
Eskapade Documentation

Release 0.6

KPMG Big Data Team

Sep 06, 2017

Contents

1	Eskapade	3
2	Release notes	5
2.1	Version 0.6	5
2.2	Version 0.5	5
2.3	Version 0.4	6
3	Contents	7
3.1	Introduction	7
3.2	Installation	8
3.3	Tutorials	12
3.4	Package structure	32
3.5	Command Line Arguments	35
3.6	Contributing	40
3.7	API	40
3.8	Appendices	133
3.9	Indices and tables	142
	Python Module Index	143

Version: 0.6

Date: Sep 06, 2017

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

Eskapade is a light-weight, python-based data analysis framework, meant for 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. The machine learning packages include:

- scikit-learn,
- Spark MLlib,
- ROOT.

Likewise, Eskapade uses a manner of different data structures to handle the data, such as:

- pandas DataFrames,
- numpy arrays,
- Spark DataFrames

and more.

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 had 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.

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!

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. Stay tuned!

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.

What is Eskapade?

Eskapade is an abbreviation for: 'Enterprise Solution KPMG Advanced Predictive Analytics Decision Engine'. It is a framework to 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 Machine Learning models that are retrained when a certain trigger is reached.

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](#).

Naming convention

Before settling on the name *Eskapade*, this project had internal naming conventions including *Decision Engine* and *Analytics Engine*. We are working on removing all old names and streamlining this all into future versions, but one might still find these in certain parts of the repository.

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

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

Eskapade with Docker

There is a Docker image available to try out Eskapade without concerns about the installation process of all necessary packages. 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.

Installing Eskapade

Make sure you have installed the latest version of the public Eskapade repository, with the command:

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

Installing Docker

Docker is a great tool to develop and automate the deployment of applications inside software containers. To install Docker, go docker.com and follow the installation instructions.

Pulling in the Eskapade environment

To pull in the environment for running Eskapade (but excluding the Eskapade repository itself), type in a shell:

```
docker pull kave/eskapade-env:0.6
```

Downloading this docker image can take a minute or two.

Running Eskapade with Docker

To start up the Eskapade docker environment, with access to the Eskapade repository, do:

```
docker run -it --name es-ktb -p 8888:8888 -v /path/to/your/local/eskapade:/opt/
↳eskapade kave/eskapade-env:0.6 bash
```

This command will start up a bash shell in the docker `kave/eskapade-env:0.6` image, and opens port 8888. The Eskapade setup file will be automatically sourced.

The option `-v /path/to/your/local/eskapade:/opt/eskapade` mounts the local directory `/path/to/your/local/eskapade` (containing your local Eskapade repository) under `/opt/eskapade` in the docker container. You can now edit the files in this directory, either locally or in the (docker) bash shell, and any updates to these files will be kept after exiting docker.

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

E.g. one can now do:

```
cd /opt/eskapade
```

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

Exit the (docker) bash shell with:

```
exit
```

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

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

```
alias eskapade_docker='docker run -it --name es-ktb -p 8888:8888 -v /path/to/local/
↳eskapade:/opt/eskapade kave/eskapade-env:0.6 bash'
```

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
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 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

Everytime 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.

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`.

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`.

Eskapade on your own machine

The repository is hosted on github, clone it to your machine with:

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

Requirements

Eskapade requires Anaconda, which can be found [here](#). Eskapade was tested with version 4.3.

It also uses some non-standard libraries, which can be found in *requirements.txt* in the repository. These can be installed by doing:

```
$ pip install -r requirements.txt
```

You are now ready to use Eskapade!

After installation

To get started, source Eskapade in the root of the repository:

```
$ source setup.sh
```

You can now call the path of Eskapade with:

```
$ echo $ESKAPADE
```

or in python with

```
import os
os.environ['ESKAPADE']
```

See the readme files in other parts of the repository for specific usage.

Installing Eskapade on MacOS

To install eskapade on MacOS, see our MacOS appendix.

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.

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:

```
$ source setup.sh
$ run_eskapade.py ./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):

```
2017-02-27 20:32:19,826 INFO [hello_world/execute]: Hello World
2017-02-27 20:32:19,828 INFO [hello_world/execute]: 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:

```
link = core_ops.HelloWorld(name='HelloWorld')
link.set_log_level(logging.DEBUG)
link.repeat = settings['n_repeat']
ch.add_link(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 $ESKAPADE/python/eskapade/core_ops/links/hello_world.py
```

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

```
self.log().info('Hello {0}'.format(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_link(hello2)
```

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.

