perform vectorization input column ‘column’ of input dataframe. Resulting dataset stored as new dataset.

initialize ()

Initialize the link.

Initialize and (further) check the assigned attributes of RecordVectorizer.

class `eskapade.analysis.links.ValueCounter` (**kwargs)

Bases: `eskapade.analysis.histogram_filling.HistogramFillerBase`

Count values in Pandas data frame.

ValueCounter does `value_counts()` on single columns of a pandas dataframe, or `groupby().size()` on multiple columns. Results of both are returned as same-style dictionaries.

Numeric and timestamp columns are converted to bin indices before the binning is applied. The binning can be provided as input.

It is possible to do cleaning of these dicts by rejecting certain keys or removing inconsistent data types. Results are stored as 1D Histograms or as ValueCounts objects.

Example is available in: `tutorials/esk302_histogram_filling_plotting.py`

__init__ (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **store_key_counts** (*str*) – key of output data to store ValueCounts objects in data store
- **store_key_hists** (*str*) – key of output data to store histograms in data store
- **columns** (*list*) – columns to pick up from input data (default is all columns)
- **bin_specs** (*dict*) – dictionaries used for rebinning numeric or timestamp columns

Example `bin_specs` dictionary is:

```
>>> bin_specs = {'x': {'bin_width': 1, 'bin_offset': 0},
>>>               'y': {'bin_edges': [0, 2, 3, 4, 5, 7, 8]},
>>>               'date': {'bin_width': np.timedelta64(30, 'D'),
>>>                       'bin_offset': np.datetime64('2010-01-04')}
```

Parameters

- **var_dtype** (*dict*) – dict of datatypes of the columns to study from dataframe. If not provided, try to determine datatypes directly from dataframe.
- **store_at_finalize** (*bool*) – Store histograms and/or ValueCount object in data-store at `finalize()`, not at `execute()`. Useful when looping over datasets. Default is False.
- **drop_inconsistent_key_types** (*bool*) – cleanup histograms and/or ValueCount objects by removing alls bins/keys with inconsistent datatypes. By default compare with data types in `var_dtype` dictionary.

- **dict** (*drop_keys*) – dictionary used for dropping specific keys from created value_counts dictionaries

Example drop_keys dictionary is:

```
>>> drop_keys = {'x': [1, 4, 8, 19],
>>>               'y': ['apple', 'pear', 'tomato'],
>>>               'x:y': [(1, 'apple'), (19, 'tomato')]}
```

drop_inconsistent_keys (*columns, obj*)

Drop inconsistent keys.

Drop inconsistent keys from a ValueCounts or Histogram object.

Parameters

- **columns** (*list*) – columns key to retrieve desired datatypes
- **obj** (*object*) – ValueCounts or Histogram object to drop inconsistent keys from

fill_histogram (*idf, columns*)

Fill input histogram with column(s) of input dataframe.

Parameters

- **idf** – input data frame used for filling histogram
- **columns** (*list*) – histogram column(s)

finalize ()

Finalize ValueCounter.

initialize ()

Initialize the link.

process_and_store ()

Make, clean, and store ValueCount objects.

process_columns (*df*)

Process columns before histogram filling.

Specifically, convert timestamp columns to integers and numeric variables are converted to indices

Parameters **df** – input (pandas) data frame

Returns output (pandas) data frame with converted timestamp columns

Return type pandas DataFrame

class eskapade.analysis.links.**WriteFromDf** (***kwargs*)

Bases: `escore.core.element.Link`

Write a DataFrame from the DataStore to disk.

__init__ (***kwargs*)

Store the configuration of the link.

Parameters

- **name** (*str*) – Name given to the link
- **key** (*str*) – the DataStore key
- **path** (*str*) – path where to save the DataFrame

- **writer** – file extension that can be written by a pandas writer function from `pd.DataFrame`, or the numpy- feather writers. For example: ‘csv’ will trigger the `DataFrame.to_csv`. To use `numpy_writer` specify one of the following:

```
{‘numpy’, ‘np’, ‘npy’, ‘npz’, }
```

To use feather specify: {‘feather’, ‘ft’} If writer is not passed the path must contain a known file extension. Valid numpy extensions {‘npy’, ‘npz’} or feather {‘ft’}

Note the numpy and feather writers will preserve the metadata such as dtypes for each column and the index if non numeric.

Parameters

- **dictionary** (*dict*) – keys (as in the arg above) and paths (as in the arg above) it will write out all the keys to the associated paths.
- **add_counter_to_name** (*bool*) – if true, add an index to the output file name. Useful when running in loops. Default is false.
- **store_index** (*bool*) – whether the index should be stored as metadata. Default is False unless the index is non-numeric
- **kwargs** – all other key word arguments are passed on to the pandas writers.

`execute()`

Execute the link.

Pick up the dataframe and write to disk.

`initialize()`

Initialize the link.

Submodules

`eskapade.analysis.correlation module`

Project: Eskapade - A python-based package for data analysis.

Created: 2018/06/23

Description: Correlation related util functions.

Convert Pearson correlation value into a chi2 value of a contingency test matrix of a bivariate gaussian, and vice-versa. Calculation uses scipy’s mvn library. Calculates correlation coefficients based on `mutual_information`, `correlation_ratio`, `pearson`, `kendall` or `spearman` methods.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

`eskapade.analysis.correlation.calculate_correlations` (*df, method*)

Calculates correlation coefficients between every column pair.

Parameters

- **df** (*pd.DataFrame*) – input data frame
- **method** (*str*) – `mutual_information`, `correlation_ratio`, `pearson`, `kendall` or `spearman`, `phik`, `significance`

Returns `pd.DataFrame`

`eskapade.analysis.correlation.chi2_from_rho` (*rho*, *n*, *subtract_from_chi2=0*, *corr0=None*,
sx=None, *sy=None*, *nx=-1*, *ny=-1*)

Calculate chi2-value of bivariate gauss having correlation value rho

Calculate no-noise chi2 value of bivar gauss with correlation rho, with respect to bivariate gauss without any correlation.

Returns float chi2 value

`eskapade.analysis.correlation.rho_from_chi2` (*chi2*, *n*, *nx*, *ny*, *sx=None*, *sy=None*)
correlation coefficient of bivariate gaussian derived from chi2-value

Chi2-value gets converted into correlation coefficient of bivariate gauss with correlation value rho, assuming giving binning and number of records. Correlation coefficient value is between 0 and 1.

Bivariate gaussian's range is set to [-5,5] by construction.

Returns float correlation coefficient

eskapade.analysis.datetime module

Project: Eskapade - A python-based package for data analysis.

Classes: TimePeriod, FreqTimePeriod

Created: 2017/03/14

Description: Time period and time period with frequency.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

modification, are permitted according to the terms listed in the file Redistribution and use in source and binary forms, with or without LICENSE.

class `eskapade.analysis.datetime.FreqTimePeriod` (**kwargs)

Bases: `eskapade.analysis.datetime.TimePeriod`

Time period with frequency.

`__init__` (**kwargs)

Initialize TimePeriod instance.

`dt_string` (*period_index*)

Convert period index into date/time string (start of period).

Parameters `period_index` (*int*) – specified period index value.

`freq`

Return frequency.

`period_index` (*dt*)

Return number of periods until date/time “dt” since 1970-01-01.

Parameters `dt` – specified date/time parameter

class `eskapade.analysis.datetime.TimePeriod` (**kwargs)

Bases: `escore.core.mixin.ArgumentsMixin`

Time period.

`__init__` (**kwargs)

Initialize TimePeriod instance.

logger

A logger that emits log messages to an observer.

The logger can be instantiated as a module or class attribute, e.g.

```

>>> logger = Logger()
>>> logger.info("I'm a module logger attribute.")
>>>
>>> class Point(object):
>>>     logger = Logger()
>>>
>>>     def __init__(self, x = 0.0, y = 0.0):
>>>         Point.logger.debug('Initializing {point} with x = {x} y = {y}',
↳ point=Point, x=x, y=y)
>>>         self._x = x
>>>         self._y = y
>>>
>>>     @property
>>>     def x(self):
>>>         self.logger.debug('Getting property x = {point._x}', point=self)
>>>         return self._x
>>>
>>>     @x.setter
>>>     def x(self, x):
>>>         self.logger.debug('Setting property y = {point._x}', point=self)
>>>         self._x = x
>>>
>>>     @property
>>>     def y(self):
>>>         self.logger.debug('Getting property y = {point._y}', point=self)
>>>         return self._y
>>>
>>>     @y.setter
>>>     def y(self, y):
>>>         self.logger.debug('Setting property y = {point._y}', point=self)
>>>         self._y = y
>>>
>>> a_point = Point(1, 2)
>>>
>>> logger.info('p_x = {point.x} p_y = {point.y}', point=a_point)
>>> logger.log_level = LogLevel.DEBUG
>>> logger.info('p_x = {point.x} p_y = {point.y}', point=a_point)

```

The logger uses PEP-3101 (Advanced String Formatting) with named placeholders, see <<https://www.python.org/dev/peps/pep-3101/>> and <<https://pyformat.info/>> for more details and examples.

Furthermore, logging events are only formatted and evaluated for logging levels that are enabled. So, there's no need to check the logging level before logging. It's also efficient.

classmethod parse_date_time (*dt*)

Try to parse specified date/time.

Parameters *dt* – specified date/time

classmethod parse_time_period (*period*)

Try to parse specified time period.

Parameters *period* – specified period

period_index (*dt*)

Get number of periods until date/time “dt”.

