
Eskapade Documentation

Release

**KPMG Advanced Analytics
Big Data team**

Feb 28, 2018

Contents

1	Eskapade	3
2	Release notes	5
2.1	Version 0.7	5
2.2	Version 0.6	6
2.3	Version 0.5	6
2.4	Version 0.4	6
3	Contents	7
3.1	Introduction	7
3.2	Installation	8
3.3	Tutorials	15
3.4	Command Line Arguments	35
3.5	Package structure	40
3.6	Developing and Contributing	43
3.7	API	44
3.8	Appendices	137
3.9	Indices and tables	144
	Python Module Index	145

Version: 0.7.0

Date: Feb 28, 2018

Web page: <http://eskapade.kave.io>

Repository: <http://github.com/kaveio/eskapade>

Code reference: [API Documentation](#)

Issues & Ideas: <https://github.com/kaveio/eskapade/issues>

Q&A Support: contact us at: kave [at] kpmg [dot] com

CHAPTER 1

Eskapade

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

In particular, Eskapade can be 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.

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.

2.1 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 git@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.

2.2 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.

2.3 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!

2.4 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.

3.1 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.

3.1.1 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.

3.1.2 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](#).

3.2 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.

3.2.1 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 git@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.

3.2.2 Eskapade with Docker

Consistent environments for Eskapade development and use can be created with docker containers. Containers are created and managed by [Docker](#). An Eskapade container contains a pre-installed Eskapade setup and runs out-of-the-box. It is possible to mount your customised Eskapade code for development inside the container.

The instructions below show how one can use a locally checked-out version of Eskapade in combination with this Docker image. This combination is a very convenient way of actively working with and/or developing code in Eskapade.

By default, Eskapade code is executed as `root` user in the container. It is possible, however, to run with reduced user privileges inside containers through user mapping with the Docker host, as described in the last section.

Required software

Docker installation instructions can be found here: <https://docs.docker.com/install/>.

Eskapade source code is optional (for development purposes), you can check it out locally with the command:

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

Getting Eskapade docker images

From DockerHub

The official Eskapade docker image is provided on [DockerHub](#).

```
$ docker pull kave/eskapade-env:0.7
```

This will download the `kave/eskapade-env:0.7` image locally. Downloading this docker image can take a minute or two.

Building from scratch

To build the docker image from scratch using the Eskapade source code, do:

```
$ cd eskapade/docker/eskapade-env && sh create_docker.sh
```

This will produce the `kave/eskapade-env:0.7` image locally.

Spinning up docker containers

Out-of-the-box

From this image, containers with the Eskapade environment set up, can be run out-of-the-box:

```
$ docker run -p 8888:8888 -it kave/eskapade-env:0.7
```

Where port 8888 is forwarded to the docker host to make Jupyter notebook available (below).

E.g. one can now do:

```
cd /opt/eskapade
eskapade_run --help
```

and run any Eskapade code. See the Tutorial section for examples.

Exit the (docker) bash shell with:

```
exit
```

See section *After you exit Docker* (below) for cleaning up obsolete docker processes.

Mounting source code

```
$ docker run -v <ESKAPADE>:/opt/eskapade -p 8888:8888 -it kave/eskapade-env:0.7
```

Where `<ESKAPADE>` specifies the path of the Eskapade source code on the docker host, and where `/opt/eskapade` is the location of the Eskapade source code inside the container.

NOTE: in case you mount a clean installation of the Eskapade source code, you have to (re-)build the libraries by executing:

```
$ pip install -e /opt/eskapade
```

Running as non-root user

For increased security in a production environment, it is recommended to run Eskapade code inside the container as non-root user. The `Dockerfile` in the `eskapade-user` directory provides an additional user-mapping layer to the `eskapade-env` image: it creates a `esdev` user that has its own virtual Python environment with Eskapade installed. The mapping of user id's between Docker host and container ensure that proper permissions are propagated when writing/reading to the mounted volume with Eskapade code.

To obtain a centrally produced Eskapade image, use:

```
$ docker pull kave/eskapade-usr:0.7
```

Or build the Eskapade docker image with `esdev` user installation, from scratch:

```
$ cd docker/eskapade-usr && docker build -t kave/eskapade-usr:0.7 .
```

This will produce the `kave/eskapade-usr:0.7` image.

From this image, containers with the Eskapade environment set up, can be run out-of-the-box:

```
$ docker run -e HOST_USER_ID=$(id -u) -e HOST_USER_GID=$(id -g) -p 8888:8888 -it kave/  
↳eskapade-usr:0.7
```

The first time you run this command it will likely take some time. The `HOST_USER_ID` and `HOST_USER_GID` environment variables are used to dynamically map user- and group id's between the host and Docker container, ensuring proper read/write permissions.

Remapping the user id

To prevent the remapping of user and group id from happening the next time you boot up the image, open another shell:

```
$ docker ps
```

Copy the top `CONTAINER-ID` string, matching the running instance of the `kave/eskapade-usr:0.7` image, and then paste it:

```
$ docker commit CONTAINER-ID kave/eskapade-usr:0.7
```

Next time when you run:

```
$ docker run -e HOST_USER_ID=$(id -u) -e HOST_USER_GID=$(id -g) -p 8888:8888 -it kave/  
↳eskapade-usr:0.7
```

the remapping of user and group id should no longer happen.

Mounting source code

Containers with the user-specific Eskapade environment setup can be run out-of-the-box, and with your own mounted (customised) source code, using:

```
$ docker run -e HOST_USER_ID=$(id -u) -e HOST_USER_GID=$(id -g) -v <ESKAPADE>:/home/
↳esdev/eskapade -p 8888:8888 -it kave/eskapade-usr:0.7
```

Where <ESKAPADE> specifies the path of the Eskapade source code.

NOTE: in case you mount a clean installation of the Eskapade source code, you have to (re-)build the libraries by executing:

```
$ pip install -e /home/esdev/eskapade
```

This combination is a great way of using and developing Eskapade code.

Consider adding a permanent alias to your local `~/.bashrc` or `~/.bash_profile` file:

```
alias eskapade_docker='docker run -e HOST_USER_ID=$(id -u) -e HOST_USER_GID=$(id -g) -
↳v <ESKAPADE>:/home/esdev/eskapade -p 8888:8888 -it kave/eskapade-usr:0.7'
```

So the next time, in a fresh shell, you can simply run the command `eskapade_docker`.

Starting Jupyter notebook

To run the Jupyter notebook on port 8888 from the docker environment:

```
cd /opt/eskapade
jupyter &
```

And press enter twice to return to the shell prompt.

The command `jupyter &` starts up Jupyter notebook in the background on port 8888 and pipes the output to the log file `nohup.out`.

In your local browser then go to address:

```
localhost:8888/
```

And you will see the familiar Jupyter environment. In case you get asked for a password, take a look at the tail end of the file `nohup.out`, where you will see the exact url address that you need to go to.

E.g. you can now do `import eskapade` (shift-enter) to get access to the Eskapade library.

Be sure to run `jupyter &` from a directory that is mounted in the docker container, such as `/opt/eskapade`. In this way any notebook(s) you create are kept after you exit the docker run.

After you exit Docker

Every time you want to have a clean Docker environment, run the following commands:

```
# --- 1. remove all exited docker processes
docker ps -a | grep Exited | awk '{print "docker stop "$1 "; docker rm "$1}' | sh

# --- 2. remove all failed docker image builds
docker images | grep "<none>" | awk '{print "docker rmi "$3}' | sh
```

```
# --- 3. remove dangling volume mounts
docker volume ls -qf dangling=true | awk '{print "docker volume rm "$1}' | sh
```

To automate this, we advise you put these commands in an executable `docker_cleanup.sh` script.

3.2.3 Eskapade on a virtual machine

Consistent environments for Eskapade development and use are created with virtual machines. Machines are created and managed by `Vagrant`, with `VirtualBox` as a provider.

Getting started

Required software

To build and run Eskapade boxes, `Vagrant` and `VirtualBox` are required. `Vagrant` can be downloaded at <https://www.vagrantup.com/downloads.html>. For Ubuntu (and other `Debian`-based systems), download the `Vagrant Debian` package, move to where you want to install this and run this command to install:

```
dpkg -i vagrant_<VERSION>_<ARCHITECTURE>.deb
```

where the version and architecture are in the file name.

`VirtualBox` can be downloaded and installed by following the instructions on <https://www.virtualbox.org/wiki/Downloads>. Make sure you install the latest version from the `VirtualBox` repository instead of an older version from the repositories of your system.

Repository

The repository, which contains the code to build your virtual environment, is hosted on `github`. Clone it to your machine with:

```
$ git clone git@github.com:KaveIO/Eskapade.git
```

All code below is executed in the root of this repository, unless otherwise stated.

Virtual environment

The environment box contains `Ubuntu` and the `KAVE Toolbox`, including `Anaconda`, `Spark`, and `ROOT`.

You can either build this box yourself, or download the image. In order to build it yourself, run the following commands. It might take a while to build (up to two hours!)

```
cd vagrant/dev
vagrant up
```

Alternatively, in order to download and import its image, do:

```
cd vagrant/dev_image
vagrant box add eskapade-dev.json
vagrant up
```


Eskapade users and developers log in as the user `esdev` on `localhost`, by default on port 2222:

```
ssh -p 2222 esdev@localhost
```

This user can log in with the password `esdev`.

By default the local directory `/path/to/your/local/escapade` is mounted (containing your local Eskapade repository) under `/opt/escapade` in your vagrant machine. You can now edit the files in this directory, either locally or in the (vagrant) bash shell, and any updates to these files will be kept after exiting the virtual machine.

You are now ready to use Eskapade!

The next time...

Simply do:

```
cd vagrant/dev
vagrant up
```

or:

```
cd vagrant/dev_image
vagrant up
```

depending on whether you built vagrant yourself or downloaded the image.

Then you can access it via ssh (password `esdev`):

```
ssh -p 2222 esdev@localhost
```

Easy log-in

To make logging in easier, the key pair `vagrant/dev/ssh/esdev_id_rsa.pub`, `vagrant/dev/ssh/esdev_id_rsa` can be used, and an example SSH configuration is provided in `vagrant/dev/ssh/config`. Put these files in your `~/.ssh/`:

```
cp vagrant/dev/ssh/* ~/.ssh/
```

You can then log in using the command:

```
ssh esdevbox
```

Vagrant boxes

Boxes are built and started with the command `vagrant up` in the directory of the `Vagrantfile` describing the box. A box can be restarted by executing `vagrant reload`. The virtual machines are administered by the `vagrant` user, which logs in by running `vagrant ssh` in the directory of the `Vagrantfile`. The `vagrant` user has root access to the system by password-less `sudo`.

Starting Jupyter notebook

To run the Jupyter notebook on port 8888 from the vagrant machine:

```
cd /opt/eskapade
jupy &
```

And press enter twice to return to the shell prompt.

The command `jupy &` starts up Jupyter notebook in the background on port 8888 and pipes the output to the log file `nohup.out`.

In your local browser then go to address:

```
localhost:8888/
```

And you will see the familiar Jupyter environment.

E.g. you can now do `import eskapade` (shift-enter) to get access to the Eskapade library.

Be sure to run `jupy &` from a directory that is mounted in the vagrant machine, such as `/opt/eskapade`. In this way any notebook(s) you create are kept after you exit the docker run.

Requirements

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

There are two optional subpackages which require external products: `root_analysis` and `spark_analysis` subpackages.

To be able to run `root_analysis`, [ROOT CERN's data analysis package](#) need to be compiled with the following flags:

```
$ -Dfftw3=ON -Dmathmore=ON -Dminuit2=ON -Droofit=ON -Dtmva=ON -Dsoversion=ON -
↳ Dthread=ON -Dpython3=ON \
$ -DPYTHON_EXECUTABLE=path_to_python_exe -DPYTHON_INCLUDE_DIR=path_to_python_include -
↳ DPYTHON_LIBRARY=path_to_python_lib
```

`spark_analysis` requires [Apache Spark](#) version 2.1.1 or higher.

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

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

Alternatively, consider installing [KaveToolbox](#) version 3.6 or higher. To install the released version:

```
$ yum -y install wget curl tar zip unzip gzip python
$ wget http://repos:kaverepos@repos.kave.io/noarch/KaveToolbox/3.6-Beta/kavetoolbox-
↳ installer-3.6-Beta.sh
$ sudo bash kavetoolbox-installer-3.6-Beta.sh [--quiet]
```

(`--quiet` is for a quieter install, remove the brackets!)

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](#)

3.2.4 Installing Eskapade on macOS

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

3.3 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 are run from the root of the repository if not otherwise stated.

3.3.1 Running your first macro

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

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 python/eskapade/tutorials/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.

3.3.2 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 python/eskapade/tutorials/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 `python/eskapade/tutorials/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 python/eskapade/tutorials/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.

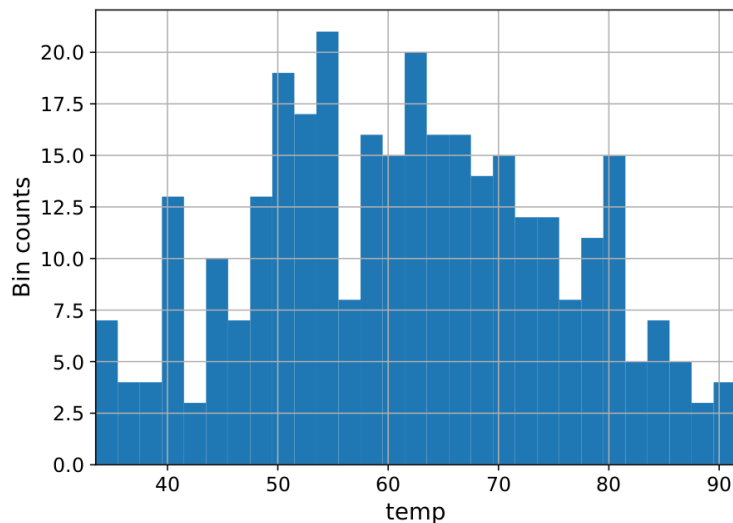
Note: If you don't have `pdflatex` installed on your machine you can install it by executing the following command: ..
code-block:: bash

```
$ 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!

3.3.3 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.

3.3.4 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>"
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.

3.3.5 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.

3.3.6 Tutorial5: using RooFit

This section provides a tutorial on how to use RooFit in Eskapade. RooFit is an advanced fitting library in ROOT, which is great for modelling all sorts of data sets. See [this tutorial](#) for a 20 min introduction into RooFit. ROOT (and RooFit) works 'out of the box' in the Eskapade docker/vagrant image.

In this tutorial we will illustrate how to define a new probability density function (pdf) in RooFit, how to compile it, and how to use it in Eskapade to simulate a dataset, fit it, and plot the results.

Note: There are many good RooFit tutorials. See the macros in the directory `$ROOTSYS/tutorials/roofit/` of your local ROOT installation. This tutorial is partially based on the RooFit tutorial `$ROOTSYS/tutorials/roofit/rf104_classfactory.C`.

Building a new probability density function

Before using a new model in Eskapade, we need to create, compile and load a probability density function model in RooFit.

Move to the directory:

```
$ cd cxx/esroofit/src/
```

Start an interactive python session and type:

```
import ROOT
ROOT.RooClassFactory.makePdf("MyPdfV2", "x,A,B", "", "A*fabs(x)+pow(x-B,2) ")
```

This command creates a RooFit skeleton probability density function class named `MyPdfV2`, with the variable `x`, `a`, `b` and the given formula expression.

Also type:

```
ROOT.RooClassFactory.makePdf("MyPdfV3", "x,A,B", "", "A*fabs(x)+pow(x-B,2) ", True, False,
↪ "x: (A/2) * (pow(x.max(rangeName), 2) + pow(x.min(rangeName), 2)) + (1./3) * (pow(x.
↪ max(rangeName) - B, 3) - pow(x.min(rangeName) - B, 3)) ")
```

This creates the RooFit p.d.f. class `MyPdfV3`, with the variable `x`, `a`, `b` and the given formula expression, and the given expression for analytical integral over `x`.

Exit python (Ctrl-D) and type:

```
$ ls -l MyPdf*
```

You will see two `cxx` files and two header files. Open the file `MyPdfV2.cxx`. You should see an `evaluate()` method in terms of `x`, `a` and `b` with the formula expression we provided.

Now open the file `MyPdfV3.cxx`. This also contains the method `analyticalIntegral()` with the expression for the analytical integral over `x` that we provided.

If no analytical integral has been provided, as in `MyPdfV2`, RooFit will try to try to compute the integral itself. (Of course this is a costly operation.) If you wish, since we know the analytical integral for `MyPdfV2`, go ahead and edit `MyPdfV2.cxx` to add the expression of the analytical integral to the class.

As another example of a simple pdf class, take a look at the expressions in the file: `cxx/esroofit/src/RooWeibull.cxx`.

Now move the header files to their correct location:

```
$ mv MyPdfV*.h ../include/
```

To make sure that these classes get picked up in Eskapade roofit library, open the file `cxx/esroofit/dict/esroofit/LinkDef.h` and add the lines:

```
#pragma link C++ class MyPdfV2+;
#pragma link C++ class MyPdfV3+;
```

Finally, let's compile the C++ code of these classes:

```
$ mkdir $ESKAPADE/build
$ cd $ESKAPADE/build
$ cmake ../cxx/esroofit
$ cmake --build .
```

You should see the compiler churning away, processing several existing classes but also MyPdfV2 and MyPdfV3.

We are now able to open the Eskapade roofit library, so we can use these classes in python:

```
from eskapade.root_analysis import roofit_utils
roofit_utils.load_libesroofit()
```

In fact, this last snippet of code is used in the tutorial macro right below.

Running the tutorial macro

Let's take a look at the steps in tutorial macro `python/eskapade/tutorials/tutorial_5.py`. The macro illustrates how do basic statistical data analysis with roofit, by making use of the `RoWorkspace` functionality. A `RoWorkspace` is a persistable container for RooFit projects. A workspace can contain and own variables, p.d.f.s, functions and datasets. The example shows how to define a pdf, simulate data, fit this data, and then plot the fit result. There are 5 sections; they are detailed in the sections below.

The next step is to run the tutorial macro.

```
$ eskapade_run python/eskapade/tutorials/tutorial_5.py
```

Let's discuss what we are seeing on the screen.

Loading the Eskapade ROOT library

The macro first checks the existence of the class MyPdfV3 that we just created in the previous section.

```
# --- 0. make sure Eskapade RooFit library is loaded

# --- load and compile the Eskapade roofit library
from eskapade.root_analysis import roofit_utils
roofit_utils.load_libesroofit()

# --- check existence of class MyPdfV3 in ROOT
pdf_name = 'MyPdfV3'
logger.info('Now checking existence of ROOT class {name}', name=pdf_name)
cl = ROOT.TClass.GetClass(pdf_name)
if not cl:
    logger.fatal('Could not find ROOT class {name}. Did you build and compile it_
↳correctly?', name=pdf_name)
    sys.exit(1)
else:
    logger.info('Successfully found ROOT class {name}', name=pdf_name)
```

In the output on the screen, look for Now checking existence of ROOT class MyPdfV3. If this was successful, it should then say Successfully found class MyPdfV3.

Instantiating a pdf

The link `WsUtils`, which stands for `RoWorkspace` utils, allows us to instantiate a pdf. Technically, one defines a model by passing strings to the `rooworkspace` factory. For examples on using the `rooworkspace` factory see [basic](#), [operations](#) and [tools](#) for more details. The entire `rooworkspace` factory syntax can be found at [commands](#).