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 -P $ESKAPADE/data/ https://s3-eu-west-1.amazonaws.com/kpmg-eskapade-share/data/
↳LAOzone.data
```

To run the macro type on your CLI:

```
$ run_eskapade.py 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:

```
* * * Welcome to Eskapade * * *
...
2017-02-10 15:24:35,968 INFO [processManager/Print]: Number of chains:    2
...
* * * 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 `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. There are two chains added to the macro, with following line you can add a chain:

```
proc_mgr.add_chain('Data')
```

This chain called `Data` is added to the `ProcessManager`, which is the object that runs the entire macro. Then the chain is fetched by:

```
proc_mgr.get_chain('Data')
```

and a `Link` is added. First the link is initialized (links are classes) and its properties are set, and finally it is inserted into the chain:

```
reader = analysis.ReadToDf(name='Read_LA_ozone', path=DATA_FILE_PATH, reader=pd.read_
↳csv, key='data')
proc_mgr.get_chain('Data').add_link(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:

```
$ run_eskapade.py 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:

```
$ make_link.sh python/eskapade/analysis/links YourLink
```

The script 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 script!

The script 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.

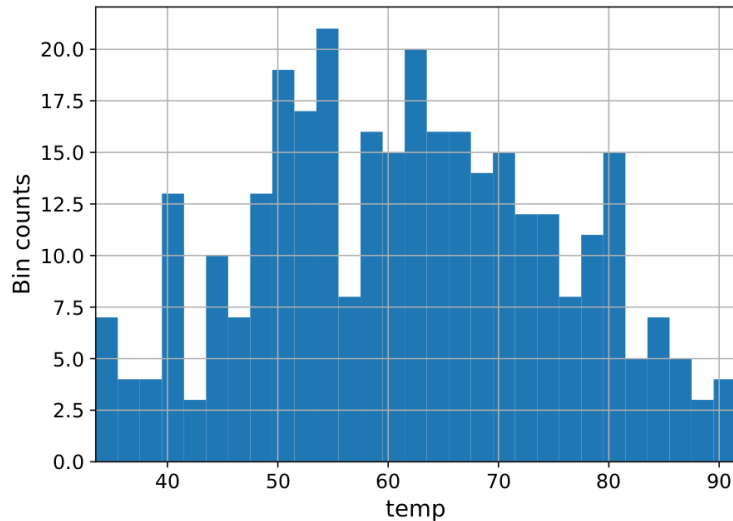
Now run the entire macro with the new code and compile the output `.tex` file. This can be done on the command line with

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

If you have pdflatex installed on your machine. 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!

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

```
$ make_macro.sh 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 = proc_mgr.add_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 sub-package.
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.

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. Do not forget to set up the environment before starting the notebook (and in case you use a virtual environment activate it):

```
$ cd $ESKAPADE
$ source setup.sh
$ cd some-working-dir
$ jupyter notebook
```

An Eskapade notebook

To run Eskapade use the `make_notebook.sh` script in `scripts/` to create a template notebook. For example:

```
$ make_notebook.sh ./ TestRun
```

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

```
import imp
import logging
imp.reload(logging)
log = logging.getLogger()
log.setLevel(logging.DEBUG) # Set the LogLevel here

from eskapade.core import execution
from eskapade import ConfigObject, DataStore, ProcessManager

# --- basic config
settings = ProcessManager().service(ConfigObject)
settings['macro'] = os.environ['ESKAPADE'] + '/tutorials/tutorial_1.py'
settings['analysisName'] = 'Tutorial_1'
settings['version'] = 0
settings['logLevel'] = logging.DEBUG # and set the LogLevel here

# --- 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.run_eskapade(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.run_eskapade(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 -u
↳storeResultsEachChain
(command line: ``-w``) on.
```

```
import imp
import logging
imp.reload(logging)
log = logging.getLogger()
log.setLevel(logging.DEBUG)

from eskapade import DataStore, ConfigObject, ProcessManager

# --- example inspecting the data store after the preprocessing chain
ds = DataStore.import_from_file(os.environ['ESKAPADE']+'/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
co = ConfigObject.import_from_file(os.environ['ESKAPADE']+'/results/Tutorial_1/proc_
↳service_data/v0/_Summary/eskapade.core.process_services.ConfigObject.pkl')
co.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:

```
$ make_notebook.sh tutorials/ tutorial_3_notebook
```

This will create a notebook in `tutorials/` with the name `tutorial_3_notebook` running macro `tutorial_1.py`. 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. Use:

```
!pwd
```

In Jupyter to find which path you are working on, and change the load path in the macro to the proper one. This can be for example:

```
os.environ['ESKAPADE'] + '/data/LAozone.data'
```

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:

```
RuntimeError: tried to add chain with existing name to process manager
```

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 make a singleton, but it already exists. 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.run_eskapade(settings)
```

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

```
ds = ProcessManager().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:

```
ESKAPADE/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(os.environ['ESKAPADE']+'/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` and `DataStore` objects if they are available. By importing and calling:

```
from eskapade import ConfigObject
settings = ConfigObject.import_from_file(os.environ['ESKAPADE']+'/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.

Tutorial 4: 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/rfl104_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 $ESKAPADE/cxx/roofit/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: `$ESKAPADE/cxx/roofit/src/RooWeibull.cxx`.

Now move the header files to their correct location:

```
$ mv MyPdfV*.h $ESKAPADE/cxx/roofit/include/
```

To make sure that these classes get picked up in Eskapade roofit library, open the file `$ESKAPADE/cxx/roofit/dict/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:

```
$ cd $ESKAPADE
$ make install
```

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 `$ESKAPADE/tutorials/tutorial_4.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.

```
$ cd $ESKAPADE
$ source setup.sh
$ run_eskapade.py tutorials/tutorial_4.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'
log.info('Now checking existence of ROOT class %s' % pdf_name)
cl = ROOT.TClass.GetClass(pdf_name)
if not cl:
    log.critical('Could not find ROOT class %s. Did you build and compile it_
↪correctly?' % pdf_name)
    sys.exit(1)
```

```
else:
    log.info('Successfully found ROOT class %s' % 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 `RooWorkspace 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 [here](#), [here](#) and [here](#) for more details. The entire `rooworkspace` factory syntax can be found [here](#).

```
ch = proc_mgr.add_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_link(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_link(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_link(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_link(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 $ESKAPADE/results/tutorial_4/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 `$ESKAPADE/tutorials` directory, e.g. see `esk401_roothist_fill_plot_convert.py` and all other 400 numbered macros.

Tutorial 5: 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 `tutorials/tutorial_5.py`.

Running the tutorial macro

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

```
$ source setup.sh
$ run_eskapade.py tutorials/tutorial_5.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:

```
proc_mgr.service(SparkManager).spark_session
```

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

```
Adding Python modules to ZIP archive /Users/gossie/git/gitlab-nl/decision-engine/
↳eskapade/es_python_modules.zip
```

This is the `SparkManager` that ensures all Eskapade source code is uploaded and available to the Spark cluster when running in a distributed environment.

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:

```
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)
proc_mgr.get_chain('Data').add_link(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:

```
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')
proc_mgr.get_chain('Data').add_link(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')
proc_mgr.get_chain('Data').add_link(transform)
```

Note that the functions defined in the macro are converted to user-defined functions with `pyspark.sql.function.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.

Histogramming

As was demonstrated in Tutorial 1, the `DfSummary` link creates LaTeX/PDF reports with histograms. Those histograms are obtained directly from a Pandas dataframe or from a dictionary of `Histogrammar` histograms. This link can be re-used for Tutorial 4. However, an additional step is needed: create histograms of Spark dataframe columns with `Histogrammar`. This step can be carried out with the `SparkHistogrammarFiller` link. The code snippet for generating a report of Spark dataframe histograms then looks like:

```
histo = spark_analysis.SparkHistogrammarFiller(name='Histogrammer', read_
↳key=transform.store_key, store_key='hist')
histo.columns = ['vis', 'vis_km', 'doy', 'date']
proc_mgr.get_chain('Summary').add_link(histo)

summarizer = visualization.DfSummary(name='Create_stats_overview', read_key=histo.
↳store_key, var_labels=VAR_LABELS, var_units=VAR_UNITS)
proc_mgr.get_chain('Summary').add_link(summarizer)
```

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.:

```
$ make_link.sh python/eskapade/spark_analysis SparkDfPrinter
```

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

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 SparkDfPrint"""

    proc_mgr = ProcessManager()
    settings = proc_mgr.service(ConfigObject)
    ds = proc_mgr.service(DataStore)

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

    return StatusCode.Success
```

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_5.py` macro. This should look similar to:

```
from eskapade.spark_analysis import SparkDfPrint
...

printer = SparkDfPrint(name='Print_spark_df', read_key=transform.store_key, nrows=42)
proc_mgr.get_chain('Summary').add_link(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:

```
$ run_eskapade.py tutorials/tutorial_5.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 examples

Example Eskapade macros using Spark can be found in the `tutorials` directory, see `esk601_spark_configuration.py` and further.