Parameters `dt` – specified date/time

class `eskapade.analysis.datetime.UniformTsTimePeriod` (**kwargs)

Bases: `eskapade.analysis.datetime.TimePeriod`

Time period with offset.

`__init__` (**kwargs)

Initialize TimePeriod instance.

offset

Get offset parameter.

period

Get period parameter.

period_index (dt)

Get number of periods until date/time “dt” since “offset”, given specified “period”.

Parameters `dt` – specified date/time

eskapade.analysis.histogram module

Project: Eskapade - A python-based package for data analysis.

Classes: ValueCounts, BinningUtil, Histogram

Created: 2017/03/14

Description: Generic 1D Histogram class.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

modification, are permitted according to the terms listed in the file Redistribution and use in source and binary forms, with or without LICENSE.

class `eskapade.analysis.histogram.BinningUtil` (**kwargs)

Bases: `object`

Helper for interpreting bin specifications.

BinningUtil is a helper class used for interpreting bin specification dictionaries. It is a base class for the Histogram class.

`__init__` (**kwargs)

Initialize link instance.

A `bin_specs` dictionary needs to be provided as input. `bin_specs` is a dict containing ‘bin_width’ and ‘bin_offset’ keys. In case bins widths are not equal, `bin_specs` contains ‘bin_edges’ (array) instead of ‘bin_width’ and ‘bin_offset’. ‘bin_width’ and ‘bin_offset’ can be numeric or numpy timestamps.

Alternatively, `bin_edges` can be provided as input to `bin_specs`.

Example `bin_specs` dictionaries are:

```
>>> bin_specs = {'bin_width': 1, 'bin_offset': 0}
>>> bin_spect = {'bin_edges': [0, 2, 3, 4, 5, 7, 8]}
>>> bin_specs = {'bin_width': np.timedelta64(30, 'D'),
                  'bin_offset': np.datetime64('2010-01-04')}
```

Parameters

- **bin_specs** (*dict*) – dictionary contains ‘bin_width’ and ‘bin_offset’ numbers or ‘bin_edges’ array. Default is None.
- **bin_edges** (*list*) – array with numpy histogram style bin_edges. Default is None.

bin_specs

Get bin_specs dictionary.

Returns bin_specs dictionary

Return type dict

get_bin_center (*bin_label*)

Return bin center for a given bin index.

Parameters **bin_label** – bin label for which to find the bin center

Returns bin center, can be float, int, timestamp

get_bin_edges ()

Return bin edges.

Returns bin edges

Return type array

get_bin_edges_range ()

Return bin range determined from bin edges.

Returns bin range

Return type tuple

get_left_bin_edge (*bin_label*)

Return left bin edge for a given bin index.

Parameters **bin_label** – bin label for which to find the left bin edge

Returns bin edge, can be float, int, timestamp

get_right_bin_edge (*bin_label*)

Return right bin edge for a given bin index.

Parameters **bin_label** – bin label for which to find the right bin edge.

Returns bin edge, can be float, int, timestamp

truncated_bin_edges (*variable_range=None*)

Bin edges corresponding to a given variable range.

Parameters **variable_range** (*list*) – variable range used for finding the right bin edges array. Optional.

Returns truncated bin edges

Return type array

value_to_bin_label (*var_value, greater_equal=False*)

Return bin index for given bin value.

Parameters

- **var_value** – variable value for which to find the bin index
- **greater_equal** (*bool*) – for float, int, timestamp, return index of bin for which value falls in range [lower edge, upper edge). If set to true, return index of bin for which value falls in range [lower edge, upper edge]. Default if false.

Returns bin index

Return type int

class eskapade.analysis.histogram.**Histogram**(*counts*, ***kwargs*)

Bases: `eskapade.analysis.histogram.BinningUtil`, `escore.core.mixin.ArgumentsMixin`

Generic 1D Histogram class.

Histogram holds bin labels (name of each bin), `value_counts` (values of the histogram) and a variable name. The bins can be categoric or numeric, where numeric includes timestamps. In case of numeric bins, `bin_specs` is set. `bin_specs` is a dict containing `bin_width` and `bin_offset`. In case bins widths are not equal, `bin_specs` contains `bin_edges` instead of `bin_width` and `bin_offset`.

`__init__`(*counts*, ***kwargs*)

Initialize Histogram instance.

A `bin_specs` dictionary can be provided as input. `bin_specs` is a dict containing 'bin_width' and 'bin_offset' keys. In case bins widths are not equal, `bin_specs` contains 'bin_edges' (array) instead of 'bin_width' and 'bin_offset'. 'bin_width' and 'bin_offset' can be numeric or numpy timestamps.

Histogram counts can be specified as a ValueCounts object, a dictionary or a tuple:

- tuple: Histogram(`(bin_values, bin_edges)`, `variable=<your_variable_name>`)
- dict: a dictionary as comes out of `pandas.series.value_counts()` or `pandas.DataFrame.groupby.size()` over one variable.
- ValueCounts: a ValueCounts object contains a `value_counts` dictionary.

Example `bin_specs` dictionaries are:

```
>>> bin_specs = { 'bin_width': 1, 'bin_offset': 0 }
>>> bin_spect = { 'bin_edges': [0,2,3,4,5,7,8] }
>>> bin_specs = { 'bin_width': np.timedelta64(30,'D'),
                  'bin_offset': np.datetime64('2010-01-04') }
```

Parameters

- **counts** – histogram counts
- **bin_specs** (*dict*) – dictionary contains 'bin_width' and 'bin_offset' numbers or 'bin_edges' array (default is None)
- **variable** (*str*) – name of the variable represented by the histogram
- **datatype** (*type*) – data type of the variable represented by the histogram (optional)

bin_centers ()

Return bin centers.

Returns array of the bin centers

Return type array

bin_edges ()

Return numpy style `bin_edges` array with uniform binning.

Returns array of all bin edges

Return type array

bin_entries ()

Return number of bin entries.

Return the bin counts of the known bins in the `value_counts` object.

Returns array of the bin counts

Return type array

bin_labels ()

Return bin labels.

Returns array of all bin labels

Return type array

classmethod combine_hists (*hists*, *labels=False*, *rel_bin_width_tol=1e-06*, ***kwargs*)

Combine a set of histograms.

Parameters

- **hists** (*array*) – array of Histograms to add up.
- **labels** (*label*) – histograms to add up have labels? (else are numeric) Default is False.
- **variable** (*str*) – name of variable described by the summed-up histogram
- **rel_bin_width_tol** (*float*) – relative tolerance between numeric bin edges.

Returns summed up histogram

Return type *Histogram*

copy (***kwargs*)

Return a copy of this histogram.

Parameters **variable** (*str*) – assign new variable name

datatype

Data type of the variable represented by the histogram.

Returns data type

Return type type

get_bin_count (*bin_label*)

Get bin count for specific bin label.

Parameters **bin_label** – a specific key to find corresponding bin.

Returns bin counter value

Return type int

get_bin_labels ()

Return all bin labels.

Returns array of all bin labels

Return type array

get_bin_range ()

Return the bin range.

Returns tuple of the bin range found

Return type tuple

get_bin_vals (*variable_range=None, combine_values=True*)

Get bin labels/edges and corresponding bin counts.

Bin values corresponding to a given variable range.

Parameters

- **variable_range** (*list*) – variable range used for finding the right bins to get values from. Optional.
- **combine_values** (*bool*) – if `bin_specs` is not set, combine existing bin labels with variable range.

Returns two arrays of bin values and bin edges

Return type array

get_hist_val (*var_value*)

Get bin count for bin by value of histogram variable.

Parameters **var_value** – a specific value to find corresponding bin.

Returns bin counter value

Return type int

get_nonone_bin_centers ()

Return bin centers.

Returns array of the bin centers

Return type array

get_nonone_bin_counts ()

Return bin counts.

Returns array of the bin counts

Return type array

get_nonone_bin_edges ()

Return numpy style bin-edges array.

Returns array of the bin edges

Return type array

get_nonone_bin_range ()

Return the bin range.

Returns tuple of the bin range found

Return type tuple

get_uniform_bin_edges ()

Return numpy style bin-edges array with uniform binning.

Returns array of all bin edges

Return type array

logger

A logger that emits log messages to an observer.

The logger can be instantiated as a module or class attribute, e.g.

```

>>> logger = Logger()
>>> logger.info("I'm a module logger attribute.")
>>>
>>> class Point(object):
>>>     logger = Logger()
>>>
>>>     def __init__(self, x = 0.0, y = 0.0):
>>>         Point.logger.debug('Initializing {point} with x = {x} y = {y}',
↵point=Point, x=x, y=y)
>>>         self._x = x
>>>         self._y = y
>>>
>>>     @property
>>>     def x(self):
>>>         self.logger.debug('Getting property x = {point._x}', point=self)
>>>         return self._x
>>>
>>>     @x.setter
>>>     def x(self, x):
>>>         self.logger.debug('Setting property y = {point._x}', point=self)
>>>         self._x = x
>>>
>>>     @property
>>>     def y(self):
>>>         self.logger.debug('Getting property y = {point._y}', point=self)
>>>         return self._y
>>>
>>>     @y.setter
>>>     def y(self, y):
>>>         self.logger.debug('Setting property y = {point._y}', point=self)
>>>         self._y = y
>>>
>>> a_point = Point(1, 2)
>>>
>>> logger.info('p_x = {point.x} p_y = {point.y}', point=a_point)
>>> logger.log_level = LogLevel.DEBUG
>>> logger.info('p_x = {point.x} p_y = {point.y}', point=a_point)

```

The logger uses PEP-3101 (Advanced String Formatting) with named placeholders, see <<https://www.python.org/dev/peps/pep-3101/>> and <<https://pyformat.info/>> for more details and examples.