```
ch = Chain('WsOps')

# --- instantiate a pdf
wsu = root_analysis.WsUtils(name = 'modeller')
wsu.factory = ["MyPdfV3::testpdf(y[-10,10],A[10,0,100],B[2,-10,10])"]
ch.add(wsu)
```

Here we use the pdf class we just created (`MyPdfV3`) to create a pdf called `testpdf`, with observable `y` and parameter `A` and `B`, having ranges $(-10, 10)$, $(0, 100)$ and $(-10, 10)$ respectively, and with initial values for `A` and `B` of 10 and 2 respectively.

Simulating data

The link `WsUtils` is then used to simulate records according to the shape of `testpdf`.

```
wsu = root_analysis.WsUtils(name = 'simulator')
wsu.add_simulate(pdf='testpdf', obs='y', num=400, key='simdata')
ch.add(wsu)
```

Here we simulate 400 records of observable `y` with pdf `testpdf` (which is of type `MyPdfV3`). The simulated data is stored in the datastore under key `simdata`.

Fitting the data

Another version of the link `WsUtils` is then used to fit the simulated records with the pdf `testpdf`.

```
wsu = root_analysis.WsUtils(name = 'fitter')
wsu.pages_key='report_pages'
wsu.add_fit(pdf='testpdf', data='simdata', key='fit_result')
ch.add(wsu)
```

The link performs a fit of pdf `testpdf` to dataset `simdata`. We store the fit result object in the datastore under key `fit_result`. The fit knows from the input dataset that the observable is `y`, so that the fit parameters are `A` and `B`.

Plotting the fit result

Finally, the last version of the link `WsUtils` is used to plot the result of the fit on top of simulated data.

```
wsu = root_analysis.WsUtils(name = 'plotter')
wsu.pages_key='report_pages'
wsu.add_plot(obs='y', data='simdata', pdf='testpdf', pdf_kwargs={'VisualizeError':
↪ 'fit_result', 'MoveToBack': ()}, key='simdata_plot')
wsu.add_plot(obs='y', pdf='testpdf', file='fit_of_simdata.pdf', key='simdata_plot')
ch.add(wsu)
```

This link is configured to do two things. First it plots the observable `y` of the the dataset `simdata` and then plots the fitted uncertainty band of the pdf `testpdf` on top of this. The plot is stored in the datastore under the

key `simdata_plot`. Then it plots the fitted pdf `testpdf` without uncertainty band on top of the same frame `simdata_plot`. The resulting plot is stored in the file `fit_of_simdata.pdf`

Fit report

The link `WsUtils` produces a summary report of the fit it has just performed. The pages of this report are stored in the datastore under the key `report_pages`. At the end of the Eskapade session, the plots and latex files produced by this tutorial are written out to disk.

The fit report can be found at:

```
$ cd results/tutorial_5/data/v0/report/
$ pdflatex report.tex
```

Take a look at the resulting fit report: `report.pdf`. It contains pages summarizing: the status and quality of the fit (including the correlation matrix), summary tables of the floating and fixed parameters in the fit, as well as the plot we have produced.

Other ROOT Examples in Eskapade

Other example Eskapade macros using ROOT and RooFit can be found in the `python/eskapade/tutorials` directory, e.g. see `esk401_roothist_fill_plot_convert.py` and all other 400 numbered macros.

3.3.7 Tutorial 6: going Spark

This section provides a tutorial on how to use Apache Spark in Eskapade. Spark works ‘out of the box’ in the Eskapade `docker/vagrant` image. For details on how to setup a custom Spark setup, see the [Spark](#) section in the Appendix.

In this tutorial we will basically redo Tutorial 1 but use Spark instead of Pandas for data processing. The following paragraphs describe step-by-step how to run a Spark job, use existing links and write your own links for Spark queries.

Note: To get familiar with Spark in Eskapade you can follow the exercises in `python/eskapade/tutorials/tutorial_6.py`.

Running the tutorial macro

The very first step to run the tutorial Spark job is:

```
$ eskapade_run python/eskapade/tutorials/tutorial_6.py
```

Eskapade will start a Spark session, do nothing, and quit - there are no chains/links defined yet. The Spark session is created via the `SparkManager` which, like the `DataStore`, is a singleton that configures and controls Spark sessions centrally. It is activated through the magic line:

```
process_manager.service(SparkManager).create_session(include_eskapade_modules=True)
```

Note that when the Spark session is created, the following line appears in logs:

```
Adding Python modules to egg archive <PATH_TO_ESKAPADE>/lib/es_python_modules.egg
```

This is the `SparkManager` that ensures all Eskapade source code is uploaded and available to the Spark cluster when running in a distributed environment. To include the Eskapade code the argument `include_eskapade_modules` need to be set to `True` (by default it is `False`).

If there was an `ImportError: No module named pyspark` then, most likely, `SPARK_HOME` and `PYTHONPATH` are not set up correctly. For details, see the [Spark](#) section in the Appendix.

Reading data

Spark can read data from various sources, e.g. local disk, HDFS, HIVE tables. Eskapade provides the `SparkDfReader` link that uses the `pyspark.sql.DataFrameReader` to read flat CSV files into Spark DataFrames, RDD's, and Pandas DataFrames. To read in the Tutorial data, the following link should be added to the `Data` chain:

```
data = Chain('Data')
reader = spark_analysis.SparkDfReader(name='Read_LA_ozone', store_key='data', read_
↳methods=['csv'])
reader.read_meth_args['csv'] = (DATA_FILE_PATH,)
reader.read_meth_kwargs['csv'] = dict(sep=',', header=True, inferSchema=True)
data.add(reader)
```

The `DataStore` holds a pointer to the Spark dataframe in (distributed) memory. This is different from a Pandas dataframe, where the entire dataframe is stored in the `DataStore`, because a Spark dataframe residing on the cluster may not fit entirely in the memory of the machine running Eskapade. This means that Spark dataframes are never written to disk in `DataStore` pickles!

Using existing links

Spark has a large set of standard functions for Spark DataFrames and RDD's. Although the purpose of Eskapade is not to duplicate this functionality, there are some links created for generic functions to facilitate specifying Spark queries directly in the macro, instead of hard-coding them in links. This is handy for bookkeeping queries at a central place and reducing code duplication, especially for smaller analysis steps. For example, the `SparkExecuteQuery` link takes any string containing SQL statements to perform a custom query with Spark on a dataframe.

Column transformations

To add two columns to the Tutorial data using the conversion functions defined earlier in the macro, two `SparkWithColumn` links need to be added to the `Data` chain, one for each additional column:

```
from pyspark.sql.functions import udf
from pyspark.sql.types import TimestampType, FloatType
...

transform = spark_analysis.SparkWithColumn(name='Transform_doy', read_key=reader.
↳store_key,
                                          store_key='transformed_data', col_select=[
↳'doy'],
                                          func=udf(comp_date, TimestampType()), new_
↳column='date')
data.add(transform)
transform = spark_analysis.SparkWithColumn(name='Transform_vis', read_key=transform.
↳store_key,
                                          store_key='transformed_data', col_select=[
↳'vis'],
```



```

func=udf(mi_to_km, FloatType()), new_
↪column='vis_km')
data.add(transform)

```

Note that the functions defined in the macro are converted to user-defined functions with `pyspark.sql.functions.udf` and their output types are explicitly specified in terms of `pyspark.sql.types`. Omitting these type definitions can lead to obscure errors when executing the job.

Creating custom links

More complex queries deserve their own links since links provide full flexibility w.r.t. specifying custom data operation. For this Tutorial the ‘complex query’ is to just print 42 rows of the Spark dataframe. Of course, more advanced Spark functions can be applied in a similar fashion. A link is created just like was done before, e.g.:

```
$ eskapade_generate_link --dir python/eskapade/spark_analysis/links SparkDfPrinter
```

This creates the link `python/eskapade/spark_analysis/links/sparkdfprinter.py`. Do not forget to include the import statements in the `__init__.py` file as indicated by the `eskapade_generate_link` command.

The next step is to add the desired functionality to the link. In this case, the Spark dataframe needs to be retrieved from the `DataStore` and a `show()` method of that dataframe needs to be executed. The `execute()` method of the link is the right location for this:

```

def execute(self):
    """Execute the link.

    :returns: status code of execution
    :rtype: StatusCode
    """
    settings = process_manager.service(ConfigObject)
    ds = process_manager.service(DataStore)

    # --- your algorithm code goes here
    self.logger.debug('Now executing link: {link}.', link=self.name)
    df = ds[self.read_key]
    df.show(self.nrows)

    return StatusCode.Success

```

There is an additional attribute `self.nrows` which should be set in the link. By default, a generated link process only the `read_key` and `store_key` arguments and fails if there are any residual `kwargs`. To set the `nrows` attribute, add `nrows` to the key-value arguments in the `__init__()` method:

```

def __init__(self, **kwargs):
    ...

    self._process_kwargs(kwargs, read_key=None, store_key=None, nrows=1)

```

In order to configure Eskapade to run this link, the link needs to be added to a chain, e.g. `Summary`, in the `tutorial/tutorial_6.py` macro. This should look similar to:

```

printer = spark_analysis.SparkDfPrinter(name='Print_spark_df', read_key=transform.
↪store_key, nrows=42)
summary.add(printer)

```

The name of the dataframe is the output name of the `transform` link and the number of rows to print is specified by the `nrows` parameter.

Eskapade should now be ready to finally execute the macro and provide the desired output:

```
$ eskapade_run python/eskapade/tutorials/tutorial_6.py

* * * Welcome to Eskapade * * *
...

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|ozone|  vh|wind|humidity|temp|  ibh|dpg|ibt|vis|doy|          date|  vis_km|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|   3|5710|  4|      28|  40|2693|-25|  87|250|  3|1976-01-03 00:00:...|  402.335|
|   5|5700|  3|      37|  45| 590|-24|128|100|  4|1976-01-04 00:00:...|  160.934|
|   5|5760|  3|      51|  54|1450| 25|139| 60|  5|1976-01-05 00:00:...|  96.5604|
...

|   6|5700|  4|      86|  55|2398| 21|121|200| 44|1976-02-13 00:00:...|  321.868|
|   4|5650|  5|      61|  41|5000| 51| 24|100| 45|1976-02-14 00:00:...|  160.934|
|   3|5610|  5|      62|  41|4281| 42| 52|250| 46|1976-02-15 00:00:...|  402.335|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 42 rows
...

* * * Leaving Eskapade. Bye! * * *
```

That's it!

Spark Streaming

Eskapade supports the use of Spark Streaming as demonstrated in the word count example `tutorials/esk610_spark_streaming_wordcount.py`. The data is processed in (near) real-time as micro batches of RDD's, so-called discretized streaming, where the stream originates from either new incoming files or network connection. As with regular Spark queries, various transformations can be defined and applied in subsequent Eskapade links.

For details on Spark Streaming, see also <https://spark.apache.org/docs/2.1.1/streaming-programming-guide.html>.

File stream

The word count example using the file stream method can be run by executing in two different terminals:

```
terminal 1 $ eskapade_run -c stream_type='file' python/eskapade/tutorials/esk610_
↳spark_streaming_wordcount.py

terminal 2 $ mkdir /tmp/eskapade_stream_test/
terminal 2 $ for ((i=0; i<=100; i++)); do echo "Hello world" > /tmp/eskapade_stream_
↳test/dummy_$(printf %05d $i); sleep 0.2; done
```

Where bash `for`-loop will create a new file containing `Hello world` in the `/tmp/eskapade_stream_test` directory every 0.2 second. Spark Streaming will pick up and process these files and in terminal 1 a word count of the processed data will be displayed. Output is stored in `results/esk610_spark_streaming/data/v0/dstream/wordcount`. Only new files in `/tmp/eskapade_stream_test` are processed, do not forget to delete this directory.

TCP stream

The word count example using the TCP stream method can be run by executing in two different terminals:

```
terminal 1 $ eskapade_run -c stream_type='tcp' python/eskapade/tutorials/esk610_spark_
↳streaming_wordcount.py

terminal 2 $ nc -lk 9999
```

Where `nc` (netcat) will stream data to port 9999 and Spark Streaming will listen to this port and process incoming data. In `terminal 2` random words can be type (followed by enter) and in `terminal 1` a word count of the processed data will be displayed. Output is stored in `results/esk610_spark_streaming/data/v0/dstream/wordcount`.

All Spark examples

All example Eskapade macros using Spark can be found in the `python/eskapade/tutorials` directory, see `esk601_spark_configuration.py` and further.

3.3.8 All available examples

Every subpackage of Eskapade contains links in its `links/` subdirectory.

- `core_ops` contains links pertaining to the core functionality of Eskapade, where the `core` package is the core framework of Eskapade.
- `analysis` contains pandas links.
- `visualization` contains plotter links.
- `root_analysis` contains ROOT links for data generation, fitting, and plotting.
- `data_quality` contains links for fixing messy data.
- `spark_analysis` contains Spark related analysis links.

The name of every link indicates its basic function. If you want to know explicitly you can read the [API documentation](#). If that does not help, read and try to understand the example macros in `python/eskapade/tutorials/`, which show the basic usage of most Eskapade functionality. (See also the [Examples](#) section right below.) If still unclear, go into the link's code to find out how it exactly works.

Many Eskapade example macros can be found in the `python/eskapade/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.
- `esk400+`: macros for processing ROOT datasets and performing fits to data using RooFit.
- `esk500+`: macros for doing data quality assessment and cleaning.
- `esk600+`: macros describing links to do basic processing of data and RDDs with Spark.

The basic Eskapade macros (`esk100+`) 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.

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.

Example esk103: Print the DataStore

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

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.

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.

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.

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.

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.

Example esk109: Debugging tips

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

Example esk110: Code profiling

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

Example esk201: Read data

Macro 201 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.)

Example esk202: Write data

Macro 202 reads a DataFrame into the data store and then writes the DataFrame to csv format on the disk.

3.3.9 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.

3.4 Command Line Arguments

3.4.1 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.

3.4.2 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 IPython 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

3.4.3 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 an IPython 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

Your 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
```

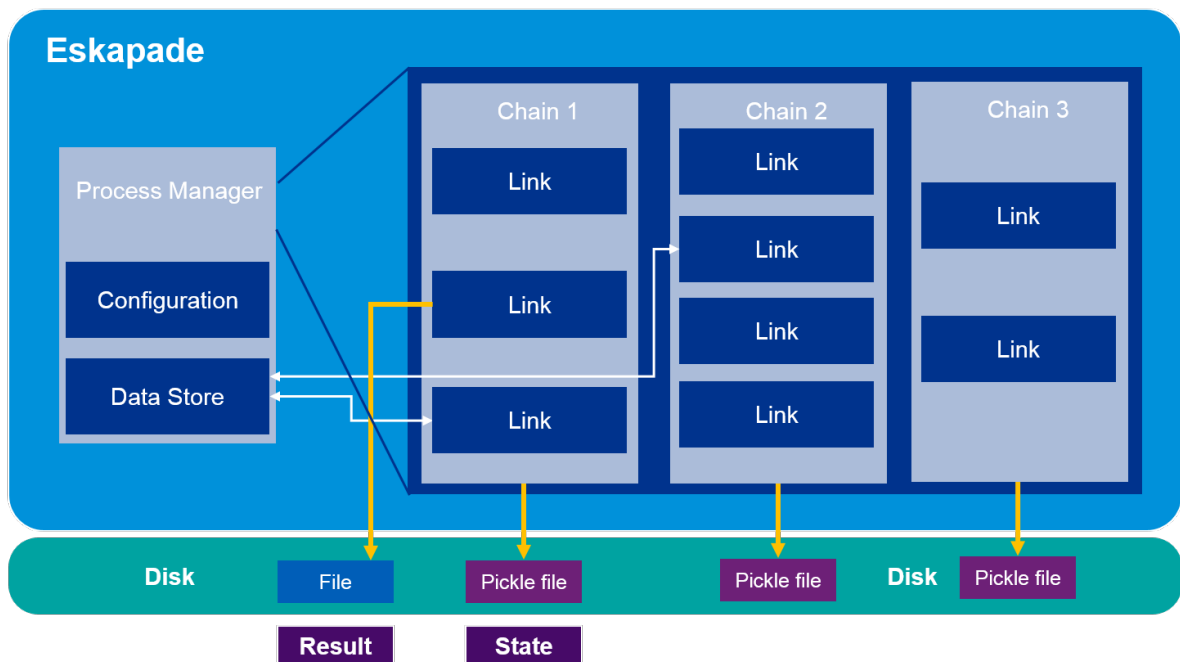
3.5 Package structure

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.

3.5.1 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` runs the macro file,
- 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 components of the architecture of Eskapade are explained in further detail in the [Tutorials section](#).

3.5.2 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
```

3.5.3 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.
- `root_analysis` contains ROOT links and code for data generation, fitting, and plotting.
- `data_quality` contains links and code for fixing messy data.
- `spark_analysis` contains Spark related analysis links and code.

3.5.4 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.

3.5.5 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`.

3.5.6 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.

3.5.7 Troubleshooting

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

3.6 Developing and Contributing

3.6.1 Preliminaries

3.6.2 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 git@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.

3.6.3 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.

3.6.4 Tips and Tricks

- Enable auto reload in ipython:

```
%load_ext autoreload
```

this will reload modules before executing any user code.

3.7 API

3.7.1 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: eskapade.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: `eskapade.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: `eskapade.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: `eskapade.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: `eskapade.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: `eskapade.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: `eskapade.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** – pandas reader is determined automatically. But can be set by hand, e.g. csv, xlsx.
- **itr_over_files** (*bool*) – Iterate over individual files, default is false. If false, are files are collected in one dataframe. NB chunksize takes priority!
- **chunksize** (*int*) – Default is none. If positive integer then will always iterate. chunk-size requires `pd.read_csv` or `pd.read_table`.
- **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.pandasReader` (*path*, *reader*, *args, **kwargs)

Pick the correct pandas 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: eskapade.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 orginal 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: `eskapade.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,
                                                             col-
                                                             umn_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 datastore 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 all 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: `eskapade.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`. For example: ‘csv’
- **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.
- **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.pandas_writer` (*path*, *writer*)

Pick the correct pandas writer.

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

Module contents

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

Bases: `eskapade.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: `eskapade.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: `eskapade.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: `eskapade.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: `eskapade.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*) – colums 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)

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

Bases: `eskapade.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<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: `eskapade.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** – pandas reader is determined automatically. But can be set by hand, e.g. csv, xlsx.
- **itr_over_files** (*bool*) – Iterate over individual files, default is false. If false, are files are collected in one dataframe. NB chunksize takes priority!
- **chunksize** (*int*) – Default is none. If positive integer then will always iterate. chunk-size requires `pd.read_csv` or `pd.read_table`.
- **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: `eskapade.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: `eskapade.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 datastore 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 all 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: `eskapade.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`. For example: ‘csv’
- **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.
- **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.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: `eskapade.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`, `eskapade.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_spec = { '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 if subkey is in different order than key. Default is true.

Returns successful projection or not

Return type bool

remove_keys_of_inconsistent_type (*prefered_key_type=None*)

Remove keys with inconsistent data type(s).

Parameters **prefered_key_type** (*tuple*) – the prefered 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: `eskapade.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
>>>
>>> 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

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

Created: 2017/02/27

Description: Definitions used in Eskapade runs: * logging levels * return-status codes * default configuration variables * user options

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.core.definitions.RandomSeeds (**kwargs)
```

Bases: object

Container for seeds of random generators.