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/latest/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 $ for ((i=0; i<=100; i++)); do echo "Hello world" > /tmp/dummy_$(printf
↪%05d ${i}); sleep 0.1; done
terminal 2 $ run_eskapade -c stream_type='tcp' $ESKAPADE/tutorials/esk610_spark_
↪streaming.py
```

Where bash `for`-loop will create a new file containing `Hello world` in the `/tmp` directory every 0.1 second. Spark Streaming will pick up and process these files and in `terminal 2` a word count of the processed data will be displayed. Output is stored in `$ESKAPADE/results/esk610_spark_streaming/data/v0/dstream/wordcount`.

TCP stream

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

```
terminal 1 $ nc -lk 9999
terminal 2 $ run_eskapade -c stream_type='tcp' $ESKAPADE/tutorials/esk610_spark_
↪streaming.py
```

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

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-docs. If that does not help, read and try to understand the example macros in `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 `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.

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 (currently 3.5.2 of Anaconda, use the interpreter of your conda environment)
 - **Under ‘Build, Execution & Deployment’ / ‘Console’ / ‘Python Console’:**
 - * The correct Python version (currently 3.5.2 of Anaconda, use the interpreter of your conda environment)
- **Run/Debug Configuration:**
 - Under ‘Python’ add new configuration
 - Script: `scripts/run_escapede.py`
 - Script parameters: `-w ../tutorials/tutorial_1.py`
 - Working directory: `$ESKAPADE`

- Python interpreter: check if it is the correct Python version (currently 3.5.2 of Anaconda, corresponding to your conda environment)
- Environment variables: should contain those defined in `$ESKAPADE/setup.sh`.

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` 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 `make_link` script.
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 = ProcessManager().service(ConfigObject)
ds = ProcessManager().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 imp
imp.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.

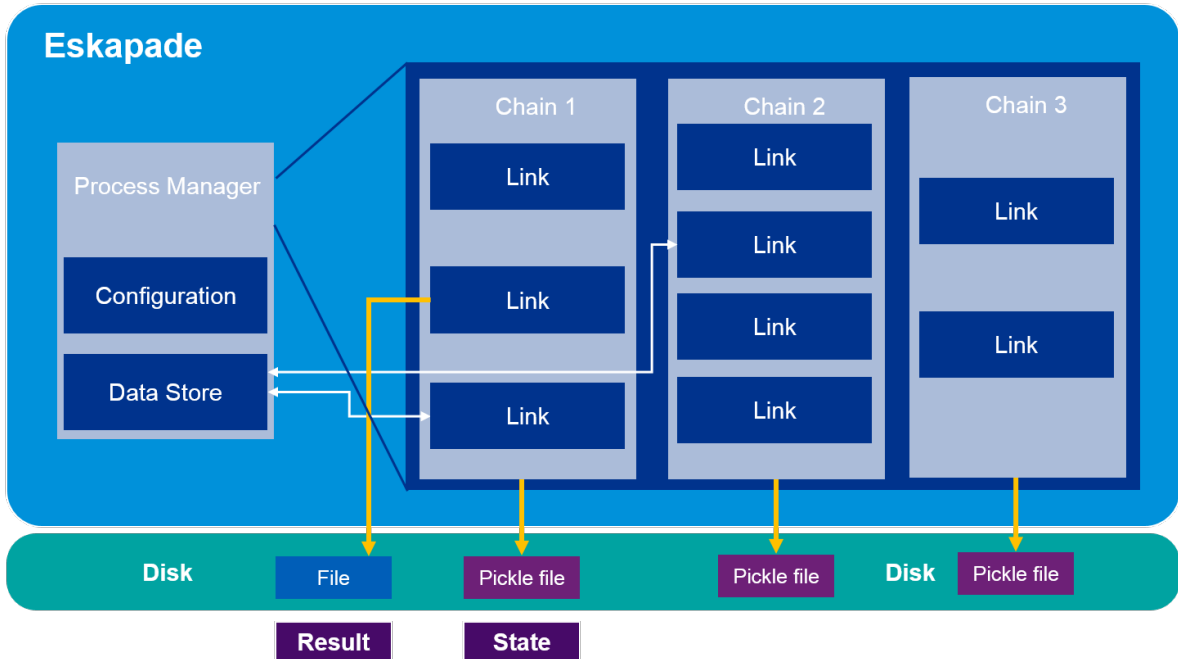
Package structure

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.

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:

- `run_eskapade.py` 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.

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, and the corresponding tests.

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
    |-tests
      |-integration
```

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`.

Unit tests are defined in modules in the `tests` directory. Ideally, there is a test module for each (link) module in the Eskapade subpackage. Optionally, integration tests are implemented in `tests/integration`. For the `eskapade.analysis` package, there is the module `test_tutorial_macros` with integration tests that run the tutorial macros corresponding to this package.

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.