Furthermore, logging events are only formatted and evaluated for logging levels that are enabled. So, there's no need to check the logging level before logging. It's also efficient.

n_bins

Number of bins in the ValueCounts object.

Returns number of bins

Return type int

n_dim

Number of histogram dimensions.

The number of histogram dimensions, which is equal to one by construction.

Returns number of dimensions

Return type int

num_bins

Number of bins.

Returns number of bins

Return type int

remove_keys_of_inconsistent_type (*preferred_key_type=None*)

Remove all keys that have inconsistent data type(s).

Parameters **preferred_key_type** (*tuple*) – the preferred key type to keep. Can be a tuple, list, or single type. E.g. str or (int,str,float). If None provided, the most common key type found is kept.

simulate (*size, *args*)

Simulate data using self (Histogram instance) as PDF.

see https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.rv_continuous.html

Parameters **size** (*int*) – number of data points to generate

Return **numpy.array generated_data** the generated data

Returns Histogram of the generated data

Return type *Histogram*

surface ()

Calculate surface of the histogram.

Returns surface

to_normalized (***kwargs*)

Return a normalized copy of this histogram.

Parameters

- **new_var_name** (*str*) – assign new variable name
- **variable_range** (*list*) – variable range used for finding the right bins to get values from.
- **combine_values** (*bool*) – if bin_specs is not set, combine existing bin labels with variable range.

variable

Name of variable represented by the histogram.

Returns variable name

Return type string

class `eskapade.analysis.histogram.ValueCounts` (*key, subkey=None, counts=None, sel=None*)

Bases: object

A dictionary of value counts.

The dictionary of value counts comes out of `pandas.series.value_counts()` for one variable or `pandas.DataFrame.groupby.size()` performed over one or multiple variables.

__init__ (*key, subkey=None, counts=None, sel=None*)

Initialize link instance.

Parameters

- **key** (*list*) – key is a tuple, list or string of (the) variable name(s), matching those and the structure of the keys in the value_counts dictionary.
- **subkey** (*list*) – subset of key. If provided, the value_counts dictionary will be projected from key onto the (subset of) subkey. E.g. use this to map a two dimensional value_counts dictionary onto one specified dimension. Default is None. Optional.
- **counts** (*dict*) – the value_counts dictionary.
- **sel** (*dict*) – Apply selections to value_counts dictionary. Default is {}. Optional.

count (*value_bin*)

Get bin count for specific bin-key value bin.

Parameters **value_bin** (*tuple*) – a specific key, and can be a list or tuple.

Returns specific bin counter value

Return type int

counts

Process and return value-counts dictionary.

Returns after processing, returns the value_counts dictionary

Return type dict

create_sub_counts (*subkey, sel=None*)

Project existing value counts onto a subset of keys.

E.g. map variables x,y onto single dimension x, so for each bin in x integrate over y.

Parameters

- **subkey** (*tuple*) – input sub-key, is a tuple, list, or string. This is the new key of variables for the returned ValueCounts object.
- **sel** (*dict*) – dictionary with selection. Optional.

Returns value_counts object where subkey has become the new key.

Return type *ValueCounts*

get_values (*val_keys=()*)

Get all key-values of a subset of keys.

E.g. give all x values in of the keys, when the value_counts object has keys (x, y).

Parameters **val_keys** (*tuple*) – a specific sub-key to get key values for.

Returns all key-values of a subset of keys.

Return type tuple

key

Process and return current value-counts key.

Returns the key

Return type tuple

nononecounts

Return value-counts dictionary without None keys.

Returns after processing, returns the value_counts dictionary without None keys

Return type dict

num_bins

Number of value-counts bins.

Returns number of bins

Return type int

num_nonone_bins

Number of not-none value-counts bins.

Returns number of not-none bins

Return type int

process_counts (*accept_equiv=True*)

Project value counts onto the existing subset of keys.

E.g. map variables x,y onto single dimension x, so for each bin in x integrate over y.

Parameters **accept_equiv** (*bool*) – accept equivalence of key and subkey if subkey is in different order than key. Default is true.

Returns successful projection or not

Return type bool

remove_keys_of_inconsistent_type (*preferred_key_type=None*)

Remove keys with inconsistent data type(s).

Parameters **preferred_key_type** (*tuple*) – the preferred key type to keep. Can be a tuple, list, or single type. E.g. str or (int, str, float). If None provided, the most common key type found is kept.

skey

Current value-counts subkey.

Returns the subkey

Return type tuple

sum_counts

Sum of counts of all value-counts bins.

Returns the sum of counts of all bins

Return type float

sum_nonone_counts

Sum of not-none counts of all value-counts bins.

Returns the sum of not-none counts of all bins

Return type float

eskapade.analysis.histogram_filling module

Project: Eskapade - A python-based package for data analysis.

Class: HistogramFillerBase

Created: 2017/03/21

Description: Algorithm to fill histogrammar sparse-bin histograms. It is possible to do cleaning of these histograms by rejecting certain keys or removing inconsistent data types. Timestamp columns are converted to nanoseconds before the binning is applied.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.histogram_filling.HistogramFillerBase` (**kwargs)
Bases: `escore.core.element.Link`

Base class link to fill histograms.

It is possible to do after-filling cleaning of these histograms by rejecting certain keys or removing inconsistent data types. Timestamp columns are converted to nanoseconds before the binning is applied. Final histograms are stored in the datastore.

`__init__` (**kwargs)
Initialize link instance.

Store and do basic check on the attributes of link HistogramFillerBase.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **store_key** (*str*) – key of output data to store histograms in data store
- **columns** (*list*) – columns to pick up from input data. (default is all columns)
- **bin_specs** (*dict*) – dictionaries used for rebinning numeric or timestamp columns

Example `bin_specs` dictionary is:

```
>>> bin_specs = {'x': {'bin_width': 1, 'bin_offset': 0},
                 'y': {'bin_edges': [0, 2, 3, 4, 5, 7, 8]}}
```

Parameters

- **var_dtype** (*dict*) – dict of datatypes of the columns to study from dataframe. If not provided, try to determine datatypes directly from dataframe.
- **store_at_finalize** (*bool*) – Store histograms in datastore at `finalize()`, not at `execute()`. Useful when looping over datasets. Default is `False`.
- **dict** (*drop_keys*) – dictionary used for dropping specific keys from bins dictionaries of histograms

Example `drop_keys` dictionary is:

```
>>> drop_keys = {'x': [1, 4, 8, 19],
                 'y': ['apple', 'pear', 'tomato'],
                 'x:y': [(1, 'apple'), (19, 'tomato')]}
```

assert_dataframe (*df*)
Check that input data is a filled pandas data frame.

Parameters *df* – input (pandas) data frame

categorize_columns (*df*)
Categorize columns of dataframe by data type.

Parameters *df* – input (pandas) data frame

drop_requested_keys (*name, counts*)

Drop requested keys from counts dictionary.

Parameters

- **name** (*string*) – key of drop_keys dict to get array of keys to be dropped
- **counts** (*dict*) – counts dictionary to drop specific keys from

Returns count dict without dropped keys

execute ()

Execute the link.

Execute() four things:

- check presence and data type of requested columns
- timestamp variables are converted to nanosec (integers)
- do the actual value counting based on categories and created indices
- then convert to histograms and add to datastore

fill_histogram (*idf, c*)

Fill input histogram with column(s) of input dataframe.

Parameters

- **idf** – input data frame used for filling histogram
- **c** (*list*) – histogram column(s)

finalize ()

Finalize the link.

Store Histograms here, if requested.

get_all_columns (*data*)

Retrieve all columns / keys from input data.

Parameters **data** – input data sample (pandas dataframe or dict)

Returns list of columns

Return type list

get_data_type (*df, col*)

Get data type of dataframe column.

Parameters

- **df** – input data frame
- **col** (*str*) – column

initialize ()

Initialize the link.

process_and_store ()

Store (and possibly process) histogram objects.

process_columns (*df*)

Process columns before histogram filling.

Specifically, convert timestamp columns to integers

Parameters **df** – input (pandas) data frame

Returns output (pandas) data frame with converted timestamp columns

Return type pandas DataFrame

`eskapade.analysis.histogram_filling.only_bool(val)`

Pass input value or array only if it is a bool.

Parameters `val` – value to be evaluated

Returns evaluated value

Return type `np.bool` or `np.ndarray`

`eskapade.analysis.histogram_filling.only_float(val)`

Pass input `val` value or array only if it is a float.

Parameters `val` – value to be evaluated

Returns evaluated value

Return type `np.float64` or `np.ndarray`

`eskapade.analysis.histogram_filling.only_int(val)`

Pass input `val` value or array only if it is an integer.

Parameters `val` – value to be evaluated

Returns evaluated value

Return type `np.int64` or `np.ndarray`

`eskapade.analysis.histogram_filling.only_str(val)`

Pass input value or array only if it is a string.

Parameters `val` – value to be evaluated

Returns evaluated value

Return type `str` or `np.ndarray`

`eskapade.analysis.histogram_filling.to_ns(x)`

Convert input timestamps to nanoseconds (integers).