Seeds are stored as key-value pairs and are accessed with `getitem` and `setitem` methods. A default seed can be accessed with the key “default”. The default seed is also returned if no seed is set for the specified key.

```
>>> import numpy as np
>>> seeds = RandomSeeds(default=999, foo=42, bar=13)
>>> seeds['NumPy'] = 100
>>> np.random.seed(seeds['NumPy'])
>>> print(seeds['nosuchseed'])
999
```

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

Initialize an instance.

Values of the specified keyword arguments must be integers, which are set as seed values for the corresponding key.

```
class eskapade.core.definitions.StatusCode
```

Bases: `enum.IntEnum`

Return status code enumeration class.

A `StatusCode` should be returned by the `initialize`, `execute`, and `finalize` methods of links, chains, and the process manager.

The enumerations are:

- Undefined (-1): Default status.
- Success (0 == `EX_OK` / `EXIT_SUCCESS`): All OK, i.e. there were no errors.
- RepeatChain (1): Repeat execution of this chain.
- SkipChain (2): Skip this chain: initialize, execute, and finalize.
- BreakChain (3): Skip the further execution of this this, but do perform finalize.
- Recoverable (4): Not OK, but can continue, i.e. there was an error, but the application can recover from it.
- Failure (5): An error occurred and the application cannot recover from it. In this case the application should just quit.

```
BreakChain = 3
```

```
Failure = 5
```

```
Recoverable = 4
```

```
RepeatChain = 1
```

```
SkipChain = 2
```

```
Success = 0
```

```
Undefined = -1
```

```
is_break_chain () → bool
```

Check if status is `StatusCode.BreakChain`.

Returns True when `StatusCode.BreakChain`, False otherwise.

Return type bool

is_failure () → bool

Check if status is *StatusCode.Failure*.

Returns True when *StatusCode.Failure*, False otherwise.

Return type bool

is_recoverable () → bool

Check if status is *StatusCode.Recoverable*.

Returns True when *StatusCode.Recoverable*, False otherwise.

Return type bool

is_repeat_chain () → bool

Check if status is *StatusCode.RepeatChain*.

Returns True when *StatusCode.RepeatChain*, False otherwise.

Return type bool

is_skip_chain () → bool

Check if status is *StatusCode.SkipChain*.

Returns True when *StatusCode.SkipChain*, False otherwise.

Return type bool

is_success () → bool

Check if status is *StatusCode.Success*.

Returns True when *StatusCode.Success*, False otherwise.

Return type bool

is_undefined () → bool

Check if status is *StatusCode.Undefined*.

Returns True when *StatusCode.Undefined*, False otherwise.

Return type bool

`eskapade.core.definitions.set_begin_end_chain_opt (opt_key, settings, args)`

Set begin/end-chain variable from user option.

`eskapade.core.definitions.set_custom_user_vars (opt_key, settings, args)`

Set custom user configuration variables.

`eskapade.core.definitions.set_log_level_opt (opt_key, settings, args)`

Set configuration log level from user option.

`eskapade.core.definitions.set_opt_var (opt_key, settings, args)`

Set configuration variable from user options.

`eskapade.core.definitions.set_seeds (opt_key, settings, args)`

Set random seeds.

`eskapade.core.definitions.set_single_chain_opt (opt_key, settings, args)`

Set single-chain variable from user option.

eskapade.core.element module

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

Created: 2017/02/27

Description: Base classes for the building blocks of an Eskapade analysis run:

- Link:
- Chain:

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.core.element.Chain` (*name*, *process_manager=None*)

Bases: `eskapade.core.meta.Processor`, `eskapade.core.meta.ProcessorSequence`, `eskapade.core.mixin.TimerMixin`

Execution Chain.

A Chain object contains a collection of links with analysis code. The links in a chain are executed in the order in which the links have been added to the chain. Typically a chain contains all links related to one topic, for example ‘validation of a model’, or ‘data preparation’, or ‘data quality checks’.

```
>>> from eskapade import process_manager
>>> from eskapade import Chain
>>> from eskapade import analysis
>>>
>>> # Create an IO chain. This is automatically registered with the process_
↳manager.
>>> io_chain = Chain('Overview')
```

And Links are added to a chain as follows:

```
>>> # add a link to the chain
>>> io_chain.add(analysis.ReadToDf(path='foo.csv', key='foo'))
>>>
>>> # Run everything.
>>> process_manager.run()
```

__init__ (*name*, *process_manager=None*)

Initialize chain.

add (*link: eskapade.core.element.Link*) → None

Add a link to the chain.

Parameters **link** (`Link`) – The link to add to the chain.

Raises

- **TypeError** – When the link is of an incompatible type.
- **KeyError** – When a Link of the same type and name already exists.

clear ()

Clear the chain.

discard (*link: eskapade.core.element.Link*) → None

Remove a link from the chain.

Parameters **link** (`Link`) –

Raises **KeyError** – When the processor does not exist.

execute () → `eskapade.core.definitions.StatusCode`

Execute links in chain.

Returns Execution status code.

Return type *StatusCode*

finalize () → `eskapade.core.definitions.StatusCode`
Finalize links and chain.

Returns Finalization status code.

Rtype *StatusCode*

get (*link_name: str*) → `eskapade.core.element.Link`
Find the link with the given name.

Parameters *link_name* (*str*) – Find a link with the given name.

Returns The chain.

Return type *Chain*

Raises **ValueError** – When the given chain name cannot be found.

initialize () → `eskapade.core.definitions.StatusCode`
Initialize chain and links.

Returns Initialization status code.

Return type *StatusCode*

n_links

Return the number of links in the chain.

Returns The number of links in the chain.

Return type `int`

class `eskapade.core.element.Link` (*name=None*)

Bases: `eskapade.core.meta.Processor`, `eskapade.core.mixin.ArgumentsMixin`,
`eskapade.core.mixin.TimerMixin`

Link base class.

A link defines the content of an algorithm. Any actual link is derived from this base class.

A link usually does three things: - takes data from the datastore - does something to it - writes data back

To take from the data store there is a simple function `load()` To write to the data store there is a simple function `store()`

Links are added to a chain as follows:

```
>>> from eskapade import process_manager
>>> from eskapade import analysis
>>>
>>> # Create a Chain instance. Note that the chain is automatically registered,
↳with process manager.
>>> io_chain = Chain('IO')
>>>
>>> # Add a link to the chain
>>> reader = analysis.ReadToDf(name='CsvReader', key='foo')
>>> reader.path = 'foo.csv'
>>> io_chain.add(reader)
>>>
>>> # Run everything.
>>> process_manager.run()
```

`__init__ (name=None)`

Initialize link.

`execute ()` → `eskapade.core.definitions.StatusCode`

Execute the Link.

This method may be overridden by the user.

Returns Status code.

Return type *StatusCode*

`finalize ()` → `eskapade.core.definitions.StatusCode`

Finalize the Link.

This method may be overridden by the user.

Returns Status code.

Return type *StatusCode*

`initialize ()` → `eskapade.core.definitions.StatusCode`

Initialize the Link.

This method may be overridden by the user.

Returns Status code.

Type `StatusCode`

`load (ds, read_key=None)`

Read all data from specified source.

Read_key can either be:

- one Link: return statuscode, [data_from_link,...]
- A list of locations: return statuscode, [data,...]
- A list of links with only one output location: return statuscode, [data,...]
- A list of links with multiple output locations: return statuscode, [data,[moredata]...]
- Any mixture of the above

Do something logical with a statuscode if this data does not exist `link.if_input_missing = statuscode`

Returns a tuple statuscode, [data in same order as read_key]

Return type (*StatusCode*,list)

`store (ds, data, store_key=None, force=False)`

Store data back to datastore.

Do something logical with a statuscode if this data already exists `link.if_output_exists = statuscode` uses `self.store_key`. If `self.store_key` is a list of locations, I must sent a list of the same length here

`summary ()`

Print a summary of the main settings of the link.

eskapade.core.exceptions module

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

Created: 2016/11/08

Description: Eskapade exceptions.

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.

exception `eskapade.core.exceptions.Error`

Bases: `Exception`

Base class for all Eskapade core exceptions.

exception `eskapade.core.exceptions.UnknownSetting`

Bases: `eskapade.core.exceptions.Error`

The user requested an unknown setting.

eskapade.core.execution module

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

Created: 2016/11/08

Description: Functions for running and resetting Eskapade machinery

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.core.execution.eskapade_run(settings=None)`

Run Eskapade.

This function is called in the script `eskapade_run` when run from the cmd line. The working principle of Eskapade is to run chains of custom code chunks (so-called 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 principle, links can be easily reused in future projects.

Parameters `settings` (`ConfigObject`) – analysis settings

Returns status of the execution

Return type `StatusCode`

`eskapade.core.execution.reset_eskapade(skip_config=False)`

Reset Eskapade objects.

Parameters `skip_config` (`bool`) – skip reset of configuration object

eskapade.core.meta module

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

Created: 2017/09/14

Description:

A collection of (generic) meta classes for some (design) patterns:

- Singleton: Meta class for the Singleton pattern.

- Processor: Meta class with abstract methods initialize, execute, and finalize.
- ProcessorSequence: A simple (processor) sequence container.

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.core.meta.Processor` (*name: str*)

Bases: `object`

Processor metaclass.

__init__ (*name: str*)

Initialize the Processor object.

execute ()

Execution logic for processor.

finalize ()

Finalization logic for processor.

initialize ()

Initialization logic for processor.

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.

name

Get the name of processor.

Returns The name of the processor.

Return type str

parent

Get the group parent.

Returns The parent/group processor sequence.

class `eskapade.core.meta.ProcessorSequence`

Bases: object

A doubly linked processor sequence.

It remembers the order in which processors are added to the sequence. It also checks if a processor already has been added to the sequence.

__init__ ()

Initialize the ProcessorSequence object.

add (*processor*: `eskapade.core.meta.Processor`) → None

Add a processor to the sequence.

Parameters **processor** (`Processor`) – The processor to add.

Raises **KeyError** – When a processor of the same type and name already exists.

clear () → None

Clear the sequence.

discard (*processor*: `eskapade.core.meta.Processor`) → None

Remove a processor from the sequence.

Parameters **processor** (`Processor`) – The processor to remove.

Raises **KeyError** – When the processor does not exist.

pop (*last*: `bool = True`) → `eskapade.core.meta.Processor`

Return the popped processor. Raise `KeyError` if empty.

By default a processor is popped from the end of the sequence.

Parameters **last** (`bool`) – Pop processor from the end of the sequence. Default is `True`.

Returns The pop processor.

Raises **KeyError** – When trying to pop from an empty list.

```
class eskapade.core.meta.Singleton
```

Bases: type

Metaclass for singletons.

Any instantiation of a Singleton class yields the exact same object, e.g.:

```
>>> class Klass(metaclass=Singleton):
>>>     pass
>>>
>>> a = Klass()
>>> b = Klass()
>>> a is b
True
```

See <https://michaelgoerz.net/notes/singleton-objects-in-python.html>.

eskapade.core.mixin module

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

Created: 2016/11/08

Classes: ArgumentsMixin, TimerMixin

Description:

Mixin classes:

- ArgumentsMixin: processes/checks arguments and sets them as attributes
- TimerMixin: keeps track of execution time
- ConfigMixin: reads and handles settings from configuration files

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.core.mixin.ArgumentsMixin
```

Bases: object

Mixin base class for argument parsing.

Class allows attributes to be accessed as dict items. Plus several argument processing helper functions.

```
check_arg_callable (*arg_names, allow_none=False)
```

Check if set of arguments has iterators.

```
check_arg_iters (*arg_names, allow_none=False)
```

Check if set of arguments has iterators.

```
check_arg_opts (allow_none=False, **name_vals)
```

Check if argument values are in set of options.

```
check_arg_types (recurse=False, allow_none=False, **name_type)
```

Check if set of arguments has correct types.

```
check_arg_vals (*arg_names, allow_none=False)
```

Check if set of arguments exists as attributes and values.

```
check_extra_kwargs (kwargs)
```

Check for residual kwargs.

check_required_args (**arg_names*)
Check if set of arguments exists as attributes.

class `eskapade.core.mixin.ConfigMixin` (*config_path=None*)

Bases: `object`

Mixin base class for configuration settings.

__init__ (*config_path=None*)
Initialize config settings.

Parameters **config_path** (*str*) – path of configuration file

config_path
Path of configuration file.

get_config (*config_path=None*)
Get settings from configuration file.

Read and return the configuration settings from a configuration file. If the path of this file is not specified as an argument, the value of the “`config_path`” property is used. If the file has already been read, return previous settings.

Parameters **config_path** (*str*) – path of configuration file

Returns configuration settings read from file

Return type `configparser.ConfigParser`

Raises **RuntimeError** – if `config_path` is not set

reset_config ()
Remove previously read settings.

class `eskapade.core.mixin.TimerMixin`

Bases: `object`

Mixin base class for timing.

__init__ ()
Initialize timer.

start_timer ()
Start run timer.

Start the timer. The timer is used to compute the run time. The returned timer start value has an undefined reference and should, therefore, only be compared to other timer values.

Returns start time in seconds

Return type `float`

stop_timer (*start_time=None*)
Stop the run timer.

Stop the timer. The timer is used to compute the run time. The elapsed time since the timer start is returned.

Parameters **start_time** (*float*) – function `start_time` input

Returns time difference with start in seconds

Return type `float`

total_time ()
Return the total run time.

Returns total time in seconds

Return type float

eskapade.core.persistence module

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

Created: 2016/11/08

Description: Utility class and functions to get correct io path, used for persistence of results

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.core.persistence.create_dir(dir_path)`
Create directory.

Parameters `dir_path` (*str*) – directory path

`eskapade.core.persistence.io_dir(io_type, io_conf=None)`
Construct directory path.

Parameters

- **io_type** (*str*) – type of result to store, e.g. data, macro, results.
- **io_conf** – IO configuration object

Returns directory path

Return type str

`eskapade.core.persistence.io_path(io_type, sub_path, io_conf=None)`
Construct directory path with sub path.

Parameters

- **io_type** (*str*) – type of result to store, e.g. data, macro, results.
- **sub_path** (*str*) – sub path to be included in io path
- **io_conf** – IO configuration object

Returns full path to directory

Return type str

`eskapade.core.persistence.record_file_number(file_name_base, file_name_ext, io_conf=None)`

Get next prediction-record file number.

Parameters

- **file_name_base** (*str*) – base file name
- **file_name_ext** (*str*) – file name extension
- **io_conf** – I/O configuration object

Returns next prediction-record file number

Return type int

`eskapade.core.persistence.repl_whites(name)`
Replace whitespace in names.

eskapade.core.process_manager module

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

Class: ProcessManager

Created: 2016/11/08

Description: The ProcessManager class is the heart of Eskapade. It performs initialization, execution, and finalization of analysis chains.

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.core.process_manager.ProcessManager`

Bases: `eskapade.core.meta.Processor`, `eskapade.core.meta.ProcessorSequence`, `eskapade.core.mixin.TimerMixin`

Eskapade run process manager.

ProcessManager is the event processing loop of Eskapade. It initializes, executes, and finalizes the analysis chains. There is, under normal circumstances, only one ProcessManager instance.

Here's an simple but illustrative analysis example:

```
>>> from eskapade import process_manager, Chain, Link, StatusCode
>>>
>>> # A chain automatically registers itself with process_manager.
>>> one_plus_one_chain = Chain('one_plus_one')
>>>
>>> class OnePlusOne(Link):
>>>     def execute(self):
>>>         self.logger.info('one plus one = {result}', result=(1+1))
>>>         return StatusCode.Success
>>>
>>> one_plus_one_chain.add(link=OnePlusOne())
>>>
>>> two_plus_two_chain = Chain('two_plus_two')
>>>
>>> class TwoPlusTwo(Link):
>>>     def execute(self):
>>>         self.logger.info('two plus two = {result}', result=(2+2))
>>>         return StatusCode.Success
>>>
>>> two_plus_two_chain.add(TwoPlusTwo())
>>>
>>> process_manager.run()
```

Ideally the user will not need to interact directly with the process manager. The magic is taken care of by the `eskapade_run` entry point.

__init__ ()

Initialize ProcessManager instance.

add (*chain*: `eskapade.core.element.Chain`) → None

Add a chain to the process manager.

Parameters **chain** (`Chain`) – The chain to add to the process manager.

Raises

- **TypeError** – When the chain is of an incompatible type.
- **KeyError** – When a chain of the same type and name already exists.

clear()

“Clear/remove all chains.

execute()

Execute all chains in order.

Returns status code of execution attempt

Return type *StatusCode*

execute_macro(*filename*, *copyfile=True*)

Execute an input python configuration file.

A copy of the configuration file is stored for bookkeeping purposes.

Parameters

- **filename** (*str*) – the path of the python configuration file
- **copyfile** (*bool*) – back up the macro for bookkeeping purposes

Raises Exception – if input configuration file cannot be found

finalize()

Finalize the process manager manager.

Returns status code of finalize attempt

Return type *StatusCode*

get(*chain_name: str*) → *eskapade.core.element.Chain*

Find the chain with the given name.

Parameters **chain_name** (*str*) – Find a chain with the given name.

Returns The chain.

Return type *Chain*

Raises ValueError – When the given chain name cannot be found.

get_service_tree()

Create tree of registered process-service classes.

Returns service tree

Return type dict

get_services()

Get set of registered process-service classes.

Returns service set

Return type set

import_services(*io_conf*, *chain=None*, *force=None*, *no_force=None*)

Import process services from files.

Parameters

- **io_conf** (*dict*) – I/O config as returned by *ConfigObject.io_conf*
- **chain** (*str*) – name of chain for which data was persisted
- **force** (*bool or list*) – force import if service already registered

- **no_force** (*list*) – do not force import of services in this list

initialize ()

Initialize the process manager.

Initializes the process manager by configuring its chains. After initialization the configuration is printed.

Returns status code of initialize attempt

Return type *StatusCode*

n_chains

Return the number of chains in the process manager.

Returns The number of links in the chain.

Return type int

persist_services (*io_conf*, *chain=None*)

Persist process services in files.

Parameters

- **io_conf** (*dict*) – I/O config as returned by ConfigObject.io_conf
- **chain** (*str*) – name of chain for which data is persisted

print_chains ()

Print all chains defined in the manager.

print_services ()

Print registered process services.

remove_all_services ()

Remove all registered process services.

remove_service (*service_cls*, *silent=False*)

Remove specified process service.

Parameters

- **service_cls** (*ProcessServiceMeta*) – service to remove
- **silent** (*bool*) – don't complain if service is not registered

reset ()

Reset the process manager.

Resetting comprises removing the chains and closing any open connections/sessions.

run () → *eskapade.core.definitions.StatusCode*

Run process manager.

Returns Status code of run execution.

Return type *StatusCode*

service (*service_spec*)

Get or register process service.

Parameters **service_spec** (*ProcessServiceMeta* or *ProcessService*) – class (instance) to register

Returns registered instance

Return type *ProcessService*

summary ()

Print process-manager summary.

Print a summary of the chains, links, and some analysis settings defined in this configuration.

eskapade.core.process_services module

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

Created: 2017/02/27

Description: Base class and core implementations of run-process services

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.core.process_services.ConfigObject`

Bases: `eskapade.core.process_services.ProcessService`

Configuration settings for Eskapade.

The ConfigObject is a dictionary meant for containing global settings of Eskapade. Settings are set in the configuration macro of an analysis, or on the command line.

The ConfigObject is a dictionary meant only for storing global settings of Eskapade. In general, it is accessed through the process manager.

Example usage:

```
>>> # first set logging output level.
>>> from eskapade.logger import Logger, LogLevel
>>> logger = Logger()
>>> logger.log_level = LogLevel.DEBUG
```

Obtain the ConfigObject from any location as follows:

```
>>> from eskapade import process_manager
>>> from eskapade import ConfigObject
>>> settings = process_manager.service(ConfigObject)
```

One can treat the ConfigObject as any other dictionary:

```
>>> settings['foo'] = 'bar'
>>> foo = settings['foo']
```

Write the ConfigObject to a pickle file with:

```
>>> settings.persist_in_file(file_path)
```

And reload from the pickle file with:

```
>>> settings = ConfigObject.import_from_file(file_path)
```

A ConfigObject pickle file can be read in by Eskapade with the command line option (-u).

class `IoConfig` (**input_config)

Bases: dict

Configuration object for I/O operations.

`__init__ (**input_config)`
 Initialize IoConfig instance.

`Print ()`
 Print a summary of the settings.

`__init__ ()`
 Initialize ConfigObject instance.

`add_macros (macro_paths)`
 Add configuration macros for Eskapade run.

`copy ()`
 Perform a shallow copy of self.

Returns copy

`get (setting: str, default: typing.Any = None) → object`
 Get value of setting. If it does not exists return the default value.

Parameters

- **setting** – The setting to get.
- **default** – The default value of the setting.

Returns The value of the setting or None if it does not exist.

`io_base_dirs () → dict`
 Get configured base directories.

Returns base directories

Return type dict

`io_conf ()`
 Get I/O configuration.

The I/O configuration contains storage locations and basic analysis info.

Returns I/O configuration

Return type *IoConfig*

`set_user_opts (parsed_args)`
 Set options specified by user on command line.

Parameters **parsed_args** (*argparse.Namespace*) – parsed user arguments

class `eskapade.core.process_services.DataStore`
 Bases: `eskapade.core.process_services.ProcessService`, dict

Store for transient data sets and related objects.

The data store is a dictionary meant for storing transient data sets or any other objects. Links can take one or several data sets as input, transform them or use them as input for a model, and store the output back again in the datastore, to be picked up again by any following link.

Example usage:

```
>>> # first set logging output level.
>>> from eskapade.logger import Logger, LogLevel
>>> logger = Logger()
>>> logger.log_level = LogLevel.DEBUG
```

Obtain the global datastore from any location as follows:

```
>>> from eskapade import process_manager
>>> from eskapade import DataStore
>>>
>>> ds = process_manager.service(DataStore)
```

One can treat the datastore as any other dict:

```
>>> ds['a'] = 1
>>> ds['b'] = 2
>>> ds['0'] = 3
>>> a = ds['a']
```

Write the datastore to a pickle file with:

```
>>> ds.persist_in_file(file_path)
```

And reload from the pickle file with:

```
>>> ds = DataStore.import_from_file(file_path)
```

Print ()

Print a summary the data store contents.

class eskapade.core.process_services.**ProcessService**

Bases: object

Base class for process services.

__init__ ()

Initialize service instance.

classmethod create ()

Create an instance of this service.

Returns service instance

Return type *ProcessService*

finish ()

Finish current processes.

This function can be implemented by a process-service implementation to finish running processes and clean up to prepare for a reset of the process manager. This would typically involve deleting large objects and closing files and database connections.

classmethod import_from_file (file_path)

Import service instance from a Pickle file.

Parameters **file_path** (*str*) – path of Pickle file

Returns imported service instance

Return type *ProcessService*

Raises RuntimeError, TypeError

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.

persist_in_file (*file_path*)

Persist service instance in Pickle file.

Parameters *file_path* (*str*) – path of Pickle file

class `eskapade.core.process_services.ProcessServiceMeta`

Bases: `type`

Meta class for process-services base class.

persist

Flag to indicate if service can be persisted.

eskapade.core.run_utils module

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

Created: 2017/04/11

Description: Utilities for Eskapade run

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.core.run_utils.create_arg_parser()
    Create parser for user arguments.
```

An argparse parser is created and returned, ready to parse arguments specified by the user on the command line.

Returns argparse.ArgumentParser

Module contents

eskapade.core_ops package

Subpackages

eskapade.core_ops.links package

Submodules

eskapade.core_ops.links.assert_in_ds module

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

Class: AssertInDs

Created: 2016/11/08

Description: Algorithm that asserts that items exists in 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.core_ops.links.assert_in_ds.AssertInDs(**kwargs)
    Bases: eskapade.core.element.Link
```

Asserts that specified item(s) exists in the datastore.

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

Store the configuration of link AssertInDs

Parameters

- **name** (*str*) – name of link
- **keySet** (*lst*) – list of keys to check

execute ()
Execute the link.

eskapade.core_ops.links.break_link module

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

Class: Break

Created: 2017/02/26

Description: Algorithm to send break signal to process manager and halt execution

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.core_ops.links.break_link.Break (**kwargs)`
Bases: `eskapade.core.element.Link`

Halt execution.

Link sends failure signal and halts execution of process manager. Break the execution of the processManager at a specific location by simply adding this link at any location in a chain.

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

Parameters *name* (*str*) – name of link

execute ()
Execute the link.

eskapade.core_ops.links.ds_object_deleter module

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

Class: DsObjectDeleter

Created: 2016/11/08

Description: Algorithm to delete objects 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.core_ops.links.ds_object_deleter.DsObjectDeleter (**kwargs)`
Bases: `eskapade.core.element.Link`

Delete objects from the datastore.

Delete objects from the DataStore by the key they are under, or keeps only the data by the specified keys.

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

Parameters

- **name** (*str*) – name of link

- **deletion_keys** (*list*) – keys to clear. Overwrites clear_all to false.
- **deletion_classes** (*list*) – delete object(s) by class type.
- **keep_only** (*list*) – keys to keep. Overwrites clear_all to false.
- **clear_all** (*bool*) – clear all key-value pairs in the datastore. Default is true.

execute()

Execute the link.

initialize()

Initialize the link.

eskapade.core_ops.links.ds_to_ds module

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

Class: DsToDs

Created: 2016/11/08

Description: Algorithm to move, copy, or remove an object in 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.core_ops.links.ds_to_ds.DsToDs` (**kwargs)

Bases: `eskapade.core.element.Link`

Link to move, copy, or remove an object in the datastore.

__init__ (**kwargs)

Initialize link instance.

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
- **move** (*bool*) – move read_key item to store_key. Default is true.
- **copy** (*bool*) – if True the read_key key, value pair will not be deleted. Default is false.
- **remove** (*bool*) – if True the item corresponding to read_key key will be deleted. Default is false.
- **columnsToAdd** (*dict*) – if the object is a pandas.DataFrame columns to add to the pandas.DataFrame. key = column name, value = column

execute()

Execute the link.

initialize()

Initialize the link.

eskapade.core_ops.links.event_looper module

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

Class: EventLooper

Created: 2016/11/08

Description: EventLooper algorithm processes input lines and reprints them, e.g. to use with map/reduce

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.core_ops.links.event_looper.EventLooper` (**kwargs)

Bases: `eskapade.core.element.Link`

Event looper algorithm processes input lines and reprints or stores them.

Input lines are taken from `sys.stdin`, processed, and printed on screen.

__init__ (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **filename** (*str*) – file name where the strings are located (txt or similar). Default is None. (optional)
- **store_key** (*str*) – key to collect in datastore. If set lines are collected. (optional)
- **line_processor_set** (*list*) – list of functions to apply to input lines. (optional)
- **sort** (*bool*) – if true, sort lines before storage (optional)
- **unique** (*bool*) – if true, keep only unique lines before storage (optional),
- **skip_line_beginning_with** (*list*) – skip line if it starts with any of the list. input is list of strings. Default is ['#'] (optional)

execute ()

Process all incoming lines.

No output is printed except for lines that are passed on, such that the output lines can be picked up again by another parser.

finalize ()

Close open file if present.

initialize ()

Perform basic checks of configured attributes.

eskapade.core_ops.links.hello_world module

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

Class: HelloWorld

Created: 2017/01/31

Description: Algorithm to do print Hello {}!

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.core_ops.links.hello_world.HelloWorld (**kwargs)
    Bases: eskapade.core.element.Link
```

Defines the content of link HelloWorld.

```
__init__ (**kwargs)
    Store the configuration of link HelloWorld.
```

Parameters

- **name** (*str*) – name assigned to the link
- **hello** (*str*) – name to print in Hello World! Defaults to ‘World’
- **repeat** (*int*) – repeat print statement N times. Default is 1

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

eskapade.core_ops.links.ipython_embed module

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

Class: IPythonEmbed

Created: 2017/02/26

Description: Link that starts up an ipython console during execution for debugging.

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.core_ops.links.ipython_embed.IPythonEmbed (**kwargs)
    Bases: eskapade.core.element.Link
```

Link to start up ipython console.

Start up an ipython console by simply adding this link at any location in a chain.

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

Parameters **name** (*str*) – name of link

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

eskapade.core_ops.links.line_printer module

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

Class: LinePrinter

Created: 2017/02/21

Description: Simple algorithm to pick up lines and reprint them.

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.core_ops.links.line_printer.LinePrinter (**kwargs)
    Bases: eskapade.core.element.Link
```

LinePrinter picks up lines from the datastore and prints them.

```
__init__ (**kwargs)
    Set up the configuration of link LinePrinter.
```

Parameters

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

```
execute ()
```

Execute the link.

No output is printed except for lines that are passed on, such that the output lines can be picked up again by another parser.

```
initialize ()
```

Initialize the link.

eskapade.core_ops.links.print_ds module

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

Class: PrintDs

Created: 2016/11/08

Description: Algorithm to print the content of 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.core_ops.links.print_ds.PrintDs (**kwargs)
    Bases: eskapade.core.element.Link
```

Print the content of the datastore.

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

Parameters

- **name** (*str*) – name of link
- **keys** (*list*) – keys of items to print explicitly.

```
execute ()
```

Execute the link.

Print overview of the datastore in current state.

eskapade.core_ops.links.repeat_chain module

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

Class: RepeatChain

Created: 2016/11/08

Description: Algorithm that sends “repeat this chain” signal to processManager, until ready.

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.core_ops.links.repeat_chain.RepeatChain` (**kwargs)

Bases: `eskapade.core.element.Link`

Algorithm that sends signal to processManager to repeat the current chain.

__init__ (**kwargs)

Link that sends signal to processManager to repeat the current chain.

Sends a RepeatChain deenums.StatusCode signal.

Parameters

- **name** (*str*) – name of link
- **listen_to** (*list*) – repeat this chain if given key is present in ConfigObject and set to true. E.g. this key is set by readtods link when looping over files.
- **maxcount** (*int*) – repeat this chain until max count has been reached. Default is -1 (off).

execute ()

Execute the link.

initialize ()

Initialize the link.

eskapade.core_ops.links.skip_chain_if_empty module

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

Class: SkipChainIfEmpty

Created: 2016/11/08

Description: Algorithm to skip to the next Chain if input dataset is empty

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.core_ops.links.skip_chain_if_empty.SkipChainIfEmpty` (**kwargs)

Bases: `eskapade.core.element.Link`

Sends a SkipChain deenums.StatusCode signal when an appointed dataset is empty.

This signal causes that the processManager to step immediately to the next Chain. Input collections can be either mongo collections or dataframes in the datastore.

__init__ (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **collection_set** (*list*) – datastore keys holding the datasets to be checked. If any of these is empty, the chain is skipped.
- **skip_chain_when_key_not_in_ds** (*bool*) – skip the chain as well if the dataframe is not present in the datastore. When True and if type is 'pandas.DataFrame', sends a SkipChain signal if key not in DataStore
- **check_at_initialize** (*bool*) – perform dataset empty is check at initialize. Default is true.
- **check_at_execute** (*bool*) – perform dataset empty is check at initialize. Default is false.

check_collection_set ()

Check existence of collection in either mongo or datastore, and check that they are not empty.

Collections need to be both present and not empty.

- For mongo collections a dedicated filter can be applied before doing the count.
- For pandas dataframes the additional option 'skip_chain_when_key_not_in_ds' exists. Meaning, skip the chain as well if the dataframe is not present in the datastore.

execute ()

Execute the link.

Skip to the next Chain if any of the input dataset is empty.

initialize ()

Initialize the link.

eskapade.core_ops.links.to_ds_dict module

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

Class: ToDsDict

Created: 2016/11/08

Description: Algorithm to store one object in the DataStore dict during run time.

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.core_ops.links.to_ds_dict.ToDsDict (**kwargs)
```

Bases: *eskapade.core.element.Link*

Stores one object in the DataStore dict during run time.

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

Link to store one external object in the DataStore dict during run time.

Parameters

- **name** (*str*) – name of link
- **store_key** (*str*) – key of object to store in data store
- **obj** – object to store

- **force** (*bool*) – overwrite if already present in datastore. default is false. (optional)
- **at_initialize** (*bool*) – store at initialize of link. Default is false.
- **at_execute** (*bool*) – store at execute of link. Default is true.
- **copydict** (*bool*) – if true and obj is a dict, copy all key value pairs into datastore. Default is false.

do_storage (*ds*)

Perform storage in datastore.

Function makes a distinction between dicts and any other object.

execute ()

Execute the link.

initialize ()

Initialize the link.

Module contents

class `eskapade.core_ops.links.AssertInDs` (**kwargs)

Bases: `eskapade.core.element.Link`

Asserts that specified item(s) exists in the datastore.

__init__ (**kwargs)

Initialize link instance.

Store the configuration of link AssertInDs

Parameters

- **name** (*str*) – name of link
- **keySet** (*lst*) – list of keys to check

execute ()

Execute the link.

class `eskapade.core_ops.links.Break` (**kwargs)

Bases: `eskapade.core.element.Link`

Halt execution.

Link sends failure signal and halts execution of process manager. Break the execution of the processManager at a specific location by simply adding this link at any location in a chain.

__init__ (**kwargs)

Initialize link instance.

Parameters **name** (*str*) – name of link

execute ()

Execute the link.

class `eskapade.core_ops.links.DsObjectDeleter` (**kwargs)

Bases: `eskapade.core.element.Link`

Delete objects from the datastore.

Delete objects from the DataStore by the key they are under, or keeps only the data by the specified keys.

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

Parameters

- **name** (*str*) – name of link
- **deletion_keys** (*list*) – keys to clear. Overwrites `clear_all` to false.
- **deletion_classes** (*list*) – delete object(s) by class type.
- **keep_only** (*list*) – keys to keep. Overwrites `clear_all` to false.
- **clear_all** (*bool*) – clear all key-value pairs in the datastore. Default is true.

`execute ()`
 Execute the link.

`initialize ()`
 Initialize the link.

class `eskapade.core_ops.links.DsToDs (**kwargs)`
 Bases: `eskapade.core.element.Link`

Link to move, copy, or remove an object in the datastore.

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

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
- **move** (*bool*) – move `read_key` item to `store_key`. Default is true.
- **copy** (*bool*) – if True the `read_key` key, value pair will not be deleted. Default is false.
- **remove** (*bool*) – if True the item corresponding to `read_key` key will be deleted. Default is false.
- **columnsToAdd** (*dict*) – if the object is a `pandas.DataFrame` columns to add to the `pandas.DataFrame`. `key` = column name, `value` = column

`execute ()`
 Execute the link.

`initialize ()`
 Initialize the link.

class `eskapade.core_ops.links.EventLooper (**kwargs)`
 Bases: `eskapade.core.element.Link`

Event looper algorithm processes input lines and reprints or stores them.

Input lines are taken from `sys.stdin`, processed, and printed on screen.

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

Parameters

- **name** (*str*) – name of link

- **filename** (*str*) – file name where the strings are located (txt or similar). Default is None. (optional)
- **store_key** (*str*) – key to collect in datastore. If set lines are collected. (optional)
- **line_processor_set** (*list*) – list of functions to apply to input lines. (optional)
- **sort** (*bool*) – if true, sort lines before storage (optional)
- **unique** (*bool*) – if true, keep only unique lines before storage (optional),
- **skip_line_beginning_with** (*list*) – skip line if it starts with any of the list. input is list of strings. Default is ['#'] (optional)

execute ()

Process all incoming lines.

No output is printed except for lines that are passed on, such that the output lines can be picked up again by another parser.

finalize ()

Close open file if present.

initialize ()

Perform basic checks of configured attributes.

class `eskapade.core_ops.links.HelloWorld` (**kwargs)

Bases: `eskapade.core.element.Link`

Defines the content of link HelloWorld.

__init__ (**kwargs)

Store the configuration of link HelloWorld.

Parameters

- **name** (*str*) – name assigned to the link
- **hello** (*str*) – name to print in Hello World! Defaults to 'World'
- **repeat** (*int*) – repeat print statement N times. Default is 1

execute ()

Execute the link.

class `eskapade.core_ops.links.IPythonEmbed` (**kwargs)

Bases: `eskapade.core.element.Link`

Link to start up ipython console.

Start up an ipython console by simply adding this link at any location in a chain.

__init__ (**kwargs)

Initialize link instance.

Parameters **name** (*str*) – name of link

execute ()

Execute the link.

class `eskapade.core_ops.links.LinePrinter` (**kwargs)

Bases: `eskapade.core.element.Link`

LinePrinter picks up lines from the datastore and prints them.

__init__ (**kwargs)

Set up the configuration of link LinePrinter.

Parameters

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

execute ()

Execute the link.

No output is printed except for lines that are passed on, such that the output lines can be picked up again by another parser.

initialize ()

Initialize the link.

class `eskapade.core_ops.links.PrintDs (**kwargs)`

Bases: `eskapade.core.element.Link`

Print the content of the datastore.

__init__ (kwargs)**

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **keys** (*list*) – keys of items to print explicitly.

execute ()

Execute the link.

Print overview of the datastore in current state.

class `eskapade.core_ops.links.RepeatChain (**kwargs)`

Bases: `eskapade.core.element.Link`

Algorithm that sends signal to processManager to repeat the current chain.

__init__ (kwargs)**

Link that sends signal to processManager to repeat the current chain.

Sends a RepeatChain deenums.StatusCode signal.

Parameters

- **name** (*str*) – name of link
- **listen_to** (*list*) – repeat this chain if given key is present in ConfigObject and set to true. E.g. this key is set by readtods link when looping over files.
- **maxcount** (*int*) – repeat this chain until max count has been reached. Default is -1 (off).

execute ()

Execute the link.

initialize ()

Initialize the link.

class `eskapade.core_ops.links.SkipChainIfEmpty (**kwargs)`

Bases: `eskapade.core.element.Link`

Sends a SkipChain deenums.StatusCode signal when an appointed dataset is empty.

This signal causes that the processManager to step immediately to the next Chain. Input collections can be either mongo collections or dataframes in the datastore.

`__init__ (**kwargs)`

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **collection_set** (*list*) – datastore keys holding the datasets to be checked. If any of these is empty, the chain is skipped.
- **skip_chain_when_key_not_in_ds** (*bool*) – skip the chain as well if the dataframe is not present in the datastore. When True and if type is ‘pandas.DataFrame’, sends a SkipChain signal if key not in DataStore
- **check_at_initialize** (*bool*) – perform dataset empty is check at initialize. Default is true.
- **check_at_execute** (*bool*) – perform dataset empty is check at initialize. Default is false.