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 ConfigObject, ProcessManager
```

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.

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/esk304_df_boxplot/data/v0/report`.

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.

Troubleshooting

One of the easiest mistakes to make when running the framework is not sourcing the right files or opening a new terminal without setting the right environment. Be careful of this, if you want to run eskapade you have to:

- Source the right virtual environment
- Source the eskapade repository
- Start your notebook / start your IDE / run the code

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

Command Line Arguments

Overview

We start this section with a short overview of a few often used arguments of the Eskapade run script `scripts/run_eskapade.py`. 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 script.

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/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:

```
$ run_escapede.py --store-all 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:

```
$ run_escapede.py --single-chain=Data 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:

```
$ run_escapede.py --begin-with=Summary tutorials/tutorial_1.py
```

Stop running after a specified chain:

```
$ run_escapede.py --end-with=Data tutorials/tutorial_1.py
```

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

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
-seed		SEED	set seed for random-number generation

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:

```
$ run_escapede.py -L DEBUG tutorials/tutorial_1.py
```

The available log levels are:

```
DEBUG,
INFO,
WARNING,
ERROR,
FATAL,
OFF
```

They correspond to the appropriate POSIX levels.

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

```
self.log().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:

```
$ run_eskapade.py --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`:

```
$ run_eskapade.py -i 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/data/VERSION/`, containing the state of Eskapade at the end of the chain:

```
$ run_eskapade.py --store-all tutorials/tutorial_1.py
```

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

```
$ run_eskapade.py --store-one=Data tutorials/tutorial_1.py
```

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

```
$ run_eskapade.py --store-none 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:

```
$ run_eskapade.py -s Summary 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.

```
$ run_eskapade.py -b Summary 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`:

```
$ run_escapede.py -e Data 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.

```
$ run_escapede.py -v 1 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:

```
$ run_escapede.py --unpickle-config results/Tutorial_1/proc_service_data/v0/latest/  
↳escapede.core.process_services.ConfigObject.pkl
```

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

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

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

Profiling your code

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

```
$ run_escapede.py --profile=cumulative 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 list of sorting options for the profiling, run:

```
$ run_eskapade.py --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:

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

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.

Testing

When contributing the original tests should succeed. To run the tests make sure you fulfill the requirements and source Eskapade in the root of the repository. Then you can run the tests:

```
$ source setup.sh  
$ run_tests.py unit  
$ run_tests.py integration
```

Contact the owners of this repository when you need help with adding tests.

API

API Documentation

Eskapade

eskapade package

Subpackages

eskapade.analysis package

Subpackages

eskapade.analysis.links package

Submodules

eskapade.analysis.links.apply_func_to_df module

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

Bases: `eskapade.core.run_elements.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)

addApplyFunc (*func, outColumn, inColumn=’, *args, **kwargs*)

Add function to be applied to dataframe

execute ()

Execute link

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

Apply groupby to dataframe

initialize ()

Initialize link

eskapade.analysis.links.apply_selection_to_df module

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

Bases: `eskapade.core.run_elements.Link`

Applies queries with sub-selections to a pandas.DataFrame

`__init__` (**kwargs)

Applies queries with sub-selections to a pandas.DataFrame

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

Parameters

- **name** (*str*) – name of link
- **readKey** (*str*) – key of data to read from data store
- **storeKey** (*str*) – key of data to store in data store. If not set readKey is overwritten.
- **querySet** (*list*) – list of strings, query expressions to evaluate in the same order, see pandas documentation

- **selectColumns** (*list*) – column names to select after querying
- **continueIfFailure** (*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 ApplySelectionToDf

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 ApplySelectionToDf

Perform checks on provided attributes

eskapade.analysis.links.assign_random_class module

class `eskapade.analysis.links.assign_random_class.AssignRandomClass` (**kwargs)

Bases: `eskapade.core.run_elements.Link`

AssignRandomClass randomly chooses a number of records that are to be added to the top X records that are returned to a client. In this way a top X also has a certain set of randomly chosen records. E.g. these records make any overtraining less likely to happen.

__init__ (**kwargs)

Store the configuration of link AssignRandomClass

Parameters

- **name** (*str*) – name of link
- **readKey** (*str*) – key of data to read from data store
- **column** (*str*) – name of new column that specifies the randomly assigned class. Default is randomclass.
- **nclasses** (*int*) – number of random classes. Needs to be set.
- **fractions** (*list*) – list of fractions of random records assigned to n classes. Needs to be set. Can be one less than n classes.
- **nevents** (*list*) – list of number of random records assigned to n classes. Can be one less than n classes. (optional instead of ‘fractions’)

execute ()

Execute AssignRandomClass

initialize ()

Check and initialize attributes of AssignRandomClass

eskapade.analysis.links.basic_generator module

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

Bases: `eskapade.core.run_elements.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 BasicGenerator

`initialize ()`

Initialize BasicGenerator

eskapade.analysis.links.df_concatenator module

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

Bases: `eskapade.core.run_elements.Link`

Concatenates multiple pandas datadrames.

`__init__ (**kwargs)`

Store the configuration of link DfConcatenator

Parameters

- **name** (*str*) – name of link
- **storeKey** (*str*) – key of data to store in data store
- **readKeys** (*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 DfConcatenator

Perform concatenation of multiple pandas datadrames.

`initialize ()`

Initialize DfConcatenator

eskapade.analysis.links.df_merger module

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

Bases: `eskapade.core.run_elements.Link`

Merges two pandas DataFrames.

`__init__ (**kwargs)`

Store the configuration of link DfMerger

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**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 HistogrammarFiller 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 pars certain columns

Example quantity dictionary is:

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

Parameters

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

Example drop_keys dictionary is:

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

construct_empty_hist (*columns*)

Create an (empty) histogram of right type

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

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

Returns created histogram

Return type histogrammar.Count

fill_histogram (*idf, columns*)

Fill input histogram with column(s) of input dataframe

Parameters

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

eskapade.analysis.links.random_sample_splitter module

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

Bases: `eskapade.core.run_elements.Link`

RandomSampleSplitter splits an input dataframe into a number of sub data-frames.

Records are assigned randomly.

__init__ (***kwargs*)

Store the configuration of link RandomSampleSplitter

Parameters

- **name** (*str*) – name of link
- **readKey** (*str*) – key of data to read from data store
- **storeKey** (*list*) – keys of datasets to store in data store. Number of sub samples equals length of storeKey list.
- **fractions** (*list*) – list of fractions (0<fraction<1) of records assigned to the sub samples. Sum can be less than 1. Needs to be set.
- **nevents** (*list*) – list of number of random records assigned to the sub samples. (optional instead of ‘fractions’)

execute ()

Execute RandomSampleSplitter

initialize ()

Check and initialize attributes of RandomSampleSplitter

eskapade.analysis.links.read_to_df module**class** `eskapade.analysis.links.read_to_df.ReadToDf` (**kwargs)Bases: `eskapade.core.run_elements.Link`

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)

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 ReadToDf

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

initialize ()

Initialize ReadToDf

isFinished ()

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

setChunkSize (*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

class `eskapade.analysis.links.record_factorizer.RecordFactorizer` (***kwargs*)
Bases: `eskapade.core.run_elements.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 RecordFactorizer 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 catergy 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 RecordFactorizer

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

initialize()

Initialize RecordFactorizer

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

eskapade.analysis.links.record_vectorizer module

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

Bases: `eskapade.core.run_elements.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 RecordVectorizer 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 RecordVectorizer

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

initialize()

Initialize RecordVectorizer

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

`eskapade.analysis.links.record_vectorizer.record_vectorizer` (*df*, *column_to_vectorize*,
column_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

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 ValueCounter 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*) – cols to pick up from input data (default is all columns)
- **bin_specs** (*dict*) – dictionaries used for rebinning numeric or timestamp columns

Example `bin_specs` dictionary is:

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

Parameters

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

Example `drop_keys` dictionary is:

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

drop_inconsistent_keys (*columns*, *obj*)

Drop inconsistent keys