Parameters `x` – value to be converted

Returns converted value

Return type `int`

`eskapade.analysis.histogram_filling.to_str(val)`

Convert input to (array of) string(s).

Parameters `val` – value to be converted

Returns converted value

Return type `str` or `np.ndarray`

`eskapade.analysis.histogram_filling.value_to_bin_center(val, **kwargs)`

Convert value to bin center.

Convert a numeric or timestamp column to a common bin center value.

Parameters

- **bin_width** – bin_width value needed to convert column to a common bin center value
- **bin_offset** – bin_offset value needed to convert column to a common bin center value

`eskapade.analysis.histogram_filling.value_to_bin_index(val, **kwargs)`

Convert value to bin index.

Convert a numeric or timestamp column to an integer bin index.

Parameters

- **bin_width** – bin_width value needed to convert column to an integer bin index
- **bin_offset** – bin_offset value needed to convert column to an integer bin index

eskapade.analysis.statistics module

Project: Eskapade - A python-based package for data analysis.

Classes: ArrayStats, GroupByStats

Created: 2017/03/21

Description: Summary of an array.

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

class `eskapade.analysis.statistics.ArrayStats` (*data*, *col_name*, *weights=None*, *unit=""*, *label=""*)

Bases: object

Create summary of an array.

Class to calculate statistics (mean, standard deviation, percentiles, etc.) and create a histogram of values in an array. The statistics can be returned as values in a dictionary, a printable string, or as a LaTeX string.

__init__ (*data*, *col_name*, *weights=None*, *unit=""*, *label=""*)

Initialize for a single column in data frame.

Parameters

- **data** (*(keys of) dict*) – Input array
- **col_name** – column name
- **weights** (*string (column of data)*) – Input array (default None)
- **unit** – Unit of column
- **label** (*str*) – Label to describe column variable

Raises TypeError

create_mpv_stat ()

Compute most probable value from histogram.

This function computes the most probable value based on the histogram from `make_histogram()`, and adds it to the statistics.

create_stats ()

Compute statistical properties of column variable.

This function computes the statistical properties of values in the specified column. It is called by other functions that use the resulting figures to create a statistical overview.

get_col_props ()

Get column properties.

Returns dict Column properties

get_latex_table (*get_stats=None, latex=True*)
Get LaTeX code string for table of stats values.

Parameters

- **get_stats** (*list*) – List of statistics that you want to filter on. (default None (all stats)) Available stats are: ‘count’, ‘filled’, ‘distinct’, ‘mean’, ‘std’, ‘min’, ‘max’, ‘p05’, ‘p16’, ‘p50’, ‘p84’, ‘p95’, ‘p99’
- **latex** (*bool*) – LaTeX output or list output (default True)

Returns str LaTeX code snippet

get_print_stats (*to_output=False*)
Get statistics in printable form.

Parameters **to_output** (*bool*) – Print statistics to output stream?

Returns str Printable statistics string

get_x_label ()
Get x label.

logger

A logger that emits log messages to an observer.

The logger can be instantiated as a module or class attribute, e.g.

```
>>> logger = Logger()
>>> logger.info("I'm a module logger attribute.")
>>>
>>> class Point(object):
>>>     logger = Logger()
>>>
>>>     def __init__(self, x = 0.0, y = 0.0):
>>>         Point.logger.debug('Initializing {point} with x = {x} y = {y}',
↳point=Point, x=x, y=y)
>>>         self._x = x
>>>         self._y = y
>>>
>>>     @property
>>>     def x(self):
>>>         self.logger.debug('Getting property x = {point._x}', point=self)
>>>         return self._x
>>>
>>>     @x.setter
>>>     def x(self, x):
>>>         self.logger.debug('Setting property y = {point._x}', point=self)
>>>         self._x = x
>>>
>>>     @property
>>>     def y(self):
>>>         self.logger.debug('Getting property y = {point._y}', point=self)
>>>         return self._y
>>>
>>>     @y.setter
>>>     def y(self, y):
>>>         self.logger.debug('Setting property y = {point._y}', point=self)
>>>         self._y = y
```

(continues on next page)

(continued from previous page)

```

>>>
>>> a_point = Point(1, 2)
>>>
>>> logger.info('p_x = {point.x} p_y = {point.y}', point=a_point)
>>> logger.log_level = LogLevel.DEBUG
>>> logger.info('p_x = {point.x} p_y = {point.y}', point=a_point)

```

The logger uses PEP-3101 (Advanced String Formatting) with named placeholders, see <<https://www.python.org/dev/peps/pep-3101/>> and <<https://pyformat.info/>> for more details and examples.

Furthermore, logging events are only formatted and evaluated for logging levels that are enabled. So, there's no need to check the logging level before logging. It's also efficient.

make_histogram (*var_bins=30, var_range=None, bin_edges=None, create_mpv_stat=True*)
Create histogram of column values.

Parameters

- **var_bins** (*int*) – Number of histogram bins
- **var_range** (*tuple*) – Range of histogram variable
- **bin_edges** (*list*) – predefined bin edges to use for histogram. Overrides `var_bins`.

class `eskapade.analysis.statistics.GroupByStats` (*data, col_name, groupby=None, weights=None, unit="", label=""*)

Bases: `eskapade.analysis.statistics.ArrayStats`

Create summary of an array in groups.

__init__ (*data, col_name, groupby=None, weights=None, unit="", label=""*)
Initialize for a single column in dataframe.

Parameters

- **data** (*(keys of) dict*) – Input array
- **col_name** – column name
- **weights** (*string (column of data)*) – Input array (default None)
- **unit** – Unit of column
- **label** (*str*) – Label to describe column variable
- **groupby** – column name

Raises `TypeError`

get_latex_table (*get_stats=None*)
Get LaTeX code string for group-by table of stats values.

Parameters **get_stats** (*list*) – same as `ArrayStats.get_latex_table` `get_stats` key word.

Returns **str** LaTeX code snippet

`eskapade.analysis.statistics.get_col_props` (*var_type*)
Get column properties.

Returns **dict** Column properties

`eskapade.analysis.statistics.weighted_quantile` (*data, weights=None, probability=0.5*)
Compute the weighted quantile of a 1D numpy array.

Weighted quantiles, inspired by: <https://github.com/nudomarinero/wquantiles/blob/master/wquantiles.py> written by Jose Sabater Here updated to return multiple quantiles in one go. Now also works when weight is None.

Parameters

- **data** (*ndarray*) – input array (one dimension).
- **weights** (*ndarray*) – array with the weights of the same size of *data*.
- **probability** (*ndarray*) – array of quantiles to compute. Each probability must have a value between 0 and 1.

Returns list of the output value(s).

Module contents

eskapade.core package

Submodules

eskapade.core.definitions module

eskapade.core.element module

eskapade.core.exceptions module

eskapade.core.execution module

eskapade.core.meta module

eskapade.core.mixin module

eskapade.core.persistence module

eskapade.core.process_manager module

eskapade.core.process_services module

eskapade.core.run_utils module

Module contents

eskapade.core_ops package

Subpackages

eskapade.core_ops.links package

Submodules

eskapade.core_ops.links.apply module

eskapade.core_ops.links.assert_in_ds module

eskapade.core_ops.links.break_link module

eskapade.core_ops.links.ds_object_deleter module

eskapade.core_ops.links.ds_to_ds module

eskapade.core_ops.links.event_looper module

eskapade.core_ops.links.hello_world module

eskapade.core_ops.links.import_data_store module

Class: FixPandasDataFrame

Created: 2017/04/07

Description: Link for fixing dirty pandas dataframe with inconsistent datatypes See example in: [tutorials/esk501_fix_pandas_dataframe.py](https://github.com/escor/core/blob/master/tutorials/esk501_fix_pandas_dataframe.py)

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

```
class eskapade.data_quality.links.fix_pandas_dataframe.FixPandasDataFrame (**kwargs)
    Bases: escore.core.element.Link
```

Fix dirty Pandas dataframe with inconsistent datatypes.

Default settings perform the following clean-up steps on an input dataframe:

- Fix all column names. E.g. remove punctuation and strange characters, and convert spaces to underscores.
- Check for various possible nans in the dataset, then make all nans consistent by turning them into `numpy.nan` (= float)
- Per column, assess dynamically the most consistent datatype (ignoring all nans in that column). E.g. bool, int, float, datetime64, string.
- Per column, make the data types of all rows consistent, by using the identified (or imposed) data type (by default ignoring all nans)

Boolean columns with contamination get converted to string columns by default. Optionally, they can be converted to integer columns as well.

The FixPandasDataFrame link can be used in a dataframe loop, in which case any data type assessed per column in the first dataframe iteration will be used for the next dataframes as well.