`check_collection_set ()`

Check existence of collection in either mongo or datastore, and check that they are not empty.

Collections need to be both present and not empty.

- For mongo collections a dedicated filter can be applied before doing the count.
- For pandas dataframes the additional option ‘skip_chain_when_key_not_in_ds’ exists. Meaning, skip the chain as well if the dataframe is not present in the datastore.

`execute ()`

Execute the link.

Skip to the next Chain if any of the input dataset is empty.

`initialize ()`

Initialize the link.

`class eskapade.core_ops.links.ToDsDict (**kwargs)`

Bases: `eskapade.core.element.Link`

Stores one object in the DataStore dict during run time.

`__init__ (**kwargs)`

Link to store one external object in the DataStore dict during run time.

Parameters

- **name** (*str*) – name of link
- **store_key** (*str*) – key of object to store in data store
- **obj** – object to store
- **force** (*bool*) – overwrite if already present in datastore. default is false. (optional)
- **at_initialize** (*bool*) – store at initialize of link. Default is false.
- **at_execute** (*bool*) – store at execute of link. Default is true.
- **copydict** (*bool*) – if true and obj is a dict, copy all key value pairs into datastore. Default is false.

`do_storage (ds)`

Perform storage in datastore.

Function makes a distinction been dicts and any other object.

execute()

Execute the link.

initialize()

Initialize the link.

Module contents

eskapade.data_quality package

Subpackages

eskapade.data_quality.links package

Submodules

eskapade.data_quality.links.fix_pandas_dataframe module

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

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](#)

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: `eskapade.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: `eskapade.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.

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)`
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.root_analysis package

Subpackages

eskapade.root_analysis.decorators package

Submodules

eskapade.root_analysis.decorators.histograms module

eskapade.root_analysis.decorators.roofit module

Module contents

eskapade.root_analysis.links package

Submodules

eskapade.root_analysis.links.add_propagated_error_to_roodataset module

eskapade.root_analysis.links.convert_dataframe_2_roodataset module

eskapade.root_analysis.links.convert_roodataset_2_dataframe module

eskapade.root_analysis.links.convert_roodataset_2_roodatahist module

eskapade.root_analysis.links.convert_root_hist_2_roodatahist module

eskapade.root_analysis.links.convert_root_hist_2_roodataset module

eskapade.root_analysis.links.print_ws module

eskapade.root_analysis.links.read_from_root_file module

eskapade.root_analysis.links.roodatahist_filler module

eskapade.root_analysis.links.roofit_percentile_binning module

eskapade.root_analysis.links.root_hist_filler module

eskapade.root_analysis.links.trunc_exp_fit module

eskapade.root_analysis.links.trunc_exp_gen module

eskapade.root_analysis.links.uncorrelation_hypothesis_tester module

Class : `correlation_summary`

Created: 2017/03/13

Description: Algorithm to do create correlation heatmaps.

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.visualization.links.correlation_summary.CorrelationSummary` (**kwargs)

Bases: `eskapade.core.element.Link`

Create a heatmap of correlations between dataframe variables.

__init__ (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input dataframe to read from data store
- **store_key** (*str*) – key of correlations dataframe in data store
- **results_path** (*str*) – path to save correlation summary pdf
- **methods** (*list*) – method(s) of computing correlations
- **pages_key** (*str*) – data store key of existing report pages

execute ()

Execute the link.

finalize ()

Finalize the link.

initialize ()

Initialize the link.

`eskapade.visualization.links.df_boxplot` module

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

Class : `DfBoxplot`

Created: 2017/02/17

Description: Link to create a boxplot of data frame columns.

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.visualization.links.df_boxplot.DfBoxplot` (**kwargs)

Bases: `eskapade.core.element.Link`

Create a boxplot of one column of a DataFrame that is grouped by values from a second column.

Creates a report page for each variable in DataFrame, containing:

- a profile of the column dataset

- a nicely scaled plot of the boxplots per group of the column

Example is available in: `tutorials/esk304_df_boxplot.py`

`__init__` (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **results_path** (*str*) – output path of summary result files
- **column** (*str*) – column pick up from input data to use as boxplot input
- **cause_columns** (*list*) – list of columns (str) to group-by, and per unique value plot a boxplot
- **statistics** (*list*) – a list of strings of the statistics you want to generate for the boxplot the full list is taken from `statistics.ArrayStats.get_latex_table` defaults to: ['count', 'mean', 'min', 'max']
- **pages_key** (*str*) – data store key of existing report pages

`execute` ()

Execute the link.

Creates a report page for each column that we group-by in the data frame.

- create statistics object for group
- create overview table of column variable
- plot boxplot of column variable per group
- store plot

`finalize` ()

Finalize the link.

`initialize` ()

Initialize the link.

eskapade.visualization.links.df_summary module

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

Class : DfSummary

Created: 2017/02/17

Description: Link to create a statistics summary of data frame columns or of a set of histograms.

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.visualization.links.df_summary.DfSummary` (**kwargs)

Bases: `eskapade.core.element.Link`

Create a summary of a dataframe.

Creates a report page for each variable in data frame, containing:

- a profile of the column dataset
- a nicely scaled plot of the column dataset

Example 1 is available in: `tutorials/esk301_dfsummary_plotter.py`

Example 2 is available in: `tutorials/esk303_histogram_filling_plotting.py` Empty histograms are automatically skipped from processing.

`__init__` (***kwargs*)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input dataframe (or histogram-dict) to read from data store
- **results_path** (*str*) – output path of summary result files
- **columns** (*list*) – columns (or histogram keys) pick up from input data to make & plot summaries for
- **hist_keys** (*list*) – alternative to columns (optional)
- **var_labels** (*dict*) – dict of column names with a label per column
- **var_units** (*dict*) – dict of column names with a unit per column
- **var_bins** (*dict*) – dict of column names with the number of bins per column. Default per column is 30.
- **hist_y_label** (*str*) – y-axis label to plot for all columns. Default is ‘Bin Counts’.
- **pages_key** (*str*) – data store key of existing report pages

`assert_data_type` (*data*)

Check type of input data.

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

`execute` ()

Execute the link.

Creates a report page for each variable in data frame.

- create statistics object for column
- create overview table of column variable
- plot histogram of column variable
- store plot

Returns execution status code

Return type *StatusCode*

`finalize` ()

Finalize the link.

`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_length (*data*)

Get length of data set.

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

Returns length of data set

get_sample (*data, key*)

Retrieve specific column or item from input data.

Parameters

- **data** – input data (pandas dataframe or dict)
- **key** (*str*) – column key

Returns data series or item

initialize ()

Initialize the link.

process_1d_histogram (*name, hist*)

Create statistics of and plot input 1d histogram.

Parameters

- **name** (*str*) – name of the histogram
- **hist** – input histogram object

process_2d_histogram (*name, hist*)

Create statistics of and plot input 2d histogram.

Parameters

- **name** (*str*) – name of the histogram
- **hist** – input histogram object

process_nan_histogram (*nphist, n_data*)

Process nans histogram.

Add nans histogram to pdf list

Parameters

- **nphist** – numpy-style input histogram, consisting of comma-separated bin_entries, bin_edges
- **n_data** (*int*) – number of entries in the processed data set

process_sample (*name, sample*)

Process various possible data samples.

Parameters

- **name** (*str*) – name of sample
- **sample** – input pandas series object or histogram

process_series (*col, sample*)

Create statistics of and plot input pandas series.

Parameters

- **col** (*str*) – name of the series

- **sample** – input pandas series object

Module contents

class `eskapade.visualization.links.CorrelationSummary` (**kwargs)

Bases: `eskapade.core.element.Link`

Create a heatmap of correlations between dataframe variables.

`__init__` (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input dataframe to read from data store
- **store_key** (*str*) – key of correlations dataframe in data store
- **results_path** (*str*) – path to save correlation summary pdf
- **methods** (*list*) – method(s) of computing correlations
- **pages_key** (*str*) – data store key of existing report pages

`execute` ()

Execute the link.

`finalize` ()

Finalize the link.

`initialize` ()

Initialize the link.

class `eskapade.visualization.links.DfBoxplot` (**kwargs)

Bases: `eskapade.core.element.Link`

Create a boxplot of one column of a DataFrame that is grouped by values from a second column.

Creates a report page for each variable in DataFrame, containing:

- a profile of the column dataset
- a nicely scaled plot of the boxplots per group of the column

Example is available in: `tutorials/esk304_df_boxplot.py`

`__init__` (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **results_path** (*str*) – output path of summary result files
- **column** (*str*) – column pick up from input data to use as boxplot input
- **cause_columns** (*list*) – list of columns (str) to group-by, and per unique value plot a boxplot

- **statistics** (*list*) – a list of strings of the statistics you want to generate for the boxplot the full list is taken from `statistics.ArrayStats.get_latex_table` defaults to: ['count', 'mean', 'min', 'max']
- **pages_key** (*str*) – data store key of existing report pages

execute()

Execute the link.

Creates a report page for each column that we group-by in the data frame.

- create statistics object for group
- create overview table of column variable
- plot boxplot of column variable per group
- store plot

finalize()

Finalize the link.

initialize()

Initialize the link.

class `eskapade.visualization.links.DfSummary` (**kwargs)

Bases: `eskapade.core.element.Link`

Create a summary of a dataframe.

Creates a report page for each variable in data frame, containing:

- a profile of the column dataset
- a nicely scaled plot of the column dataset

Example 1 is available in: `tutorials/esk301_dfsummary_plotter.py`

Example 2 is available in: `tutorials/esk303_histogram_filling_plotting.py` Empty histograms are automatically skipped from processing.

__init__ (**kwargs)

Initialize link instance.

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input dataframe (or histogram-dict) to read from data store
- **results_path** (*str*) – output path of summary result files
- **columns** (*list*) – columns (or histogram keys) pick up from input data to make & plot summaries for
- **hist_keys** (*list*) – alternative to columns (optional)
- **var_labels** (*dict*) – dict of column names with a label per column
- **var_units** (*dict*) – dict of column names with a unit per column
- **var_bins** (*dict*) – dict of column names with the number of bins per column. Default per column is 30.
- **hist_y_label** (*str*) – y-axis label to plot for all columns. Default is 'Bin Counts'.
- **pages_key** (*str*) – data store key of existing report pages

assert_data_type (*data*)

Check type of input data.

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

execute ()

Execute the link.

Creates a report page for each variable in data frame.

- create statistics object for column
- create overview table of column variable
- plot histogram of column variable
- store plot

Returns execution status code

Return type *StatusCode*

finalize ()

Finalize the link.

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_length (*data*)

Get length of data set.

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

Returns length of data set

get_sample (*data, key*)

Retrieve specific column or item from input data.

Parameters

- **data** – input data (pandas dataframe or dict)
- **key** (*str*) – column key

Returns data series or item

initialize ()

Initialize the link.

process_1d_histogram (*name, hist*)

Create statistics of and plot input 1d histogram.

Parameters

- **name** (*str*) – name of the histogram
- **hist** – input histogram object

process_2d_histogram (*name, hist*)

Create statistics of and plot input 2d histogram.

Parameters

- **name** (*str*) – name of the histogram
- **hist** – input histogram object

process_nan_histogram (*nphist, n_data*)

Process nans histogram.

Add nans histogram to pdf list

Parameters

- **nphist** – numpy-style input histogram, consisting of comma-separated *bin_entries*, *bin_edges*
- **n_data** (*int*) – number of entries in the processed data set

process_sample (*name, sample*)

Process various possible data samples.

Parameters

- **name** (*str*) – name of sample
- **sample** – input pandas series object or histogram

process_series (*col, sample*)

Create statistics of and plot input pandas series.

Parameters

- **col** (*str*) – name of the series
- **sample** – input pandas series object

Submodules**eskapade.visualization.vis_utils module**

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

Created: 2017/02/28

Description: Utility functions to collect Eskapade python modules e.g. functions to get correct Eskapade file paths and env variables

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.visualization.vis_utils.box_plot (df,          cause_col,          result_col='cost',
                                             pdf_file_name="",          ylim_quant=0.95,
                                             ylim_high=None,          ylim_low=0,          rot=90,
                                             statlim=400,          label_dict=None,          title_add="",
                                             top=20)
```

Make box plot.

Function that plots the boxplot of the column `df[result_col]` in groups of `cause_col`. This means that the DataFrame is grouped-by on the cause column and then the distribution per group is plotted in a boxplot using the standard pandas functionality. Boxplots with less than `statlim` (default=400) entries in it are automatically removed.

Parameters

- **df** – pandas DataFrame
- **cause_col** (*str*) – name of the column to group on. This can technically be a number, but that is uncommon.
- **result_col** (*str*) – column to do the boxplot on
- **pdf_file_name** (*str*) – if set, will store the plot in a pdf file
- **ylim_quant** (*float*) – the quantile of the y upper limit
- **ylim_high** (*float*) – when defined, this limit is used, when not defined, defaults to None and ylim_high is determined by ylim_quant
- **ylim_low** (*float*) – matplotlib set_ylim lower bound
- **rot** (*int*) – matplotlib rot
- **statlim** (*int*) – the number of entries that a group is required to have in order to be plotted
- **label_dict** (*dict*) – dictionary with labels for the columns, usage example: label_dict={'col_x': 'Time'}
- **title_add** (*str*) – string that is added to the automatic title (the y column name)
- **top** (*int*) – only print the top 20 characters of x-labels and y-labels. (default is 20)

`eskapade.visualization.vis_utils.delete_smallstat(df, group_col, statlim=400)`

Remove low-statistics groups from dataframe.

Function to make a new DataFrame that removes all groups of group_col that have less than statlim entries.

Parameters

- **df** – pandas DataFrame
- **group_col** (*str*) – name of the column to group on
- **statlim** (*int*) – number of entries a group has to have to be statistically significant

Returns smaller DataFrame and the number of removed categories

Return type tuple

`eskapade.visualization.vis_utils.plot_2d_histogram(hist, x_lim, y_lim, title, x_label, y_label, pdf_file_name)`

Plot 2d histogram with matplotlib.

Parameters

- **hist** – input numpy histogram = x_bin_edges, y_bin_edges, bin_entries_2dgrid
- **x_lim** (*tuple*) – range tuple of x-axis (min,max)
- **y_lim** (*tuple*) – range tuple of y-axis (min,max)
- **title** (*str*) – title of plot
- **x_label** (*str*) – Label for histogram x-axis
- **y_label** (*str*) – Label for histogram y-axis
- **pdf_file_name** (*str*) – if set, will store the plot in a pdf file

```

eskapade.visualization.vis_utils.plot_correlation_matrix(matrix_colors,
                                                         x_labels,      y_labels,
                                                         pdf_file_name="",  title=
                                                         'correlation',
                                                         vmin=-1,         vmax=1,
                                                         color_map='RdYlGn',
                                                         x_label="",     y_label="",
                                                         top=20,         ma-
                                                         trix_numbers=None,
                                                         print_both_numbers=True)

```

Create and plot correlation matrix.

Parameters

- **matrix_colors** – input correlation matrix
- **x_labels** (*list*) – Labels for histogram x-axis bins
- **y_labels** (*list*) – Labels for histogram y-axis bins
- **pdf_file_name** (*str*) – if set, will store the plot in a pdf file
- **title** (*str*) – if set, title of the plot
- **vmin** (*float*) – minimum value of color legend (default is -1)
- **vmax** (*float*) – maximum value of color legend (default is +1)
- **x_label** (*str*) – Label for histogram x-axis
- **y_label** (*str*) – Label for histogram y-axis
- **color_map** (*str*) – color map passed to matplotlib pcolormesh. (default is ‘RdYlGn’)
- **top** (*int*) – only print the top 20 characters of x-labels and y-labels. (default is 20)
- **matrix_numbers** – input matrix used for plotting numbers. (default it matrix_colors)