Drop inconsistent keys from a ValueCounts or Histogram object.

Parameters

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

fill_histogram (*idf*, *columns*)

Fill input histogram with column(s) of input dataframe

Parameters

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

initialize ()

Initialize ValueCounter

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

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

Bases: `eskapade.core.run_elements.Link`

Write a DataFrame from the DataStore to disk.

__init__ (**kwargs)

Store the configuration of link WriteFromDf

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 WriteFromDf

Pick up the dataframe and write to disk.

initialize ()

Initialize WriteFromDf

`eskapade.analysis.links.write_from_df.pandasWriter` (*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.run_elements.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)

addApplyFunc (*func*, *outColumn*, *inColumn*=’, *args, **kwargs)

Add function to be applied to dataframe

execute ()

Execute link

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

Apply groupby to dataframe

initialize ()

Initialize link

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

Bases: `eskapade.core.run_elements.Link`

Applies queries with sub-selections to a pandas.DataFrame

`__init__` (**kwargs)

Applies queries with sub-selections to a pandas.DataFrame

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

Parameters

- **name** (*str*) – name of link
- **readKey** (*str*) – key of data to read from data store
- **storeKey** (*str*) – key of data to store in data store. If not set readKey is overwritten.
- **querySet** (*list*) – list of strings, query expressions to evaluate in the same order, see pandas documentation
- **selectColumns** (*list*) – column names to select after querying
- **continueIfFailure** (*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 ApplySelectionToDf

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 ApplySelectionToDf

Perform checks on provided attributes

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

Bases: `eskapade.core.run_elements.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 BasicGenerator

`initialize` ()

Initialize BasicGenerator

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

Bases: `eskapade.core.run_elements.Link`

Concatenates multiple pandas datadrames.

`__init__ (**kwargs)`

Store the configuration of link DfConcatenator

Parameters

- **name** (*str*) – name of link
- **storeKey** (*str*) – key of data to store in data store
- **readKeys** (*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 DfConcatenator

Perform concatenation of multiple pandas datadrames.

`initialize ()`

Initialize DfConcatenator

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

Bases: `eskapade.core.run_elements.Link`

Merges two pandas DataFrames.

`__init__ (**kwargs)`

Store the configuration of link DfMerger

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.ReadToDf (**kwargs)`

Bases: `eskapade.core.run_elements.Link`

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)

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. chunksize requires `pd.read_csv` or `pd.read_table`.
- **kwargs** – all other key word arguments are passed on to the pandas reader.

`execute` ()

Execute ReadToDf

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

`initialize` ()

Initialize ReadToDf

`isFinished` ()

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

`setChunkSize` (*size*)

`sum_data_length` ()

Return sum length of all datasets processed sofar

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

Bases: `eskapade.core.run_elements.Link`

Vectorize data-frame columns

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

`__init__` (**kwargs)

Initialize RecordVectorizer 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 RecordVectorizer

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

initialize ()

Initialize RecordVectorizer

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

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

Bases: `eskapade.core.run_elements.Link`

Write a DataFrame from the DataStore to disk.

__init__ (**kwargs)

Store the configuration of link WriteFromDf

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 WriteFromDf

Pick up the dataframe and write to disk.

initialize ()

Initialize WriteFromDf

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

Bases: `eskapade.core.run_elements.Link`

AssignRandomClass randomly chooses a number of records that are to be added to the top X records that are returned to a client. In this way a top X also has a certain set of randomly chosen records. E.g. these records make any overtraining less likely to happen.

__init__ (**kwargs)

Store the configuration of link AssignRandomClass

Parameters

- **name** (*str*) – name of link
- **readKey** (*str*) – key of data to read from data store
- **column** (*str*) – name of new column that specifies the randomly assigned class. Default is randomclass.
- **nclasses** (*int*) – number of random classes. Needs to be set.
- **fractions** (*list*) – list of fractions of random records assigned to n classes. Needs to be set. Can be one less than n classes.
- **nevents** (*list*) – list of number of random records assigned to n classes. Can be one less than n classes. (optional instead of ‘fractions’)

execute ()

Execute AssignRandomClass

initialize ()

Check and initialize attributes of AssignRandomClass

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

Bases: `eskapade.core.run_elements.Link`

RandomSampleSplitter splits an input dataframe into a number of sub data-frames.

Records are assigned randomly.

__init__ (**kwargs)

Store the configuration of link RandomSampleSplitter

Parameters

- **name** (*str*) – name of link
- **readKey** (*str*) – key of data to read from data store
- **storeKey** (*list*) – keys of datasets to store in data store. Number of sub samples equals length of storeKey list.
- **fractions** (*list*) – list of fractions (0<fraction<1) of records assigned to the sub samples. Sum can be less than 1. Needs to be set.
- **nevents** (*list*) – list of number of random records assigned to the sub samples. (optional instead of ‘fractions’)

execute ()

Execute RandomSampleSplitter

initialize ()

Check and initialize attributes of RandomSampleSplitter

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

Bases: `eskapade.core.run_elements.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 RecordFactorizer 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 RecordFactorizer

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

initialize ()

Initialize RecordFactorizer

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

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 ValueCounter 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)

initialize ()

Initialize ValueCounter

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.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 HistogrammarFiller instance

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

Parameters

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

Example `bin_specs` dictionary is:

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

Parameters

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

Example `quantity` dictionary is:

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

Parameters

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

Example `drop_keys` dictionary is:

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

construct_empty_hist (*columns*)

Create an (empty) histogram of right type

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

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

Returns created histogram

Return type histogrammar.Count