The default settings should work pretty well in many circumstances, by can be configured pretty flexibly. Optionally:

- Instead of dynamically assessed, the data type can also be imposed per column
- All nans in a column can be converted to a value consistent with the data type of that column. E.g. for integer columns, nan -> -999
- An alternative nan can be set per column and datatype
- Modifications can be applied inplace, i.e. directly to the input dataframe

```
__init__ (**kwargs)
    Initialize link instance.
```

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **copy_columns_from_df** (*bool*) – if true, copy all columns from the dataframe (default is true)
- **original_columns** (*list*) – original (unfixed) column names to pick up from input data (required if `copy_columns_from_df` is set to false)
- **contaminated_columns** (*list*) – (original) columns that are known to have mistakes and that should be fixed (optional)
- **fix_column_names** (*bool*) – if true, fix column names (default is true)

- **strip_hive_prefix** (*bool*) – if true, strip table-name (hive) prefix from column names, e.g. table.bla -> bla (default is false)
- **convert_inconsistent_dtypes** (*bool*) – fix column datatypes in case of data type inconsistencies in rows (default is true)
- **var_dtype** (*dict*) – dict forcing columns to certain datatypes, e.g. {'A': int} (optional)
- **var_convert_inconsistent_dtypes** (*dict*) – dict allowing one to overwrite if certain columns datatypes should be fixed, e.g. {'A': False} (optional)
- **var_convert_func** (*dict*) – dict with datatype conversion functions for certain columns
- **check_nan_func** – boolean return function to check for nans in columns. (default is None, in which case a standard checker function gets picked up)
- **convert_inconsistent_nans** (*bool*) – if true, convert all nans to data type consistent with rest of column (default is false)
- **var_convert_inconsistent_nans** (*dict*) – dict allowing one to overwrite if certain column nans should be fixed, e.g. {'A': False} (optional)
- **var_nan** (*dict*) – dict with nans for certain columns (optional)
- **nan_dtype_map** (*dict*) – dictionary of nans for given data types, e.g. { int: -999 }
- **nan_default** – default nan value to which all nans found get converted (default is numpy.nan)
- **var_bool_to_int** (*list*) – convert boolean column to int (default is conversion of boolean to string)
- **inplace** (*bool*) – replace original columns; overwrites store_key to read_key (default is False)
- **store_key** (*str*) – key of output data to store in data store
- **drop_dup_rec** (*bool*) – if true, drop duplicate records from data frame after other fixes (default is false)
- **strip_string_columns** (*bool*) – if true, apply strip command to string columns (default is true)
- **cleanup_string_columns** (*list*) – boolean or list. apply cleaning-up to list of selected or all string columns. More aggressive than strip. Default is empty (= false).

execute ()

Execute the link.

Fixing the Pandas dataframe consists of four steps:

- Fix all column names. E.g. remove punctuation and strange characters, and convert spaces to underscores.
- Check existing nans in that dataset, and make all nans consistent, for easy conversion later on.
- Assess most consistent datatype for each column (ignoring all nans)
- Make data types in each row consistent (by default ignoring all nans)

initialize ()

Initialize the link.

`eskapade.data_quality.links.fix_pandas_dataframe.determine_preferred_dtype (dtype_cnt)`
Determine preferred column data type.

Module contents

class `eskapade.data_quality.links.FixPandasDataFrame` (**kwargs)

Bases: `escore.core.element.Link`

Fix dirty Pandas dataframe with inconsistent datatypes.

Default settings perform the following clean-up steps on an input dataframe:

- Fix all column names. E.g. remove punctuation and strange characters, and convert spaces to underscores.
- Check for various possible nans in the dataset, then make all nans consistent by turning them into `numpy.nan` (= float)
- Per column, assess dynamically the most consistent datatype (ignoring all nans in that column). E.g. bool, int, float, `datetime64`, string.
- Per column, make the data types of all rows consistent, by using the identified (or imposed) data type (by default ignoring all nans)

Boolean columns with contamination get converted to string columns by default. Optionally, they can be converted to integer columns as well.

The `FixPandasDataFrame` link can be used in a dataframe loop, in which case any data type assessed per column in the first dataframe iteration will be used for the next dataframes as well.

The default settings should work pretty well in many circumstances, but can be configured pretty flexibly. Optionally:

- Instead of dynamically assessed, the data type can also be imposed per column
- All nans in a column can be converted to a value consistent with the data type of that column. E.g. for integer columns, `nan` -> `-999`
- An alternative nan can be set per column and datatype
- Modifications can be applied inplace, i.e. directly to the input dataframe

__init__ (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **copy_columns_from_df** (*bool*) – if true, copy all columns from the dataframe (default is true)
- **original_columns** (*list*) – original (unfixed) column names to pick up from input data (required if `copy_columns_from_df` is set to false)
- **contaminated_columns** (*list*) – (original) columns that are known to have mistakes and that should be fixed (optional)
- **fix_column_names** (*bool*) – if true, fix column names (default is true)
- **strip_hive_prefix** (*bool*) – if true, strip table-name (hive) prefix from column names, e.g. `table.bla` -> `bla` (default is false)
- **convert_inconsistent_dtypes** (*bool*) – fix column datatypes in case of data type inconsistencies in rows (default is true)
- **var_dtype** (*dict*) – dict forcing columns to certain datatypes, e.g. `{‘A’: int}` (optional)

- **var_convert_inconsistent_dtypes** (*dict*) – dict allowing one to overwrite if certain columns datatypes should be fixed, e.g. {'A': False} (optional)
- **var_convert_func** (*dict*) – dict with datatype conversion functions for certain columns
- **check_nan_func** – boolean return function to check for nans in columns. (default is None, in which case a standard checker function gets picked up)
- **convert_inconsistent_nans** (*bool*) – if true, convert all nans to data type consistent with rest of column (default is false)
- **var_convert_inconsistent_nans** (*dict*) – dict allowing one to overwrite if certain column nans should be fixed, e.g. {'A': False} (optional)
- **var_nan** (*dict*) – dict with nans for certain columns (optional)
- **nan_dtype_map** (*dict*) – dictionary of nans for given data types, e.g. { int: -999 }
- **nan_default** – default nan value to which all nans found get converted (default is numpy.nan)
- **var_bool_to_int** (*list*) – convert boolean column to int (default is conversion of boolean to string)
- **inplace** (*bool*) – replace original columns; overwrites store_key to read_key (default is False)
- **store_key** (*str*) – key of output data to store in data store
- **drop_dup_rec** (*bool*) – if true, drop duplicate records from data frame after other fixes (default is false)
- **strip_string_columns** (*bool*) – if true, apply strip command to string columns (default is true)
- **cleanup_string_columns** (*list*) – boolean or list. apply cleaning-up to list of selected or all string columns. More aggressive than strip. Default is empty (= false).

execute ()

Execute the link.

Fixing the Pandas dataframe consists of four steps:

- Fix all column names. E.g. remove punctuation and strange characters, and convert spaces to underscores.
- Check existing nans in that dataset, and make all nans consistent, for easy conversion later on.
- Assess most consistent datatype for each column (ignoring all nans)
- Make data types in each row consistent (by default ignoring all nans)

initialize ()

Initialize the link.

Submodules

eskapade.data_quality.dq_helper module

Project: Eskapade - A Python-based package for data analysis.

Module: data_quality.dq_helper

Created: 2017/04/11

Description: Data-quality helper functions

Authors: KPMG Advanced Analytics & Big Data team, Amstelveen, The Netherlands

Redistribution and use in source and binary forms, with or without modification, are permitted according to the terms listed in the file LICENSE.

`eskapade.data_quality.dq_helper.bool_to_int(val, **kwargs)`

Convert input boolean to int.

Parameters `val` – value to be evaluated

Returns evaluated value

Return type `np.int64`

`eskapade.data_quality.dq_helper.bool_to_str(val, **kwargs)`

Convert input boolean to str.

Parameters `val` – value to be evaluated

Returns evaluated value

Return type `str`

`eskapade.data_quality.dq_helper.check_nan(val)`

Check input value for not a number.

Parameters `val` – value to be checked for nan

Returns true if nan

Return type `bool`

`eskapade.data_quality.dq_helper.cleanup_string(col)`

Cleanup input string.

Parameters `col` – string to be cleaned up

Returns cleaned up string

Return type `str`

`eskapade.data_quality.dq_helper.convert(val)`

Convert input to interpreted data type.

Parameters `val` – value to be interpreted

Returns interpreted value

`eskapade.data_quality.dq_helper.to_date_time(val, **kwargs)`

Convert input to `numpy.datetime64`.

Parameters `val` – value to be evaluated

Returns evaluated value

Return type `numpy.datetime64`

`eskapade.data_quality.dq_helper.to_float(val, **kwargs)`

Convert input to float.

Parameters `val` – value to be evaluated

Returns evaluated value

Return type `np.float64`

`eskapade.data_quality.dq_helper.to_int(val, **kwargs)`

Convert input to int.

Parameters `val` – value to be evaluated

Returns evaluated value

Return type `np.int64`

`eskapade.data_quality.dq_helper.to_str(val, **kwargs)`

Convert input to string.

Parameters `val` – value to be converted

Returns converted value

Return type `str`

Module contents

eskapade.logger package

Module contents

eskapade.visualization package

Subpackages

eskapade.visualization.links package

Submodules

eskapade.visualization.links.correlation_summary module

eskapade.visualization.links.df_boxplot module

eskapade.visualization.links.df_summary module

Module contents

Submodules

eskapade.visualization.vis_utils module

Module contents

Submodules

eskapade.entry_points module

eskapade.exceptions module

eskapade.helpers module

eskapade.resources module

eskapade.utils module

eskapade.version module

THIS FILE IS AUTO-GENERATED BY ESKAPADE SETUP.PY.

Module contents

6.10 Appendices

6.10.1 Miscellaneous

Collection of miscellaneous Eskapade related items.

- See *Migration Tips* to migrate between Eskapade versions.
- See *macOS* to get started with Eskapade on a mac.

Migration Tips

From version 0.8 to 0.9

In Eskapade v0.9 the core functionality has been migrated to the separate package Eskapade-Core. We have tried to make this transition as seamless as possible, but you may well run into any migration issues. In case you do below we list the changes needed to migrate from Eskapade version 0.8 to version 0.9.