```

eskapade.visualization.vis_utils.plot_histogram(hist,      x_label,      y_label=None,
                                                is_num=True,      is_ts=False,
                                                pdf_file_name="", top=20)

```

Create and plot histogram of column values.

Parameters

- **hist** – input numpy histogram = values, bin_edges
- **x_label** (*str*) – Label for histogram x-axis
- **y_label** (*str*) – Label for histogram y-axis
- **is_num** (*bool*) – True if observable to plot is numeric
- **is_ts** (*bool*) – True if observable to plot is a timestamp
- **pdf_file_name** (*str*) – if set, will store the plot in a pdf file
- **top** (*int*) – only print the top 20 characters of x-labels and y-labels. (default is 20)

Module contents

Submodules

eskapade.entry_points module

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

Created: 2017-08-08

Description: Collection of eskapade entry points

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.entry_points.eskapade_bootstrap()`
Generate Eskapade project structure.

`eskapade.entry_points.eskapade_generate_link()`
Generate Eskapade link.

By default does not create init file.

`eskapade.entry_points.eskapade_generate_macro()`
Generate Eskapade macro.

`eskapade.entry_points.eskapade_generate_notebook()`
Generate Eskapade notebook.

`eskapade.entry_points.eskapade_ignite()`
Log info message.

`eskapade.entry_points.eskapade_run()`
Run Eskapade.

Top-level entry point for an Eskapade run started from the command line. Arguments specified by the user are parsed and converted to settings in the configuration object. Optionally, an interactive IPython session is started when the run is finished.

`eskapade.entry_points.eskapade_trial()`
Run Eskapade tests.

We will keep this here until we've completed switch to pytest or nose and tox. We could also keep it, but I don't like the fact that packages etc. are hard coded. Gotta come up with a better solution.

eskapade.exceptions module

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

Created: 2017/03/31

Description: Eskapade exceptions

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.

exception `eskapade.exceptions.MissingPackageError` (*message*="", *required_by*="")

Bases: `Exception`

Exception raised if third-party package is missing.

`__init__` (*message*="", *required_by*="")

Set missing-package arguments.

Parameters

- **message** (*str*) – message to show when raised
- **required_by** (*str*) – info on component that requires the package

exception `eskapade.exceptions.MissingPy4jError` (*message*="", *required_by*="")

Bases: `eskapade.exceptions.MissingPackageError`

Exception raised if Py4J is missing.

`__init__` (*message*="", *required_by*="")

Set missing-package arguments.

Parameters

- **message** (*str*) – message to show when raised
- **required_by** (*str*) – info on component that requires the package

exception `eskapade.exceptions.MissingRooFitError` (*message*="", *required_by*="")

Bases: `eskapade.exceptions.MissingPackageError`

Exception raised if RooFit is missing.

`__init__` (*message*="", *required_by*="")

Set missing-package arguments.

Parameters

- **message** (*str*) – message to show when raised
- **required_by** (*str*) – info on component that requires the package

exception `eskapade.exceptions.MissingRooStatsError` (*message*="", *required_by*="")

Bases: `eskapade.exceptions.MissingPackageError`

Exception raised if RooStats is missing.

`__init__` (*message*="", *required_by*="")

Set missing-package arguments.

Parameters

- **message** (*str*) – message to show when raised
- **required_by** (*str*) – info on component that requires the package

exception `eskapade.exceptions.MissingRootError` (*message*="", *required_by*="")

Bases: `eskapade.exceptions.MissingPackageError`

Exception raised if ROOT is missing.

`__init__` (*message*="", *required_by*="")

Set missing-package arguments.

Parameters

- **message** (*str*) – message to show when raised

- **required_by** (*str*) – info on component that requires the package

exception `eskapade.exceptions.MissingSparkError` (*message=*, *required_by=*)

Bases: `eskapade.exceptions.MissingPackageError`

Exception raised if Spark is missing.

__init__ (*message=*, *required_by=*)

Set missing-package arguments.

Parameters

- **message** (*str*) – message to show when raised
- **required_by** (*str*) – info on component that requires the package

eskapade.helpers module

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

Module: helpers

Created: 2017/05/17

Description: Helper functions for Eskapade

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.helpers.apply_transform_funcs` (*obj*, *trans_funcs*, *func_args=None*,
func_kwargs=None)

Transform object by applying transformation functions.

Transformation functions are applied sequentially to the output of the previous function, starting with the specified object. The final resulting object is returned.

Functions are specified with the “`trans_funcs`” argument. This is an iterable of either the functions themselves or (function, positional arguments, keyword arguments) combinations. The function arguments can also be specified by the “`func_args`” and “`func_kwargs`” arguments. The functions are called with the object to be transformed as a first argument, followed by the specified positional and keyword arguments.

A specified function can be either a callable object or a string. The latter is interpreted as a method of the type of the object on which the function is applied.

Parameters

- **obj** – object to transform
- **trans_funcs** (*iterable*) – functions to apply, specified the function or a (function, args, kwargs) tuple
- **func_args** (*dict*) – function positional arguments, specified as (function, arguments tuple) pairs
- **func_kwargs** (*dict*) – function keyword arguments, specified as (function, kwargs dict) pairs

Returns transformed object

`eskapade.helpers.obj_repr` (*obj*)

Get generic string representation of object.

`eskapade.helpers.process_transform_funcs` (*trans_funcs*, *func_args=None*,
func_kwargs=None)

Process input of the `apply_transform_funcs` function.

Parameters

- **trans_funcs** (*iterable*) – functions to apply, specified the function or a (function, args, kwargs) tuple
- **func_args** (*dict*) – function positional arguments, specified as (function, arguments tuple) pairs
- **func_kwargs** (*dict*) – function keyword arguments, specified as (function, kwargs dict) pairs

Returns transformation functions for `apply_transform_funcs` function

Return type list

eskapade.resources module

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

Created: 2017/08/23

Description: Collection of helper functions to get fixtures, i.e. test data, ROOT/RooFit libs, and tutorials. These are mostly used by the (integration) tests.

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.resources.config` (*name: str*) → str

Return the absolute path of a config.

Parameters **name** (*str*) – The name of the config.

Returns The absolute path of the config.

Raises **FileNotFoundError** – If the config cannot be found.

`eskapade.resources.fixture` (*name: str*) → str

Return the full path filename of a fixture data set.

Parameters **name** (*str*) – The name of the fixture.

Returns The full path filename of the fixture data set.

Return type str

Raises **FileNotFoundError** – If the fixture cannot be found.

`eskapade.resources.lib` (*name: str*) → str

Return the full path filename of a library.

Parameters **name** (*str*) – The name of the library.

Returns The full path filename of the library.

Return type str

Raises **FileNotFoundError** – If the library cannot be found.

`eskapade.resources.template` (*name: str*) → str

Return the full path filename of a tutorial.

Parameters `name` (*str*) – The name of the template.

Returns The full path filename of the tutorial.

Return type `str`

Raises `FileNotFoundError` – If the template cannot be found.

`eskapade.resources.tutorial` (*name: str*) → `str`
Return the full path filename of a tutorial.

Parameters `name` (*str*) – The name of the tutorial.

Returns The full path filename of the tutorial.

Return type `str`

Raises `FileNotFoundError` – If the tutorial cannot be found.

eskapade.utils module

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

Created: 2016/11/08

Description: Utility functions to collect Eskapade python modules e.g. functions to get correct Eskapade file paths and env variables

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.utils.collect_python_modules` ()
Collect Eskapade Python modules.

`eskapade.utils.get_env_var` (*key*)
Retrieve Eskapade-specific environment variables.

Parameters `key` (*str*) – Eskapade-specific key to variable

Returns environment variable value

Return type `str`

`eskapade.utils.set_matplotlib_backend` (*backend=None, batch=None, silent=True*)
Set Matplotlib backend.

Parameters

- **backend** (*str*) – backend to set
- **batch** (*bool*) – require backend to be non-interactive
- **silent** (*bool*) – do not raise exception if backend cannot be set

Raises `RuntimeError`

eskapade.version module

THIS FILE IS AUTO-GENERATED BY ESKAPADE SETUP.PY.

Module contents

3.8 Appendices

3.8.1 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.
- See *Apache Spark* for details on using Spark with Eskapade.

Migration Tips

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

The Eskapade source code can be obtained from git:

```
$ git clone git@github.com:KaveIO/Eskapade.git 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 its dependencies) by simply running:

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

Apache Spark

Eskapade supports the use of [Apache Spark](#) for parallel processing of large data volumes. Jobs can run on a single laptop using Spark libraries as well as on a Spark/Hadoop cluster in combination with YARN. This section describes

how to setup and configure Spark for use with Eskapade. For examples on running Spark jobs with Eskapade, see the [Spark tutorial](#).

Note: Eskapade supports both batch and streaming processing with Apache Spark.

Requirements

A working setup of the Apache Spark libraries is included in both the Eskapade docker and vagrant image (see section [Installation](#)). For installation of Spark libraries in a custom setup, please refer to the [Spark documentation](#).

Spark installation

The environment variables `SPARK_HOME` and `PYTHONPATH` need be set and to point to the location of the Spark installation and the Python libraries of Spark and `py4j` (dependency). In the Eskapade docker, for example, it is set to:

```
$ echo $SPARK_HOME
/opt/spark/pro/
$ echo $PYTHONPATH
/opt/spark/pro/python:/opt/spark/pro/python/lib/py4j-0.10.4-src.zip:...
```

Configuration

The Spark configuration can be set in two ways:

1. an Eskapade macro (preferred)
2. an Eskapade link

This is demonstrated in the following tutorial macro:

```
$ eskapade_run python/eskapade/tutorials/esk601_spark_configuration.py
```

Both methods are described below. For a full explanation of Spark configuration settings, see [Spark Configuration](#). In case configuration settings seem not to be picked up correctly, please check [Notes](#) at the end of this section.

Eskapade macro (preferred)

This method allows to specify settings per macro, i.e. per analysis, and is therefore the preferred way for bookkeeping analysis-specific settings.

The most easy way to start a Spark session is:

```
from eskapade import process_manager
from eskapade.spark_analysis import SparkManager

spark = sm.create_session(eskapade_settings=settings)
sc = spark.sparkContext
```

The default Spark configuration file `python/eskapade/config/spark/spark.cfg` will be picked up. It contains the following settings:

```
[spark]
spark.app.name=es_spark
spark.jars.packages=org.diana-hep:histogrammar-sparksql_2.11:1.0.4
spark.master=local[*]
spark.driver.host=localhost
```

The default Spark settings can be adapted here for all macros at once. In case, alternative settings are only relevant for a single analysis, those settings can also be specified in the macro using the argument variables in the `create_session` method of the `SparkManager`:

```
from eskapade import process_manager
from eskapade.spark_analysis import SparkManager

spark = sm.create_session(spark_settings=[('spark.app.name', 'es_spark_alt_config'), (
↳ 'spark.master', 'local[42]')])

sm = process_manager.service(SparkManager)
spark = sm.create_session(eskapade_settings=settings,
                          spark_settings=spark_settings,
                          config_path='/path/to/alternative/spark.cfg',
                          enable_hive_support=False,
                          include_eskapade_modules=False
                          )
```

Where all arguments are optional:

- `eskapade_settings` default configuration file as specified by the `sparkCfgFile` key in `ConfigObject` (i.e. `spark.cfg`)
- `config_path` alternative path to configuration file
- `spark_settings` list of key-value pairs to specify additional Spark settings
- `enable_hive_support`: switch to disable/enable Spark Hive support
- `include_eskapade_modules`: switch to include/exclude Eskapade modules in Spark job submission (e.g. for user-defined functions)

Eskapade link

This method allows to (re-)start Spark sessions from within a `SparkConfigurator` link. This means that by specifying multiple instances of this link in a macro, multiple Spark sessions with different settings can sequentially be run. This can be useful for larger analysis jobs that contain multiple Spark queries with very different CPU/memory needs - although the recently introduced *Dynamic allocation* feature is a more elegant way to achieve this behaviour.

Configurations for Spark jobs are set via the `SparkConf` class that holds a list of key/value pairs with settings, e.g.:

```
from eskapade import Chain
from eskapade.spark_analysis import SparkConfigurator

conf_link = SparkConfigurator(name='SparkConfigurator', spark_settings=[('spark.master
↳ ', 'local[3]')])
conf_link.log_level = 'INFO'
config = Chain('Config')
config.add(conf_link)
```

Note that the `SparkConfigurator` stops any existing Spark session before starting a new one. This means that the user should make sure all relevant data is stored at this point, since all cached Spark data will be cleared from memory.

Parameters

The most important parameters to play with for optimal performance:

- `num-executors`
- `executor-cores`
- `executor-memory`
- `driver-memory`

Dynamic allocation

Since version 2.1, Spark allows for [dynamic resource allocation](#). This requires the following settings:

- `spark.dynamicAllocation.enabled=true`
- `spark.shuffle.service.enabled=true`

Depending on the mode (standalone, YARN, Mesos), an additional shuffle service needs to be set up. See the documentation for details.

Logging

The logging level of Spark can be controlled in two ways:

1. through `$SPARK_HOME/conf/log4j.properties`

```
log4j.logger.org.apache.spark.api.python.PythonGatewayServer=INFO
```

2. through the `SparkContext` in Python:

```
spark = process_manager.service(SparkManager).get_session()
spark.sparkContext.setLogLevel('INFO')
```

PS: the loggers in Python can be controlled through:

```
import logging
print(logging.Logger.manager.loggerDict) # obtain list of all registered loggers
logging.getLogger('py4j').setLevel('INFO')
logging.getLogger('py4j.java_gateway').setLevel('INFO')
```

However, not all Spark-related loggers are available here (as they are JAVA-based).

Notes

There are a few pitfalls w.r.t. setting up Spark correctly:

1. If the environment variable `PYSPARK_SUBMIT_ARGS` is defined, its settings may override those specified in the macro/link. This can be prevented by unsetting the variable:

```
$ unset PYSPARK_SUBMIT_ARGS
```

or in the macro:

```
import os
del os.environ['PYSPARK_SUBMIT_ARGS']
```

The former will clear the variable from the shell session, whereas the latter will only clear it in the Python session.

2. In client mode not all driver options set via `SparkConf` are picked up at job submission because the JVM has already been started. Those settings should therefore be passed through the `SPARK_OPTS` environment variable, instead of using `SparkConf` in an Eskapade macro or link:

```
SPARK_OPTS=--driver-java-options=-Xms1024M --driver-java-options=-Xmx4096M --driver-
↪java-options=-Dlog4j.logLevel=info --driver-memory 2g
```

3. In case a Spark machine is not connected to a network, setting the `SPARK_LOCAL_HOSTNAME` environment variable or the `spark.driver.host` key in `SparkConf` to the value `localhost` may fix DNS resolution timeouts which prevent Spark from starting jobs.

3.9 Indices and tables

- [genindex](#)
- [modindex](#)

e

- eskapade, 3
- eskapade.analysis, 82
- eskapade.analysis.datetime, 65
- eskapade.analysis.histogram, 67
- eskapade.analysis.histogram_filling, 76
- eskapade.analysis.links, 56
- eskapade.analysis.links.apply_func_to_df, 44
- eskapade.analysis.links.apply_selection_to_df, 45
- eskapade.analysis.links.basic_generator, 46
- eskapade.analysis.links.df_concatenator, 46
- eskapade.analysis.links.df_merger, 47
- eskapade.analysis.links.histogrammar_filler, 48
- eskapade.analysis.links.random_sample_splitter, 49
- eskapade.analysis.links.read_to_df, 50
- eskapade.analysis.links.record_factorizer, 51
- eskapade.analysis.links.record_vectorizer, 53
- eskapade.analysis.links.value_counter, 54
- eskapade.analysis.links.write_from_df, 56
- eskapade.analysis.statistics, 79
- eskapade.core, 101
- eskapade.core.definitions, 82
- eskapade.core.element, 84
- eskapade.core.exceptions, 87
- eskapade.core.execution, 88
- eskapade.core.meta, 88
- eskapade.core.mixin, 91
- eskapade.core.persistence, 93
- eskapade.core.process_manager, 94
- eskapade.core.process_services, 97
- eskapade.core.run_utils, 101
- eskapade.core_ops, 114
- eskapade.core_ops.links, 109
- eskapade.core_ops.links.assert_in_ds, 101
- eskapade.core_ops.links.break_link, 102
- eskapade.core_ops.links.ds_object_deleter, 102
- eskapade.core_ops.links.ds_to_ds, 103
- eskapade.core_ops.links.event_looper, 104
- eskapade.core_ops.links.hello_world, 104
- eskapade.core_ops.links.ipython_embed, 105
- eskapade.core_ops.links.line_printer, 105
- eskapade.core_ops.links.print_ds, 106
- eskapade.core_ops.links.repeat_chain, 107
- eskapade.core_ops.links.skip_chain_if_empty, 107
- eskapade.core_ops.links.to_ds_dict, 108
- eskapade.data_quality, 121
- eskapade.data_quality.dq_helper, 118
- eskapade.data_quality.links, 116
- eskapade.data_quality.links.fix_pandas_dataframe, 114
- eskapade.entry_points, 132
- eskapade.exceptions, 132
- eskapade.helpers, 134
- eskapade.logger, 121
- eskapade.resources, 135
- eskapade.utils, 136
- eskapade.version, 136
- eskapade.visualization, 132
- eskapade.visualization.links, 126
- eskapade.visualization.links.correlation_summary, 121
- eskapade.visualization.links.df_boxplot,

122

`eskapade.visualization.links.df_summary,`

123

`eskapade.visualization.vis_utils,` 129

Symbols

- `__init__()` (eskapade.analysis.datetime.FreqTimePeriod method), 66
- `__init__()` (eskapade.analysis.datetime.TimePeriod method), 66
- `__init__()` (eskapade.analysis.datetime.UniformTsTimePeriod method), 67
- `__init__()` (eskapade.analysis.histogram.BinningUtil method), 68
- `__init__()` (eskapade.analysis.histogram.Histogram method), 69
- `__init__()` (eskapade.analysis.histogram.ValueCounts method), 74
- `__init__()` (eskapade.analysis.histogram_filling.HistogramFillerBase method), 76
- `__init__()` (eskapade.analysis.links.ApplyFuncToDf method), 57
- `__init__()` (eskapade.analysis.links.ApplySelectionToDf method), 57
- `__init__()` (eskapade.analysis.links.BasicGenerator method), 58
- `__init__()` (eskapade.analysis.links.DfConcatenator method), 58
- `__init__()` (eskapade.analysis.links.DfMerger method), 58
- `__init__()` (eskapade.analysis.links.HistogrammarFiller method), 59
- `__init__()` (eskapade.analysis.links.RandomSampleSplitter method), 60
- `__init__()` (eskapade.analysis.links.ReadToDf method), 61
- `__init__()` (eskapade.analysis.links.RecordFactorizer method), 62
- `__init__()` (eskapade.analysis.links.RecordVectorizer method), 63
- `__init__()` (eskapade.analysis.links.ValueCounter method), 63
- `__init__()` (eskapade.analysis.links.WriteFromDf method), 65
- `__init__()` (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 44
- `__init__()` (eskapade.analysis.links.apply_selection_to_df.ApplySelectionToDf method), 45
- `__init__()` (eskapade.analysis.links.basic_generator.BasicGenerator method), 46
- `__init__()` (eskapade.analysis.links.df_concatenator.DfConcatenator method), 47
- `__init__()` (eskapade.analysis.links.df_merger.DfMerger method), 47
- `__init__()` (eskapade.analysis.links.histogrammar_filler.HistogrammarFiller method), 48
- `__init__()` (eskapade.analysis.links.random_sample_splitter.RandomSampleSplitter method), 50
- `__init__()` (eskapade.analysis.links.read_to_df.ReadToDf method), 51
- `__init__()` (eskapade.analysis.links.record_factorizer.RecordFactorizer method), 52
- `__init__()` (eskapade.analysis.links.record_vectorizer.RecordVectorizer method), 53
- `__init__()` (eskapade.analysis.links.value_counter.ValueCounter method), 54
- `__init__()` (eskapade.analysis.links.write_from_df.WriteFromDf method), 56
- `__init__()` (eskapade.analysis.statistics.ArrayStats method), 79
- `__init__()` (eskapade.analysis.statistics.GroupByStats method), 81
- `__init__()` (eskapade.core.definitions.RandomSeeds method), 83
- `__init__()` (eskapade.core.element.Chain method), 85
- `__init__()` (eskapade.core.element.Link method), 86
- `__init__()` (eskapade.core.meta.Processor method), 89
- `__init__()` (eskapade.core.meta.ProcessorSequence method), 90
- `__init__()` (eskapade.core.mixin.ConfigMixin method), 92
- `__init__()` (eskapade.core.mixin.TimerMixin method), 92
- `__init__()` (eskapade.core.process_manager.ProcessManager method), 94

[__init__\(\) \(eskapade.core.process_services.ConfigObject method\), 98](#)
[__init__\(\) \(eskapade.core.process_services.ConfigObject.IoConfig method\), 97](#)
[__init__\(\) \(eskapade.core.process_services.ProcessService method\), 99](#)
[__init__\(\) \(eskapade.core_ops.links.AssertInDs method\), 109](#)
[__init__\(\) \(eskapade.core_ops.links.Break method\), 109](#)
[__init__\(\) \(eskapade.core_ops.links.DsObjectDeleter method\), 109](#)
[__init__\(\) \(eskapade.core_ops.links.DsToDs method\), 110](#)
[__init__\(\) \(eskapade.core_ops.links.EventLooper method\), 110](#)
[__init__\(\) \(eskapade.core_ops.links.HelloWorld method\), 111](#)
[__init__\(\) \(eskapade.core_ops.links.IPythonEmbed method\), 111](#)
[__init__\(\) \(eskapade.core_ops.links.LinePrinter method\), 111](#)
[__init__\(\) \(eskapade.core_ops.links.PrintDs method\), 112](#)
[__init__\(\) \(eskapade.core_ops.links.RepeatChain method\), 112](#)
[__init__\(\) \(eskapade.core_ops.links.SkipChainIfEmpty method\), 112](#)
[__init__\(\) \(eskapade.core_ops.links.ToDsDict method\), 113](#)
[__init__\(\) \(eskapade.core_ops.links.assert_in_ds.AssertInDs method\), 101](#)
[__init__\(\) \(eskapade.core_ops.links.break_link.Break method\), 102](#)
[__init__\(\) \(eskapade.core_ops.links.ds_object_deleter.DsObjectDeleter method\), 102](#)
[__init__\(\) \(eskapade.core_ops.links.ds_to_ds.DsToDs method\), 103](#)
[__init__\(\) \(eskapade.core_ops.links.event_looper.EventLooper method\), 104](#)
[__init__\(\) \(eskapade.core_ops.links.hello_world.HelloWorld method\), 105](#)
[__init__\(\) \(eskapade.core_ops.links.ipython_embed.IPythonEmbed method\), 105](#)
[__init__\(\) \(eskapade.core_ops.links.line_printer.LinePrinter method\), 106](#)
[__init__\(\) \(eskapade.core_ops.links.print_ds.PrintDs method\), 106](#)
[__init__\(\) \(eskapade.core_ops.links.repeat_chain.RepeatChain method\), 107](#)
[__init__\(\) \(eskapade.core_ops.links.skip_chain_if_empty.SkipChainIfEmpty method\), 107](#)
[__init__\(\) \(eskapade.core_ops.links.to_ds_dict.ToDsDict method\), 108](#)
[__init__\(\) \(eskapade.data_quality.links.FixPandasDataFrame method\), 117](#)
[__init__\(\) \(eskapade.data_quality.links.fix_pandas_dataframe.FixPandasDataFrame method\), 115](#)
[__init__\(\) \(eskapade.exceptions.MissingPackageError method\), 133](#)
[__init__\(\) \(eskapade.exceptions.MissingPy4jError method\), 133](#)
[__init__\(\) \(eskapade.exceptions.MissingRooFitError method\), 133](#)
[__init__\(\) \(eskapade.exceptions.MissingRooStatsError method\), 133](#)
[__init__\(\) \(eskapade.exceptions.MissingRootError method\), 133](#)
[__init__\(\) \(eskapade.exceptions.MissingSparkError method\), 134](#)
[__init__\(\) \(eskapade.visualization.links.CorrelationSummary method\), 126](#)
[__init__\(\) \(eskapade.visualization.links.DfBoxplot method\), 126](#)
[__init__\(\) \(eskapade.visualization.links.DfSummary method\), 127](#)
[__init__\(\) \(eskapade.visualization.links.correlation_summary.CorrelationSummary method\), 122](#)
[__init__\(\) \(eskapade.visualization.links.df_boxplot.DfBoxplot method\), 123](#)
[__init__\(\) \(eskapade.visualization.links.df_summary.DfSummary method\), 124](#)

A

[add\(\) \(eskapade.core.element.Chain method\), 85](#)
[add\(\) \(eskapade.core.meta.ProcessorSequence method\), 90](#)
[add\(\) \(eskapade.core.process_manager.ProcessManager method\), 94](#)
[add_apply_func\(\) \(eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method\), 45](#)
[add_apply_func\(\) \(eskapade.analysis.links.ApplyFuncToDf method\), 57](#)
[add_macros\(\) \(eskapade.core.process_services.ConfigObject method\), 98](#)
[apply_transform_funcs\(\) \(in module eskapade.helpers\), 134](#)
[ApplyFuncToDf \(class in eskapade.analysis.links\), 56](#)
[ApplyFuncToDf \(class in eskapade.analysis.links.apply_func_to_df\), 44](#)
[ApplySelectionToDf \(class in eskapade.analysis.links\), 57](#)
[ApplySelectionToDf \(class in eskapade.analysis.links.apply_selection_to_df\), 45](#)
[ArgumentsMixin \(class in eskapade.core.mixin\), 91](#)
[ArrayStats \(class in eskapade.analysis.statistics\), 79](#)

`assert_data_type()` (eskapade.visualization.links.df_summary.DfSummary method), 124
`assert_data_type()` (eskapade.visualization.links.DfSummary method), 127
`assert_dataframe()` (eskapade.analysis.histogram_filling.HistogramFillerBase method), 77
`AssertInDs` (class in eskapade.core_ops.links), 109
`AssertInDs` (class in eskapade.core_ops.links.assert_in_ds), 101

B

`BasicGenerator` (class in eskapade.analysis.links), 58
`BasicGenerator` (class in eskapade.analysis.links.basic_generator), 46
`bin_centers()` (eskapade.analysis.histogram.Histogram method), 70
`bin_edges()` (eskapade.analysis.histogram.Histogram method), 70
`bin_entries()` (eskapade.analysis.histogram.Histogram method), 70
`bin_labels()` (eskapade.analysis.histogram.Histogram method), 70
`bin_specs` (eskapade.analysis.histogram.BinningUtil attribute), 68
`BinningUtil` (class in eskapade.analysis.histogram), 68
`bool_to_int()` (in module eskapade.data_quality.dq_helper), 118
`bool_to_str()` (in module eskapade.data_quality.dq_helper), 118
`box_plot()` (in module eskapade.visualization.vis_utils), 129
`Break` (class in eskapade.core_ops.links), 109
`Break` (class in eskapade.core_ops.links.break_link), 102
`BreakChain` (eskapade.core.definitions.StatusCode attribute), 83

C

`categorize_columns()` (eskapade.analysis.histogram_filling.HistogramFillerBase method), 77
`Chain` (class in eskapade.core.element), 85
`check_arg_callable()` (eskapade.core.mixin.ArgumentsMixin method), 91
`check_arg_iters()` (eskapade.core.mixin.ArgumentsMixin method), 91
`check_arg_opts()` (eskapade.core.mixin.ArgumentsMixin method), 91
`check_arg_types()` (eskapade.core.mixin.ArgumentsMixin method), 91
`check_arg_vals()` (eskapade.core.mixin.ArgumentsMixin method), 91
`check_collection_set()` (eskapade.core_ops.links.skip_chain_if_empty.SkipChainIfEmpty method), 108
`check_collection_set()` (eskapade.core_ops.links.SkipChainIfEmpty method), 113
`check_extra_kwargs()` (eskapade.core.mixin.ArgumentsMixin method), 91
`check_nan()` (in module eskapade.data_quality.dq_helper), 118
`check_required_args()` (eskapade.core.mixin.ArgumentsMixin method), 92
`cleanup_string()` (in module eskapade.data_quality.dq_helper), 118
`clear()` (eskapade.core.element.Chain method), 85
`clear()` (eskapade.core.meta.ProcessorSequence method), 90
`clear()` (eskapade.core.process_manager.ProcessManager method), 95
`collect_python_modules()` (in module eskapade.utils), 136
`combine_hists()` (eskapade.analysis.histogram.Histogram class method), 70
`config()` (in module eskapade.resources), 135
`config_path` (eskapade.core.mixin.ConfigMixin attribute), 92
`ConfigMixin` (class in eskapade.core.mixin), 92
`ConfigObject` (class in eskapade.core.process_services), 97
`ConfigObject.IoConfig` (class in eskapade.core.process_services), 97
`construct_empty_hist()` (eskapade.analysis.links.histogrammar_filler.HistogrammarFiller method), 49
`construct_empty_hist()` (eskapade.analysis.links.HistogrammarFiller method), 60
`convert()` (in module eskapade.data_quality.dq_helper), 119
`copy()` (eskapade.analysis.histogram.Histogram method), 70
`copy()` (eskapade.core.process_services.ConfigObject method), 98
`CorrelationSummary` (class in eskapade.visualization.links), 126
`CorrelationSummary` (class in eskapade.visualization.links.correlation_summary), 122
`count()` (eskapade.analysis.histogram.ValueCounts method), 74

counts (eskapade.analysis.histogram.ValueCounts attribute), 74
 create() (eskapade.core.process_services.ProcessService class method), 99
 create_arg_parser() (in module eskapade.core.run_utils), 101
 create_dir() (in module eskapade.core.persistence), 93
 create_mpv_stat() (eskapade.analysis.statistics.ArrayStats method), 80
 create_stats() (eskapade.analysis.statistics.ArrayStats method), 80
 create_sub_counts() (eskapade.analysis.histogram.ValueCounts method), 74

D

DataStore (class in eskapade.core.process_services), 98
 datatype (eskapade.analysis.histogram.Histogram attribute), 71
 delete_smallstat() (in module eskapade.visualization.vis_utils), 130
 determine_preferred_dtype() (in module eskapade.data_quality.links.fix_pandas_dataframe), 116
 DfBoxplot (class in eskapade.visualization.links), 126
 DfBoxplot (class in eskapade.visualization.links.df_boxplot), 122
 DfConcatenator (class in eskapade.analysis.links), 58
 DfConcatenator (class in eskapade.analysis.links.df_concatenator), 46
 DfMerger (class in eskapade.analysis.links), 58
 DfMerger (class in eskapade.analysis.links.df_merger), 47
 DfSummary (class in eskapade.visualization.links), 127
 DfSummary (class in eskapade.visualization.links.df_summary), 123
 discard() (eskapade.core.element.Chain method), 85
 discard() (eskapade.core.meta.ProcessorSequence method), 90
 do_storage() (eskapade.core_ops.links.to_ds_dict.ToDsDict method), 109
 do_storage() (eskapade.core_ops.links.ToDsDict method), 113
 drop_inconsistent_keys() (eskapade.analysis.links.value_counter.ValueCounter method), 55
 drop_inconsistent_keys() (eskapade.analysis.links.ValueCounter method), 64
 drop_requested_keys() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 77
 DsObjectDeleter (class in eskapade.core_ops.links), 109

DsObjectDeleter (class in eskapade.core_ops.links.ds_object_deleter), 102
 DsToDs (class in eskapade.core_ops.links), 110
 DsToDs (class in eskapade.core_ops.links.ds_to_ds), 103
 dt_string() (eskapade.analysis.datetime.FreqTimePeriod method), 66

E

Error, 88
 eskapade (module), 1, 137
 eskapade.analysis (module), 82
 eskapade.analysis.datetime (module), 65
 eskapade.analysis.histogram (module), 67
 eskapade.analysis.histogram_filling (module), 76
 eskapade.analysis.links (module), 56
 eskapade.analysis.links.apply_func_to_df (module), 44
 eskapade.analysis.links.apply_selection_to_df (module), 45
 eskapade.analysis.links.basic_generator (module), 46
 eskapade.analysis.links.df_concatenator (module), 46
 eskapade.analysis.links.df_merger (module), 47
 eskapade.analysis.links.histogrammer_filler (module), 48
 eskapade.analysis.links.random_sample_splitter (module), 49
 eskapade.analysis.links.read_to_df (module), 50
 eskapade.analysis.links.record_factorizer (module), 51
 eskapade.analysis.links.record_vectorizer (module), 53
 eskapade.analysis.links.value_counter (module), 54
 eskapade.analysis.links.write_from_df (module), 56
 eskapade.analysis.statistics (module), 79
 eskapade.core (module), 101
 eskapade.core.definitions (module), 82
 eskapade.core.element (module), 84
 eskapade.core.exceptions (module), 87
 eskapade.core.execution (module), 88
 eskapade.core.meta (module), 88
 eskapade.core.mixin (module), 91
 eskapade.core.persistence (module), 93
 eskapade.core.process_manager (module), 94
 eskapade.core.process_services (module), 97
 eskapade.core.run_utils (module), 101
 eskapade.core_ops (module), 114
 eskapade.core_ops.links (module), 109
 eskapade.core_ops.links.assert_in_ds (module), 101
 eskapade.core_ops.links.break_link (module), 102
 eskapade.core_ops.links.ds_object_deleter (module), 102
 eskapade.core_ops.links.ds_to_ds (module), 103
 eskapade.core_ops.links.event_looper (module), 104
 eskapade.core_ops.links.hello_world (module), 104
 eskapade.core_ops.links.ipython_embed (module), 105
 eskapade.core_ops.links.line_printer (module), 105
 eskapade.core_ops.links.print_ds (module), 106
 eskapade.core_ops.links.repeat_chain (module), 107

- eskapade.core_ops.links.skip_chain_if_empty (module), 107
- eskapade.core_ops.links.to_ds_dict (module), 108
- eskapade.data_quality (module), 121
- eskapade.data_quality.dq_helper (module), 118
- eskapade.data_quality.links (module), 116
- eskapade.data_quality.links.fix_pandas_dataframe (module), 114
- eskapade.entry_points (module), 132
- eskapade.exceptions (module), 132
- eskapade.helpers (module), 134
- eskapade.logger (module), 121
- eskapade.resources (module), 135
- eskapade.utils (module), 136
- eskapade.version (module), 136
- eskapade.visualization (module), 132
- eskapade.visualization.links (module), 126
- eskapade.visualization.links.correlation_summary (module), 121
- eskapade.visualization.links.df_boxplot (module), 122
- eskapade.visualization.links.df_summary (module), 123
- eskapade.visualization.vis_utils (module), 129
- eskapade_bootstrap() (in module eskapade.entry_points), 132
- eskapade_generate_link() (in module eskapade.entry_points), 132
- eskapade_generate_macro() (in module eskapade.entry_points), 132
- eskapade_generate_notebook() (in module eskapade.entry_points), 132
- eskapade_ignite() (in module eskapade.entry_points), 132
- eskapade_run() (in module eskapade.core.execution), 88
- eskapade_run() (in module eskapade.entry_points), 132
- eskapade_trial() (in module eskapade.entry_points), 132
- EventLooper (class in eskapade.core_ops.links), 110
- EventLooper (class in eskapade.core_ops.links.event_looper), 104
- execute() (eskapade.analysis.histogram_filling.HistogramFilling method), 77
- execute() (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 45
- execute() (eskapade.analysis.links.apply_selection_to_df.ApplySelectionToDf method), 45
- execute() (eskapade.analysis.links.ApplyFuncToDf method), 57
- execute() (eskapade.analysis.links.ApplySelectionToDf method), 57
- execute() (eskapade.analysis.links.basic_generator.BasicGenerator method), 46
- execute() (eskapade.analysis.links.BasicGenerator method), 58
- execute() (eskapade.analysis.links.df_concatenator.DfConcatenator method), 47
- execute() (eskapade.analysis.links.df_merger.DfMerger method), 48
- execute() (eskapade.analysis.links.DfConcatenator method), 58
- execute() (eskapade.analysis.links.DfMerger method), 59
- execute() (eskapade.analysis.links.random_sample_splitter.RandomSampleSplitter method), 50
- execute() (eskapade.analysis.links.RandomSampleSplitter method), 61
- execute() (eskapade.analysis.links.read_to_df.ReadToDf method), 51
- execute() (eskapade.analysis.links.ReadToDf method), 61
- execute() (eskapade.analysis.links.record_factorizer.RecordFactorizer method), 52
- execute() (eskapade.analysis.links.record_vectorizer.RecordVectorizer method), 53
- execute() (eskapade.analysis.links.RecordFactorizer method), 62
- execute() (eskapade.analysis.links.RecordVectorizer method), 63
- execute() (eskapade.analysis.links.write_from_df.WriteFromDf method), 56
- execute() (eskapade.analysis.links.WriteFromDf method), 65
- execute() (eskapade.core.element.Chain method), 85
- execute() (eskapade.core.element.Link method), 87
- execute() (eskapade.core.meta.Processor method), 89
- execute() (eskapade.core.process_manager.ProcessManager method), 95
- execute() (eskapade.core_ops.links.assert_in_ds.AssertInDs method), 101
- execute() (eskapade.core_ops.links.AssertInDs method), 109
- execute() (eskapade.core_ops.links.Break method), 109
- execute() (eskapade.core_ops.links.break_link.Break method), 102
- execute() (eskapade.core_ops.links.ds_object_deleter.DsObjectDeleter method), 103
- execute() (eskapade.core_ops.links.ds_to_ds.DsToDs method), 103
- execute() (eskapade.core_ops.links.DsObjectDeleter method), 110
- execute() (eskapade.core_ops.links.DsToDs method), 110
- execute() (eskapade.core_ops.links.event_looper.EventLooper method), 104
- execute() (eskapade.core_ops.links.EventLooper method), 111
- execute() (eskapade.core_ops.links.hello_world.HelloWorld method), 105
- execute() (eskapade.core_ops.links.HelloWorld method), 111
- execute() (eskapade.core_ops.links.ipython_embed.IPythonEmbed method), 105
- execute() (eskapade.core_ops.links.IPythonEmbed method), 105

method), 111

execute() (eskapade.core_ops.links.line_printer.LinePrinter method), 106

execute() (eskapade.core_ops.links.LinePrinter method), 112

execute() (eskapade.core_ops.links.print_ds.PrintDs method), 106

execute() (eskapade.core_ops.links.PrintDs method), 112

execute() (eskapade.core_ops.links.repeat_chain.RepeatChain method), 107

execute() (eskapade.core_ops.links.RepeatChain method), 112

execute() (eskapade.core_ops.links.skip_chain_if_empty.SkipChainIfEmpty method), 108

execute() (eskapade.core_ops.links.SkipChainIfEmpty method), 113

execute() (eskapade.core_ops.links.to_ds_dict.ToDsDict method), 109

execute() (eskapade.core_ops.links.ToDsDict method), 113

execute() (eskapade.data_quality.links.fix_pandas_dataframe.FixPandasDataFrame method), 116

execute() (eskapade.data_quality.links.FixPandasDataFrame method), 118

execute() (eskapade.visualization.links.correlation_summary.CorrelationSummary method), 122

execute() (eskapade.visualization.links.CorrelationSummary method), 126

execute() (eskapade.visualization.links.df_boxplot.DfBoxplot method), 123

execute() (eskapade.visualization.links.df_summary.DfSummary method), 124

execute() (eskapade.visualization.links.DfBoxplot method), 127

execute() (eskapade.visualization.links.DfSummary method), 128

execute_macro() (eskapade.core.process_manager.ProcessManager method), 95

finalize() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 77

finalize() (eskapade.analysis.links.value_counter.ValueCounter method), 55

finalize() (eskapade.analysis.links.ValueCounter method), 64

finalize() (eskapade.core.element.Chain method), 86

finalize() (eskapade.core.element.Link method), 87

finalize() (eskapade.core.meta.Processor method), 89

finalize() (eskapade.core.process_manager.ProcessManager method), 95

finalize() (eskapade.core_ops.links.event_looper.EventLooper method), 104

finalize() (eskapade.core_ops.links.EventLooper method), 111

finalize() (eskapade.visualization.links.correlation_summary.CorrelationSummary method), 122

finalize() (eskapade.visualization.links.CorrelationSummary method), 126

finalize() (eskapade.visualization.links.df_boxplot.DfBoxplot method), 123

finalize() (eskapade.visualization.links.df_summary.DfSummary method), 124

finalize() (eskapade.visualization.links.DfBoxplot method), 127

finalize() (eskapade.visualization.links.DfSummary method), 128

finish() (eskapade.core.process_services.ProcessService method), 99

FixPandasDataFrame (class in eskapade.data_quality.links), 116

FixPandasDataFrame (class in eskapade.data_quality.links.fix_pandas_dataframe), 114

fixture() (in module eskapade.resources), 135

freq (eskapade.analysis.datetime.FreqTimePeriod attribute), 66

FreqTimePeriod (class in eskapade.analysis.datetime), 66

F

Failure (eskapade.core.definitions.StatusCode attribute), 83

fill_histogram() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 77

fill_histogram() (eskapade.analysis.links.histogrammar_filler.HistogrammarFiller method), 49

fill_histogram() (eskapade.analysis.links.HistogrammarFiller method), 60

fill_histogram() (eskapade.analysis.links.value_counter.ValueCounter method), 55

fill_histogram() (eskapade.analysis.links.ValueCounter method), 64

get() (eskapade.core.element.Chain method), 86

get() (eskapade.core.process_manager.ProcessManager method), 95

get() (eskapade.core.process_services.ConfigObject method), 98

get_all_columns() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 77

get_all_columns() (eskapade.visualization.links.df_summary.DfSummary method), 124

get_all_columns() (eskapade.visualization.links.DfSummary method), 128

G

[get_bin_center\(\)](#) (eskapade.analysis.histogram.BinningUtil method), 68
[get_bin_count\(\)](#) (eskapade.analysis.histogram.Histogram method), 71
[get_bin_edges\(\)](#) (eskapade.analysis.histogram.BinningUtil method), 68
[get_bin_edges_range\(\)](#) (eskapade.analysis.histogram.BinningUtil method), 68
[get_bin_labels\(\)](#) (eskapade.analysis.histogram.Histogram method), 71
[get_bin_range\(\)](#) (eskapade.analysis.histogram.Histogram method), 71
[get_bin_vals\(\)](#) (eskapade.analysis.histogram.Histogram method), 71
[get_col_props\(\)](#) (eskapade.analysis.statistics.ArrayStats method), 80
[get_col_props\(\)](#) (in module eskapade.analysis.statistics), 82
[get_config\(\)](#) (eskapade.core.mixin.ConfigMixin method), 92
[get_data_type\(\)](#) (eskapade.analysis.histogram_filling.HistogramFillerBase method), 78
[get_env_var\(\)](#) (in module eskapade.utils), 136
[get_hist_val\(\)](#) (eskapade.analysis.histogram.Histogram method), 71
[get_latex_table\(\)](#) (eskapade.analysis.statistics.ArrayStats method), 80
[get_latex_table\(\)](#) (eskapade.analysis.statistics.GroupByStats method), 82
[get_left_bin_edge\(\)](#) (eskapade.analysis.histogram.BinningUtil method), 69
[get_length\(\)](#) (eskapade.visualization.links.df_summary.DfSummary method), 125
[get_length\(\)](#) (eskapade.visualization.links.DfSummary method), 128
[get_nonone_bin_centers\(\)](#) (eskapade.analysis.histogram.Histogram method), 71
[get_nonone_bin_counts\(\)](#) (eskapade.analysis.histogram.Histogram method), 71
[get_nonone_bin_edges\(\)](#) (eskapade.analysis.histogram.Histogram method), 72
[get_nonone_bin_range\(\)](#) (eskapade.analysis.histogram.Histogram method), 72
[get_print_stats\(\)](#) (eskapade.analysis.statistics.ArrayStats method), 80
[get_right_bin_edge\(\)](#) (eskapade.analysis.histogram.BinningUtil method), 69
[get_sample\(\)](#) (eskapade.visualization.links.df_summary.DfSummary method), 125
[get_sample\(\)](#) (eskapade.visualization.links.DfSummary method), 128
[get_service_tree\(\)](#) (eskapade.core.process_manager.ProcessManager method), 95
[get_services\(\)](#) (eskapade.core.process_manager.ProcessManager method), 95
[get_uniform_bin_edges\(\)](#) (eskapade.analysis.histogram.Histogram method), 72
[get_values\(\)](#) (eskapade.analysis.histogram.ValueCounts method), 75
[get_x_label\(\)](#) (eskapade.analysis.statistics.ArrayStats method), 80
[groupbyapply\(\)](#) (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 45
[groupbyapply\(\)](#) (eskapade.analysis.links.ApplyFuncToDf method), 57
[GroupByStats](#) (class in eskapade.analysis.statistics), 81