- Whenever a line with `import eskapade` fails, simply replace `eskapade` with `escore`.
- For example: from `eskapade import core` change to `from escore import core`

That's it.

From version 0.6 to 0.7

Below we list the API changes needed to migrate from Eskapade version 0.6 to version 0.7.

Links

- Process manager definition:
 - `proc_mgr` . change to `process_manager` .
 - `ProcessManager` change to `process_manager`
 - Delete line: `proc_mgr = ProcessManager()`
- Logger:
 - Change `log()` . to `logger` .

Macros

- Process manager definition:
 - `proc_mgr` . change to `process_manager` .
 - `ProcessManager` change to `process_manager`
 - Delete line: `proc_mgr = ProcessManager()`
- Logger:

- import logging **change to** from eskapade.logger import Logger, LogLevel
- log. **change to** logger.
- log = logging.getLogger('macro.cpf_analysis') **change to** logger = Logger()
- logging **change to** LogLevel

- Settings:

Remove `os.environ['WORKDIRROOT']`, since the environment variable WORKDIRROOT is no longer defined, define explicitly the data and macro paths, or execute the macros and tests from the root directory of the project, resulting in something like:

- settings['resultsDir'] = os.getcwd() + 'es_results'
- settings['macrosDir'] = os.getcwd() + 'es_macros'
- settings['dataDir'] = os.getcwd() + 'data'

- Chain definition in macros:

- To import the Chain object add from eskapade import Chain
- **Change** `process_manager.add_chain('chain_name')` **to** `<chain_name> = Chain('chain_name')`
- `process_manager.get_chain('ReadCSV').add_link` **to** `<chain_name>.add`

Tests

- Process manager definition:

- **Change** `ProcessManager()` **to** `process_manager`
- **Change** `process_manager.get_chain` **to** `process_manager.get`

- Settings:

Remove `os.environ['WORKDIRROOT']`, since the environment variable WORKDIRROOT is no longer defined, define explicitly the data and macro paths, or execute the macros and tests from the root directory of the project, resulting in something like:

- settings['macrosDir'] = os.getcwd() + '/es_macros'
- settings['dataDir'] = os.getcwd() + '/data'

- StatusCode:

- **Change** `status.isSkipChain()` **to** `status.is_skip_chain()`

macOS

To install Eskapade on macOS there are basically four steps:

- Setting up Python 3.6
- Setting up Apache Spark 2.x
- Setting up ROOT 6.10.08
- Setting up Eskapade

Note: This installation guide is written for [macOS High Sierra with Homebrew](#), and [fish](#).

Setting up Python 3.6

Homebrew provides Python 3.6 for macOS:

```
$ brew install python3
```

To create an isolated Python installation use `virtualenv`:

```
$ virtualenv venv/eskapade --python=python3 --system-site-packages
```

Each time a new terminal is started, set up the virtual python environment:

```
$ . ~/venv/eskapade/bin/activate.fish
```

Setting up ROOT 6

Clone ROOT from the git repository:

```
git clone http://root.cern.ch/git/root.git
cd root
git checkout -b v6-10-08 v6-10-08
```

Then compile it with the additional flags to ensure the desired functionality:

```
$ mkdir ~/root_v06-10-08_p36m && cd ~/root_v06-10-08_p36m
$ cmake -Dfftw3=ON -Dmathmore=ON -Dminuit2=ON -Droofit=ON -Dtmva=ON -Dsoversion=ON -
↳ Dthread=ON -Dpython3=ON -DPYTHON_EXECUTABLE=/usr/local/opt/python3/Frameworks/
↳ Python.framework/Versions/3.6/bin/python3.6m -DPYTHON_INCLUDE_DIR=/usr/local/opt/
↳ python3/Frameworks/Python.framework/Versions/3.6/include/python3.6m/ -DPYTHON_
↳ LIBRARY=/usr/local/opt/python3/Frameworks/Python.framework/Versions/3.6/lib/
↳ libpython3.6m.dylib $HOME/root
$ cmake --build . -- -j7
```

PS: make sure all the flags are picked up correctly (for example, `-Dfftw3` requires `fftw` to be installed with Homebrew).

To setup the ROOT environment each time a new shell is started, set the following environment variables:

```
set -xg ROOTSYS "$HOME/root_v06-10-08_p36m"
set -xg PATH $ROOTSYS/bin $PATH
set -xg LD_LIBRARY_PATH "$ROOTSYS/lib:$LD_LIBRARY_PATH"
set -xg DYLD_LIBRARY_PATH "$ROOTSYS/lib:$DYLD_LIBRARY_PATH"
set -xg LIBPATH "$ROOTSYS/lib:$LIBPATH"
set -xg SHLIB_PATH "$ROOTSYS/lib:$SHLIB_PATH"
set -xg PYTHONPATH "$ROOTSYS/lib:$PYTHONPATH"
```

Note that for bash shells this can be done by sourcing the script in `root_v06-10-08_p36m/bin/thisroot.sh`.

Finally, install the Python packages for ROOT bindings:

```
$ pip install rootpy==1.0.1 root-numpy=4.7.3
```

Setting up Apache Spark 2.x

Apache Spark is provided through Homebrew:

```
$ brew install apache-spark
```

The `py4j` package is needed to support access to Java objects from Python:

```
$ pip install py4j==0.10.4
```

To set up the Spark environment each time a new terminal is started set:

```
set -xg SPARK_HOME (brew --prefix apache-spark)/libexec
set -xg SPARK_LOCAL_HOSTNAME "localhost"
set -xg PYTHONPATH "$SPARK_HOME/python:$PYTHONPATH"
```

Setting up Eskapade

To set up the Eskapade environment (Python, Spark, ROOT) each time a new terminal is started, source a shell script (e.g. `setup_eskapade.fish`) that contains set the environment variables as described above:

```
# --- setup Python
. ~/venv/eskapade/bin/activate.fish

# --- setup ROOT
set -xg ROOTSYS "${HOME}/root_v06-10-08_p36m"
set -xg PATH $ROOTSYS/bin $PATH
set -xg LD_LIBRARY_PATH "$ROOTSYS/lib:$LD_LIBRARY_PATH"
set -xg DYLD_LIBRARY_PATH "$ROOTSYS/lib:$DYLD_LIBRARY_PATH"
set -xg LIBPATH "$ROOTSYS/lib:$LIBPATH"
set -xg SHLIB_PATH "$ROOTSYS/lib:$SHLIB_PATH"
set -xg PYTHONPATH "$ROOTSYS/lib:$PYTHONPATH"

# --- setup Spark
set -xg SPARK_HOME (brew --prefix apache-spark)/libexec
set -xg SPARK_LOCAL_HOSTNAME "localhost"
set -xg PYTHONPATH "$SPARK_HOME/python:$PYTHONPATH"

# --- setup Eskapade
cd /path/to/eskapade
```

Finally, install Eskapade (and it's dependencies) by simply running:

```
$ pip install Eskapade
```

6.11 Indices and tables

- [genindex](#)
- [modindex](#)

e

eskapade, 100
eskapade.analysis, 92
eskapade.analysis.correlation, 72
eskapade.analysis.datetime, 73
eskapade.analysis.histogram, 75
eskapade.analysis.histogram_filling, 83
eskapade.analysis.links, 63
eskapade.analysis.links.apply_func_to_df,
49
eskapade.analysis.links.apply_selection_to_df,
50
eskapade.analysis.links.basic_generator,
51
eskapade.analysis.links.df_concatenator,
51
eskapade.analysis.links.df_merger, 52
eskapade.analysis.links.histogrammar_filler,
53
eskapade.analysis.links.random_sample_splitter,
54
eskapade.analysis.links.read_to_df, 55
eskapade.analysis.links.record_factorizer,
57
eskapade.analysis.links.record_vectorizer,
58
eskapade.analysis.links.value_counter,
59
eskapade.analysis.links.write_from_df,
61
eskapade.analysis.statistics, 87
eskapade.data_quality, 99
eskapade.data_quality.dq_helper, 96
eskapade.data_quality.links, 95
eskapade.data_quality.links.fix_pandas_dataframe,
92
eskapade.logger, 99
eskapade.version, 99