H

[HelloWorld](#) (class in eskapade.core_ops.links), 111
[HelloWorld](#) (class in eskapade.core_ops.links.hello_world), 105
[Histogram](#) (class in eskapade.analysis.histogram), 69
[HistogramFillerBase](#) (class in eskapade.analysis.histogram_filling), 76
[HistogrammarFiller](#) (class in eskapade.analysis.links), 59
[HistogrammarFiller](#) (class in eskapade.analysis.links.histogrammar_filler), 48

I

[import_from_file\(\)](#) (eskapade.core.process_services.ProcessService class method), 99
[import_services\(\)](#) (eskapade.core.process_manager.ProcessManager method), 95
[initialize\(\)](#) (eskapade.analysis.histogram_filling.HistogramFillerBase method), 78
[initialize\(\)](#) (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 45
[initialize\(\)](#) (eskapade.analysis.links.apply_selection_to_df.ApplySelectionToDf method), 46
[initialize\(\)](#) (eskapade.analysis.links.ApplyFuncToDf method), 57
[initialize\(\)](#) (eskapade.analysis.links.ApplySelectionToDf method), 58

initialize() (eskapade.analysis.links.basic_generator.BasicGenerator method), 106
 method), 46
 initialize() (eskapade.analysis.links.BasicGenerator method), 58
 initialize() (eskapade.analysis.links.df_concatenator.DfConcatenator method), 107
 method), 47
 initialize() (eskapade.analysis.links.df_merger.DfMerger method), 48
 initialize() (eskapade.analysis.links.DfConcatenator method), 58
 initialize() (eskapade.analysis.links.DfMerger method), 59
 initialize() (eskapade.analysis.links.random_sample_splitter.RandomSampleSplitter method), 50
 initialize() (eskapade.analysis.links.RandomSampleSplitter method), 61
 initialize() (eskapade.analysis.links.read_to_df.ReadToDf method), 51
 initialize() (eskapade.analysis.links.ReadToDf method), 61
 initialize() (eskapade.analysis.links.record_factorizer.RecordFactorizer method), 122
 method), 52
 initialize() (eskapade.analysis.links.record_vectorizer.RecordVectorizer method), 126
 method), 53
 initialize() (eskapade.analysis.links.RecordFactorizer method), 63
 initialize() (eskapade.analysis.links.RecordVectorizer method), 63
 initialize() (eskapade.analysis.links.value_counter.ValueCounter method), 55
 initialize() (eskapade.analysis.links.ValueCounter method), 65
 initialize() (eskapade.analysis.links.write_from_df.WriteFromDf method), 56
 initialize() (eskapade.analysis.links.WriteFromDf method), 65
 initialize() (eskapade.core.element.Chain method), 86
 initialize() (eskapade.core.element.Link method), 87
 initialize() (eskapade.core.meta.Processor method), 89
 initialize() (eskapade.core.process_manager.ProcessManager method), 96
 initialize() (eskapade.core_ops.links.ds_object_deleter.DsObjectDeleter method), 83
 method), 103
 initialize() (eskapade.core_ops.links.ds_to_ds.DsToDs method), 103
 initialize() (eskapade.core_ops.links.DsObjectDeleter method), 110
 initialize() (eskapade.core_ops.links.DsToDs method), 110
 initialize() (eskapade.core_ops.links.event_looper.EventLooper method), 104
 initialize() (eskapade.core_ops.links.EventLooper method), 111
 initialize() (eskapade.core_ops.links.line_printer.LinePrinter method), 106
 initialize() (eskapade.core_ops.links.LinePrinter method), 112
 initialize() (eskapade.core_ops.links.repeat_chain.RepeatChain method), 112
 initialize() (eskapade.core_ops.links.RepeatChain method), 112
 initialize() (eskapade.core_ops.links.skip_chain_if_empty.SkipChainIfEmpty method), 108
 initialize() (eskapade.core_ops.links.SkipChainIfEmpty method), 113
 initialize() (eskapade.core_ops.links.to_ds_dict.ToDsDict method), 108
 initialize() (eskapade.core_ops.links.ToDsDict method), 114
 initialize() (eskapade.data_quality.links.fix_pandas_dataframe.FixPandasDataFrame method), 116
 initialize() (eskapade.data_quality.links.FixPandasDataFrame method), 118
 initialize() (eskapade.visualization.links.correlation_summary.CorrelationSummary method), 122
 initialize() (eskapade.visualization.links.CorrelationSummary method), 126
 initialize() (eskapade.visualization.links.df_boxplot.DfBoxplot method), 123
 initialize() (eskapade.visualization.links.df_summary.DfSummary method), 125
 initialize() (eskapade.visualization.links.DfBoxplot method), 127
 initialize() (eskapade.visualization.links.DfSummary method), 128
 io_base_dirs() (eskapade.core.process_services.ConfigObject method), 98
 io_conf() (eskapade.core.process_services.ConfigObject method), 98
 io_dir() (in module eskapade.core.persistence), 93
 io_path() (in module eskapade.core.persistence), 93
 IPythonEmbed (class in eskapade.core_ops.links), 111
 IPythonEmbed (class in eskapade.core_ops.links.ipython_embed), 105
 is_break_chain() (eskapade.core.definitions.StatusCode method), 83
 is_failure() (eskapade.core.definitions.StatusCode method), 83
 is_finished() (eskapade.analysis.links.read_to_df.ReadToDf method), 51
 is_finished() (eskapade.analysis.links.ReadToDf method), 62
 is_recoverable() (eskapade.core.definitions.StatusCode method), 84
 is_repeat_chain() (eskapade.core.definitions.StatusCode method), 84
 is_skip_chain() (eskapade.core.definitions.StatusCode method), 84

is_success() (eskapade.core.definitions.StatusCode method), 84
 is_undefined() (eskapade.core.definitions.StatusCode method), 84

K

key (eskapade.analysis.histogram.ValueCounts attribute), 75

L

latest_data_length() (eskapade.analysis.links.read_to_df.ReadToDf method), 51
 latest_data_length() (eskapade.analysis.links.ReadToDf method), 62
 lib() (in module eskapade.resources), 135
 LinePrinter (class in eskapade.core_ops.links), 111
 LinePrinter (class in eskapade.core_ops.links.line_printer), 106
 Link (class in eskapade.core.element), 86
 load() (eskapade.core.element.Link method), 87
 logger (eskapade.analysis.datetime.TimePeriod attribute), 66
 logger (eskapade.analysis.histogram.Histogram attribute), 72
 logger (eskapade.analysis.statistics.ArrayStats attribute), 80
 logger (eskapade.core.meta.Processor attribute), 89
 logger (eskapade.core.process_services.ProcessService attribute), 99

M

make_histogram() (eskapade.analysis.statistics.ArrayStats method), 81
 MissingPackageError, 132
 MissingPy4jError, 133
 MissingRooFitError, 133
 MissingRooStatsError, 133
 MissingRootError, 133
 MissingSparkError, 134

N

n_bins (eskapade.analysis.histogram.Histogram attribute), 73
 n_chains (eskapade.core.process_manager.ProcessManager attribute), 96
 n_dim (eskapade.analysis.histogram.Histogram attribute), 73
 n_links (eskapade.core.element.Chain attribute), 86
 name (eskapade.core.meta.Processor attribute), 90
 nononecounts (eskapade.analysis.histogram.ValueCounts attribute), 75

num_bins (eskapade.analysis.histogram.Histogram attribute), 73
 num_bins (eskapade.analysis.histogram.ValueCounts attribute), 75
 num_nonone_bins (eskapade.analysis.histogram.ValueCounts attribute), 75

O

obj_repr() (in module eskapade.helpers), 134
 offset (eskapade.analysis.datetime.UniformTsTimePeriod attribute), 67
 only_bool() (in module eskapade.analysis.histogram_filling), 78
 only_float() (in module eskapade.analysis.histogram_filling), 78
 only_int() (in module eskapade.analysis.histogram_filling), 78
 only_str() (in module eskapade.analysis.histogram_filling), 78

P

pandas_writer() (in module eskapade.analysis.links.write_from_df), 56
 pandasReader() (in module eskapade.analysis.links.read_to_df), 51
 parent (eskapade.core.meta.Processor attribute), 90
 parse_date_time() (eskapade.analysis.datetime.TimePeriod class method), 67
 parse_time_period() (eskapade.analysis.datetime.TimePeriod class method), 67
 period (eskapade.analysis.datetime.UniformTsTimePeriod attribute), 67
 period_index() (eskapade.analysis.datetime.FreqTimePeriod method), 66
 period_index() (eskapade.analysis.datetime.TimePeriod method), 67
 period_index() (eskapade.analysis.datetime.UniformTsTimePeriod method), 67
 persist (eskapade.core.process_services.ProcessServiceMeta attribute), 100
 persist_in_file() (eskapade.core.process_services.ProcessService method), 100
 persist_services() (eskapade.core.process_manager.ProcessManager method), 96
 plot_2d_histogram() (in module eskapade.visualization.vis_utils), 130
 plot_correlation_matrix() (in module eskapade.visualization.vis_utils), 130
 plot_histogram() (in module eskapade.visualization.vis_utils), 131

pop() (eskapade.core.meta.ProcessorSequence method), 90

Print() (eskapade.core.process_services.ConfigObject method), 98

Print() (eskapade.core.process_services.DataStore method), 99

print_chains() (eskapade.core.process_manager.ProcessManager method), 96

print_services() (eskapade.core.process_manager.ProcessManager method), 96

PrintDs (class in eskapade.core_ops.links), 112

PrintDs (class in eskapade.core_ops.links.print_ds), 106

process_1d_histogram() (eskapade.visualization.links.df_summary.DfSummary method), 125

process_1d_histogram() (eskapade.visualization.links.DfSummary method), 128

process_2d_histogram() (eskapade.visualization.links.df_summary.DfSummary method), 125

process_2d_histogram() (eskapade.visualization.links.DfSummary method), 128

process_and_store() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 78

process_and_store() (eskapade.analysis.links.value_counter.ValueCounter method), 55

process_and_store() (eskapade.analysis.links.ValueCounter method), 65

process_columns() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 78

process_columns() (eskapade.analysis.links.value_counter.ValueCounter method), 55

process_columns() (eskapade.analysis.links.ValueCounter method), 65

process_counts() (eskapade.analysis.histogram.ValueCounts method), 75

process_nan_histogram() (eskapade.visualization.links.df_summary.DfSummary method), 125

process_nan_histogram() (eskapade.visualization.links.DfSummary method), 129

process_sample() (eskapade.visualization.links.df_summary.DfSummary method), 125

process_sample() (eskapade.visualization.links.DfSummary method), 129

process_series() (eskapade.visualization.links.df_summary.DfSummary method), 125

process_series() (eskapade.visualization.links.DfSummary method), 129

process_transform_funcs() (in module eskapade.helpers), 134

ProcessManager (class in eskapade.core.process_manager), 94

Processor (class in eskapade.core.meta), 89

ProcessorSequence (class in eskapade.core.meta), 90

ProcessService (class in eskapade.core.process_services), 99

ProcessServiceMeta (class in eskapade.core.process_services), 100

R

RandomSampleSplitter (class in eskapade.analysis.links), 60

RandomSampleSplitter (class in eskapade.analysis.links.random_sample_splitter), 50

RandomSeeds (class in eskapade.core.definitions), 82

ReadToDf (class in eskapade.analysis.links), 61

ReadToDf (class in eskapade.analysis.links.read_to_df), 50

record_file_number() (in module eskapade.core.persistence), 93

record_vectorizer() (in module eskapade.analysis.links.record_vectorizer), 53

RecordFactorizer (class in eskapade.analysis.links), 62

RecordFactorizer (class in eskapade.analysis.links.record_factorizer), 52

RecordVectorizer (class in eskapade.analysis.links), 63

RecordVectorizer (class in eskapade.analysis.links.record_vectorizer), 53

Recoverable (eskapade.core.definitions.StatusCode attribute), 83

remove_all_services() (eskapade.core.process_manager.ProcessManager method), 96

remove_keys_of_inconsistent_type() (eskapade.analysis.histogram.Histogram method), 73

remove_keys_of_inconsistent_type() (eskapade.analysis.histogram.ValueCounts method), 75

remove_service() (eskapade.core.process_manager.ProcessManager method), 96

RepeatChain (class in eskapade.core_ops.links), 112

- RepeatChain (class in eskapade.core_ops.links.repeat_chain), 107
- RepeatChain (eskapade.core.definitions.StatusCode attribute), 83
- repl_whites() (in module eskapade.core.persistence), 93
- reset() (eskapade.core.process_manager.ProcessManager method), 96
- reset_config() (eskapade.core.mixin.ConfigMixin method), 92
- reset_escaped() (in module eskapade.core.execution), 88
- run() (eskapade.core.process_manager.ProcessManager method), 96
- ## S
- service() (eskapade.core.process_manager.ProcessManager method), 96
- set_begin_end_chain_opt() (in module eskapade.core.definitions), 84
- set_chunk_size() (eskapade.analysis.links.read_to_df.ReadToDf method), 51
- set_chunk_size() (eskapade.analysis.links.ReadToDf method), 62
- set_custom_user_vars() (in module eskapade.core.definitions), 84
- set_log_level_opt() (in module eskapade.core.definitions), 84
- set_matplotlib_backend() (in module eskapade.utils), 136
- set_opt_var() (in module eskapade.core.definitions), 84
- set_seeds() (in module eskapade.core.definitions), 84
- set_single_chain_opt() (in module eskapade.core.definitions), 84
- set_user_opts() (eskapade.core.process_services.ConfigObject method), 98
- simulate() (eskapade.analysis.histogram.Histogram method), 73
- Singleton (class in eskapade.core.meta), 90
- skey (eskapade.analysis.histogram.ValueCounts attribute), 75
- SkipChain (eskapade.core.definitions.StatusCode attribute), 83
- SkipChainIfEmpty (class in eskapade.core_ops.links), 112
- SkipChainIfEmpty (class in eskapade.core_ops.links.skip_chain_if_empty), 107
- start_timer() (eskapade.core.mixin.TimerMixin method), 92
- StatusCode (class in eskapade.core.definitions), 83
- stop_timer() (eskapade.core.mixin.TimerMixin method), 92
- store() (eskapade.core.element.Link method), 87
- Success (eskapade.core.definitions.StatusCode attribute), 83
- sum_counts (eskapade.analysis.histogram.ValueCounts attribute), 76
- sum_data_length() (eskapade.analysis.links.read_to_df.ReadToDf method), 51
- sum_data_length() (eskapade.analysis.links.ReadToDf method), 62
- sum_nonone_counts (eskapade.analysis.histogram.ValueCounts attribute), 76
- summary() (eskapade.core.element.Link method), 87
- summary() (eskapade.core.process_manager.ProcessManager method), 96
- surface() (eskapade.analysis.histogram.Histogram method), 73
- ## T
- template() (in module eskapade.resources), 135
- TimePeriod (class in eskapade.analysis.datetime), 66
- TimerMixin (class in eskapade.core.mixin), 92
- to_date_time() (in module eskapade.data_quality.dq_helper), 119
- to_float() (in module eskapade.data_quality.dq_helper), 119
- to_int() (in module eskapade.data_quality.dq_helper), 119
- to_normalized() (eskapade.analysis.histogram.Histogram method), 73
- to_ns() (in module eskapade.analysis.histogram_filling), 79
- to_str() (in module eskapade.analysis.histogram_filling), 79
- to_str() (in module eskapade.data_quality.dq_helper), 119
- ToDsDict (class in eskapade.core_ops.links), 113
- ToDsDict (class in eskapade.core_ops.links.to_ds_dict), 108
- total_time() (eskapade.core.mixin.TimerMixin method), 92
- truncated_bin_edges() (eskapade.analysis.histogram.BinningUtil method), 69
- tutorial() (in module eskapade.resources), 136
- ## U
- Undefined (eskapade.core.definitions.StatusCode attribute), 83
- UniformTsTimePeriod (class in eskapade.analysis.datetime), 67
- UnknownSetting, 88
- ## V
- value_to_bin_center() (in module eskapade.analysis.histogram_filling), 79

`value_to_bin_index()` (in module `eskapade.analysis.histogram_filling`), 79

`value_to_bin_label()` (`eskapade.analysis.histogram.BinningUtil` method), 69

`ValueCounter` (class in `eskapade.analysis.links`), 63

`ValueCounter` (class in `eskapade.analysis.links.value_counter`), 54

`ValueCounts` (class in `eskapade.analysis.histogram`), 74

`variable` (`eskapade.analysis.histogram.Histogram` attribute), 74

W

`weighted_quantile()` (in module `eskapade.analysis.statistics`), 82

`WriteFromDf` (class in `eskapade.analysis.links`), 65

`WriteFromDf` (class in `eskapade.analysis.links.write_from_df`), 56