Symbols

- `__init__()` (eskapade.analysis.datetime.FreqTimePeriod method), 73
- `__init__()` (eskapade.analysis.datetime.TimePeriod method), 73
- `__init__()` (eskapade.analysis.datetime.UniformTsTimePeriod method), 75
- `__init__()` (eskapade.analysis.histogram.BinningUtil method), 75
- `__init__()` (eskapade.analysis.histogram.Histogram method), 77
- `__init__()` (eskapade.analysis.histogram.ValueCounts method), 81
- `__init__()` (eskapade.analysis.histogram_filling.HistogramFillerBase method), 84
- `__init__()` (eskapade.analysis.links.ApplyFuncToDf method), 63
- `__init__()` (eskapade.analysis.links.ApplySelectionToDf method), 63
- `__init__()` (eskapade.analysis.links.BasicGenerator method), 64
- `__init__()` (eskapade.analysis.links.DfConcatenator method), 64
- `__init__()` (eskapade.analysis.links.DfMerger method), 65
- `__init__()` (eskapade.analysis.links.HistogrammarFiller method), 65
- `__init__()` (eskapade.analysis.links.RandomSampleSplitter method), 67
- `__init__()` (eskapade.analysis.links.ReadToDf method), 67
- `__init__()` (eskapade.analysis.links.RecordFactorizer method), 68
- `__init__()` (eskapade.analysis.links.RecordVectorizer method), 69
- `__init__()` (eskapade.analysis.links.ValueCounter method), 70
- `__init__()` (eskapade.analysis.links.WriteFromDf method), 71
- `__init__()` (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 49
- `__init__()` (eskapade.analysis.links.apply_selection_to_df.ApplySelectionToDf method), 50
- `__init__()` (eskapade.analysis.links.basic_generator.BasicGenerator method), 51
- `__init__()` (eskapade.analysis.links.df_concatenator.DfConcatenator method), 52
- `__init__()` (eskapade.analysis.links.df_merger.DfMerger method), 52
- `__init__()` (eskapade.analysis.links.histogrammar_filler.HistogrammarFiller method), 53
- `__init__()` (eskapade.analysis.links.random_sample_splitter.RandomSampleSplitter method), 55
- `__init__()` (eskapade.analysis.links.read_to_df.ReadToDf method), 56
- `__init__()` (eskapade.analysis.links.record_factorizer.RecordFactorizer method), 57
- `__init__()` (eskapade.analysis.links.record_vectorizer.RecordVectorizer method), 59
- `__init__()` (eskapade.analysis.links.value_counter.ValueCounter method), 60
- `__init__()` (eskapade.analysis.links.write_from_df.WriteFromDf method), 62
- `__init__()` (eskapade.analysis.statistics.ArrayStats method), 87
- `__init__()` (eskapade.analysis.statistics.GroupByStats method), 89
- `__init__()` (eskapade.data_quality.links.FixPandasDataFrame method), 95
- `__init__()` (eskapade.data_quality.links.fix_pandas_dataframe.FixPandasDataFrame method), 93

A

- `add_apply_func()` (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 50
- `add_apply_func()` (eskapade.analysis.links.ApplyFuncToDf method), 63

ApplyFuncToDf (class in eskapade.analysis.links), 63
 ApplyFuncToDf (class in eskapade.analysis.links.apply_func_to_df), 49
 ApplySelectionToDf (class in eskapade.analysis.links), 63
 ApplySelectionToDf (class in eskapade.analysis.links.apply_selection_to_df), 50
 ArrayStats (class in eskapade.analysis.statistics), 87
 assert_dataframe() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 84

B

BasicGenerator (class in eskapade.analysis.links), 64
 BasicGenerator (class in eskapade.analysis.links.basic_generator), 51
 bin_centers() (eskapade.analysis.histogram.Histogram method), 77
 bin_edges() (eskapade.analysis.histogram.Histogram method), 77
 bin_entries() (eskapade.analysis.histogram.Histogram method), 77
 bin_labels() (eskapade.analysis.histogram.Histogram method), 78
 bin_specs (eskapade.analysis.histogram.BinningUtil attribute), 76
 BinningUtil (class in eskapade.analysis.histogram), 75
 bool_to_int() (in module eskapade.data_quality.dq_helper), 97
 bool_to_str() (in module eskapade.data_quality.dq_helper), 97

C

calculate_correlations() (in module eskapade.analysis.correlation), 72
 categorize_columns() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 84
 check_nan() (in module eskapade.data_quality.dq_helper), 97
 chi2_from_rho() (in module eskapade.analysis.correlation), 72
 cleanup_string() (in module eskapade.data_quality.dq_helper), 97
 combine_hists() (eskapade.analysis.histogram.Histogram class method), 78
 construct_empty_hist() (eskapade.analysis.links.histogrammar_filler.HistogrammarFiller method), 54
 construct_empty_hist() (eskapade.analysis.links.HistogrammarFiller method), 66

convert() (in module eskapade.data_quality.dq_helper), 97
 copy() (eskapade.analysis.histogram.Histogram method), 78
 count() (eskapade.analysis.histogram.ValueCounts method), 82
 counts (eskapade.analysis.histogram.ValueCounts attribute), 82
 create_mpv_stat() (eskapade.analysis.statistics.ArrayStats method), 87
 create_stats() (eskapade.analysis.statistics.ArrayStats method), 87
 create_sub_counts() (eskapade.analysis.histogram.ValueCounts method), 82

D

datatype (eskapade.analysis.histogram.Histogram attribute), 78
 determine_preferred_dtype() (in module eskapade.data_quality.links.fix_pandas_dataframe), 94
 DfConcatenator (class in eskapade.analysis.links), 64
 DfConcatenator (class in eskapade.analysis.links.df_concatenator), 51
 DfMerger (class in eskapade.analysis.links), 65
 DfMerger (class in eskapade.analysis.links.df_merger), 52
 drop_inconsistent_keys() (eskapade.analysis.links.value_counter.ValueCounter method), 61
 drop_inconsistent_keys() (eskapade.analysis.links.ValueCounter method), 71
 drop_requested_keys() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 84
 dt_string() (eskapade.analysis.datetime.FreqTimePeriod method), 73

E

eskapade (module), 100
 eskapade.analysis (module), 92
 eskapade.analysis.correlation (module), 72
 eskapade.analysis.datetime (module), 73
 eskapade.analysis.histogram (module), 75
 eskapade.analysis.histogram_filling (module), 83
 eskapade.analysis.links (module), 63
 eskapade.analysis.links.apply_func_to_df (module), 49
 eskapade.analysis.links.apply_selection_to_df (module), 50
 eskapade.analysis.links.basic_generator (module), 51
 eskapade.analysis.links.df_concatenator (module), 51

- eskapade.analysis.links.df_merger (module), 52
 - eskapade.analysis.links.histogrammar_filler (module), 53
 - eskapade.analysis.links.random_sample_splitter (module), 54
 - eskapade.analysis.links.read_to_df (module), 55
 - eskapade.analysis.links.record_factorizer (module), 57
 - eskapade.analysis.links.record_vectorizer (module), 58
 - eskapade.analysis.links.value_counter (module), 59
 - eskapade.analysis.links.write_from_df (module), 61
 - eskapade.analysis.statistics (module), 87
 - eskapade.data_quality (module), 99
 - eskapade.data_quality.dq_helper (module), 96
 - eskapade.data_quality.links (module), 95
 - eskapade.data_quality.links.fix_pandas_dataframe (module), 92
 - eskapade.logger (module), 99
 - eskapade.version (module), 99
 - execute() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 85
 - execute() (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 50
 - execute() (eskapade.analysis.links.apply_selection_to_df.ApplySelectionToDf method), 50
 - execute() (eskapade.analysis.links.ApplyFuncToDf method), 63
 - execute() (eskapade.analysis.links.ApplySelectionToDf method), 64
 - execute() (eskapade.analysis.links.basic_generator.BasicGenerator method), 51
 - execute() (eskapade.analysis.links.BasicGenerator method), 64
 - execute() (eskapade.analysis.links.df_concatenator.DfConcatenator method), 52
 - execute() (eskapade.analysis.links.df_merger.DfMerger method), 53
 - execute() (eskapade.analysis.links.DfConcatenator method), 65
 - execute() (eskapade.analysis.links.DfMerger method), 65
 - execute() (eskapade.analysis.links.random_sample_splitter.RandomSampleSplitter method), 55
 - execute() (eskapade.analysis.links.RandomSampleSplitter method), 67
 - execute() (eskapade.analysis.links.read_to_df.ReadToDf method), 56
 - execute() (eskapade.analysis.links.ReadToDf method), 68
 - execute() (eskapade.analysis.links.record_factorizer.RecordFactorizer method), 58
 - execute() (eskapade.analysis.links.record_vectorizer.RecordVectorizer method), 59
 - execute() (eskapade.analysis.links.RecordFactorizer method), 69
 - execute() (eskapade.analysis.links.RecordVectorizer method), 70
 - execute() (eskapade.analysis.links.write_from_df.WriteFromDf method), 62
 - execute() (eskapade.analysis.links.WriteFromDf method), 72
 - execute() (eskapade.data_quality.links.fix_pandas_dataframe.FixPandasDataFrame method), 94
 - execute() (eskapade.data_quality.links.FixPandasDataFrame method), 96
- ## F
- feather_reader() (in module eskapade.analysis.links.read_to_df), 57
 - feather_writer() (in module eskapade.analysis.links.write_from_df), 62
 - fill_histogram() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 85
 - fill_histogram() (eskapade.analysis.links.histogrammar_filler.HistogrammarFillerBase method), 54
 - fill_histogram() (eskapade.analysis.links.HistogrammarFiller method), 66
 - fill_histogram() (eskapade.analysis.links.value_counter.ValueCounter method), 61
 - fill_histogram() (eskapade.analysis.links.ValueCounter method), 71
 - finalize() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 85
 - finalize() (eskapade.analysis.links.value_counter.ValueCounter method), 61
 - finalize() (eskapade.analysis.links.ValueCounter method), 71
 - FixPandasDataFrame (class in eskapade.data_quality.links), 95
 - fix_pandas_dataframe (class in eskapade.data_quality.links.fix_pandas_dataframe), 93
 - freq (eskapade.analysis.datetime.FreqTimePeriod attribute), 73
 - FreqTimePeriod (class in eskapade.analysis.datetime), 73
- ## G
- get_all_columns() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 85
 - get_bin_center() (eskapade.analysis.histogram.BinningUtil method), 76
 - get_bin_count() (eskapade.analysis.histogram.Histogram method), 78
 - get_bin_edges() (eskapade.analysis.histogram.BinningUtil method), 76
 - get_bin_edges_range() (eskapade.analysis.histogram.BinningUtil method), 76
 - get_bin_labels() (eskapade.analysis.histogram.Histogram method), 78

get_bin_range() (eskapade.analysis.histogram.Histogram method), 78	HistogramFillerBase (class in eskapade.analysis.histogram_filling), 84
get_bin_vals() (eskapade.analysis.histogram.Histogram method), 78	HistogrammarFiller (class in eskapade.analysis.links), 65
get_col_props() (eskapade.analysis.statistics.ArrayStats method), 87	HistogrammarFiller (class in eskapade.analysis.links.histogrammar_filler), 53
get_col_props() (in module eskapade.analysis.statistics), 89	
get_data_type() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 85	initialize() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 85
get_hist_val() (eskapade.analysis.histogram.Histogram method), 79	initialize() (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 50
get_latex_table() (eskapade.analysis.statistics.ArrayStats method), 88	initialize() (eskapade.analysis.links.apply_selection_to_df.ApplySelectionToDf method), 51
get_latex_table() (eskapade.analysis.statistics.GroupByStats method), 89	initialize() (eskapade.analysis.links.ApplyFuncToDf method), 63
get_left_bin_edge() (eskapade.analysis.histogram.BinningUtil method), 76	initialize() (eskapade.analysis.links.ApplySelectionToDf method), 64
get_nonone_bin_centers() (eskapade.analysis.histogram.Histogram method), 79	initialize() (eskapade.analysis.links.basic_generator.BasicGenerator method), 51
get_nonone_bin_counts() (eskapade.analysis.histogram.Histogram method), 79	initialize() (eskapade.analysis.links.BasicGenerator method), 64
get_nonone_bin_edges() (eskapade.analysis.histogram.Histogram method), 79	initialize() (eskapade.analysis.links.df_concatenator.DfConcatenator method), 52
get_nonone_bin_range() (eskapade.analysis.histogram.Histogram method), 79	initialize() (eskapade.analysis.links.df_merger.DfMerger method), 53
get_print_stats() (eskapade.analysis.statistics.ArrayStats method), 88	initialize() (eskapade.analysis.links.DfConcatenator method), 65
get_right_bin_edge() (eskapade.analysis.histogram.BinningUtil method), 76	initialize() (eskapade.analysis.links.DfMerger method), 65
get_uniform_bin_edges() (eskapade.analysis.histogram.Histogram method), 79	initialize() (eskapade.analysis.links.random_sample_splitter.RandomSampleSplitter method), 55
get_values() (eskapade.analysis.histogram.ValueCounts method), 82	initialize() (eskapade.analysis.links.RandomSampleSplitter method), 67
get_writer() (in module eskapade.analysis.links.write_from_df), 62	initialize() (eskapade.analysis.links.read_to_df.ReadToDf method), 56
get_x_label() (eskapade.analysis.statistics.ArrayStats method), 88	initialize() (eskapade.analysis.links.ReadToDf method), 68
groupbyapply() (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 50	initialize() (eskapade.analysis.links.record_factorizer.RecordFactorizer method), 58
groupbyapply() (eskapade.analysis.links.ApplyFuncToDf method), 63	initialize() (eskapade.analysis.links.record_vectorizer.RecordVectorizer method), 59
GroupByStats (class in eskapade.analysis.statistics), 89	initialize() (eskapade.analysis.links.RecordFactorizer method), 69
	initialize() (eskapade.analysis.links.RecordVectorizer method), 70
H	initialize() (eskapade.analysis.links.value_counter.ValueCounter method), 61
Histogram (class in eskapade.analysis.histogram), 77	initialize() (eskapade.analysis.links.ValueCounter method), 71
	initialize() (eskapade.analysis.links.write_from_df.WriteFromDf method), 62
	initialize() (eskapade.analysis.links.WriteFromDf method), 72

initialize() (eskapade.data_quality.links.fix_pandas_dataframe_to_fix_pandas_dataframe module method), 94

initialize() (eskapade.data_quality.links.FixPandasDataFrame.only_float() (in module eskapade.analysis.histogram_filling), 86 method), 96

is_finished() (eskapade.analysis.links.read_to_df.ReadToDf.only_int() (in module eskapade.analysis.histogram_filling), 86 method), 56

is_finished() (eskapade.analysis.links.ReadToDf.only_str() (in module eskapade.analysis.histogram_filling), 86 method), 68

K

key (eskapade.analysis.histogram.ValueCounts attribute), 82

L

latest_data_length() (eskapade.analysis.links.read_to_df.ReadToDf method), 56

latest_data_length() (eskapade.analysis.links.ReadToDf method), 68

logger (eskapade.analysis.datetime.TimePeriod attribute), 73

logger (eskapade.analysis.histogram.Histogram attribute), 79

logger (eskapade.analysis.statistics.ArrayStats attribute), 88

M

make_histogram() (eskapade.analysis.statistics.ArrayStats method), 89

N

n_bins (eskapade.analysis.histogram.Histogram attribute), 80

n_dim (eskapade.analysis.histogram.Histogram attribute), 80

nononecounts (eskapade.analysis.histogram.ValueCounts attribute), 82

num_bins (eskapade.analysis.histogram.Histogram attribute), 80

num_bins (eskapade.analysis.histogram.ValueCounts attribute), 82

num_nonone_bins (eskapade.analysis.histogram.ValueCounts attribute), 83

numpy_reader() (in module eskapade.analysis.links.read_to_df), 57

numpy_writer() (in module eskapade.analysis.links.write_from_df), 62

O

offset (eskapade.analysis.datetime.UniformTsTimePeriod attribute), 75

P

parse_date_time() (eskapade.analysis.datetime.TimePeriod class method), 74

parse_time_period() (eskapade.analysis.datetime.TimePeriod class method), 74

period (eskapade.analysis.datetime.UniformTsTimePeriod attribute), 75

period_index() (eskapade.analysis.datetime.FreqTimePeriod method), 73

period_index() (eskapade.analysis.datetime.TimePeriod method), 74

period_index() (eskapade.analysis.datetime.UniformTsTimePeriod method), 75

process_and_store() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 85

process_and_store() (eskapade.analysis.links.value_counter.ValueCounter method), 61

process_and_store() (eskapade.analysis.links.ValueCounter method), 71

process_columns() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 85

process_columns() (eskapade.analysis.links.value_counter.ValueCounter method), 61

process_columns() (eskapade.analysis.links.ValueCounter method), 71

process_counts() (eskapade.analysis.histogram.ValueCounts method), 83

R

RandomSampleSplitter (class in eskapade.analysis.links), 67

RandomSampleSplitter (class in eskapade.analysis.links.random_sample_splitter), 55

ReadToDf (class in eskapade.analysis.links), 67

ReadToDf (class in eskapade.analysis.links.read_to_df), 55

record_vectorizer() (in module eskapade.analysis.links.record_vectorizer), 59

RecordFactorizer (class in eskapade.analysis.links), 68

RecordFactorizer (class in eskapade.analysis.links.record_factorizer), 57

RecordVectorizer (class in eskapade.analysis.links), 69

RecordVectorizer (class in eskapade.analysis.links.record_vectorizer), 58

remove_keys_of_inconsistent_type() (eskapade.analysis.histogram.Histogram method), 81

remove_keys_of_inconsistent_type() (eskapade.analysis.histogram.ValueCounts method), 83

rho_from_chi2() (in module eskapade.analysis.correlation), 73

S

set_chunk_size() (eskapade.analysis.links.read_to_df.ReadToDf method), 56

set_chunk_size() (eskapade.analysis.links.ReadToDf method), 68

set_reader() (in module eskapade.analysis.links.read_to_df), 57

simulate() (eskapade.analysis.histogram.Histogram method), 81

skey (eskapade.analysis.histogram.ValueCounts attribute), 83

sum_counts (eskapade.analysis.histogram.ValueCounts attribute), 83

sum_data_length() (eskapade.analysis.links.read_to_df.ReadToDf method), 56

sum_data_length() (eskapade.analysis.links.ReadToDf method), 68

sum_nonone_counts (eskapade.analysis.histogram.ValueCounts attribute), 83

surface() (eskapade.analysis.histogram.Histogram method), 81

T

TimePeriod (class in eskapade.analysis.datetime), 73

to_date_time() (in module eskapade.data_quality.dq_helper), 97

to_float() (in module eskapade.data_quality.dq_helper), 97

to_int() (in module eskapade.data_quality.dq_helper), 97

to_normalized() (eskapade.analysis.histogram.Histogram method), 81

to_ns() (in module eskapade.analysis.histogram_filling), 86

to_str() (in module eskapade.analysis.histogram_filling), 86

to_str() (in module eskapade.data_quality.dq_helper), 98

truncated_bin_edges() (eskapade.analysis.histogram.BinningUtil method), 76

U

UniformTsTimePeriod (class in eskapade.analysis.datetime), 75

V

value_to_bin_center() (in module eskapade.analysis.histogram_filling), 86

value_to_bin_index() (in module eskapade.analysis.histogram_filling), 86

value_to_bin_label() (eskapade.analysis.histogram.BinningUtil method), 76

ValueCounter (class in eskapade.analysis.links), 70

ValueCounter (class in eskapade.analysis.links.value_counter), 60

ValueCounts (class in eskapade.analysis.histogram), 81

variable (eskapade.analysis.histogram.Histogram attribute), 81

W

weighted_quantile() (in module eskapade.analysis.statistics), 89

WriteFromDf (class in eskapade.analysis.links), 71

WriteFromDf (class in eskapade.analysis.links.write_from_df), 61