fill_histogram (*idf, columns*)

Fill input histogram with column(s) of input dataframe

Parameters

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

Submodules

eskapade.analysis.datetime module

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

Bases: `eskapade.analysis.datetime.TimePeriod`

__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.mixins.ArgumentsMixin`, `eskapade.mixins.LoggingMixin`

__init__ (**kwargs)

initialize TimePeriod instance

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`

`__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

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 BinningUtil 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_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

- **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=[]*)

Bin edges corresponding to a given variable range

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

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.mixins.ArgumentsMixin`, `eskapade.mixins.LoggingMixin`

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_specs = { '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 ()

Retrun 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=[]*, *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.
- **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

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 (*prefered_key_type=None*)

Remove all keys that have inconsistent data type(s)

Parameters **prefered_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={}, sel={}*)

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={}, sel={}*)

Initialize ValueCounts 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.
- **counts** (*dict*) – 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.
- **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

Value-counts dictionary

Returns after processing, returns the `value_counts` dictionary

Return type `dict`

create_sub_counts (*subkey, sel={}*)

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. Default is `{}`.

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

Returns all key-values of a subset of keys.

Return type `tuple`

key

Current value-counts key

Returns the key

Return type `tuple`

nononecounts

Value-counts dictionary without `None` keys

Returns after processing, returns the `value_counts` dictionary without `None` keys

Return type `dict`

num_bins

Number of value-counts bins

Returns number of bins

Return type int

num_nonone_bins

Number of not-none value-counts bins

Returns number of not-none bins

Return type int

process_counts (*accept_equiv=True*)

Project value counts onto the existing subset of keys

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

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

Returns successful projection or not

Return type bool

remove_keys_of_inconsistent_type (*prefered_key_type=None*)

Remove keys with inconsistent data type(s)

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

skey

Current value-counts subkey

Returns the subkey

Return type tuple

sum_counts

Sum of counts of all value-counts bins

Returns the sum of counts of all bins

Return type float

sum_nonone_counts

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

Returns the sum of not-none counts of all bins

Return type float

eskapade.analysis.histogram_filling module

class `eskapade.analysis.histogram_filling.HistogramFillerBase` (**kwargs)

Bases: `eskapade.core.run_elements.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 HistogramFillerBase 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 HistogramFillerBase

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 HistogramFillerBase

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 HistogramFillerBase

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

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

Bases: `eskapade.mixins.LoggingMixin`

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` ()

`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 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
- **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

Subpackages

Submodules

eskapade.core.definitions module**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.

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

__init__ (**kwargs)Initialize `RandomSeeds` 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.Enum`

Return-status codes for Eskapade run

A `StatusCode` object is returned after each `initialize`, `execute`, and `finalize` function call of links, chains, and the process manager.

Failure = 4**Recoverable = 2****RepeatChain = 5****SkipChain = 3****Success = 1****Undefined = 6****isFailure** ()

Check if status is “Failure”

Return type `bool`**isRecoverable** ()

Check if status is “Recoverable”

Return type `bool`**isRepeatChain** ()

Check if status is “RepeatChain”

Return type `bool`**isSkipChain** ()

Check if status is “SkipChain”

Return type `bool`**isSuccess** ()

Check if status is “Success”

Return type bool

isUndefined()

Check if status is “Undefined”

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.execution module

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

Reset Eskapade objects

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

`eskapade.core.execution.run_eskapade(settings=None)`

Run Eskapade

This function is called in the script `scripts/run_eskapade.py` 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.persistence module

class `eskapade.core.persistence.IoConfig(**input_config)`

Bases: dict

Configuration object for I/O operations

__init__ (***input_config*)

Initialize IoConfig instance

`eskapade.core.persistence.create_dir(dir_path)`

Function to create directory

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

`eskapade.core.persistence.io_dir` (*io_type, io_conf*)

Functions to construct I/O paths

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, io_conf, sub_path*)

Functions to construct IO paths

Parameters

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

Returns full path to directory

Return type `str`

`eskapade.core.persistence.record_file_number` (*io_conf, file_name_base, file_name_ext*)

Function to get next prediction-record file number

Parameters

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

Returns next prediction-record file number

Return type `int`

`eskapade.core.persistence.repl_whites` (*name*)

eskapade.core.process_manager module

class `eskapade.core.process_manager.ProcessManager`

Bases: `eskapade.mixins.LoggingMixin`, `eskapade.mixins.TimerMixin`

Eskapade run-process manager

The `processManager` singleton class forms the core of Eskapade. It performs initialization, execution, and finalizing of the configured chains. Chains are added to the `processManager` (PM) thusly:

```
>>> proc_mgr = ProcessManager()
>>> proc_mgr.add_chain('Data')
>>> proc_mgr.add_chain('MyOverview')
```

The function:

```
>>> proc_mgr.execute_all()
```


executes all chains. This function is called by the `run_escaped.py` (the executable script of this project). The chains are executed in the order in which they have been added to the PM. One can begin and end the execution at chains specified in the configuration.

The `run_escaped.py` (the main of this project) script does the following:

- Imports and instantiates `processManager`
- Parses and process the cmd line arguments
- Executes the python configuration macro -> chains and links are defined in the PM
- Executes the PM

To be precise, `proc_mgr.execute_all()` does the following:

- initialize():** For each Chain, set name of previous Chain. Needed to pick up correct versions of persisted service instances.
- execute_all():** For each added Chain:
 - Initialize Chain: instantiates internal variables, and initialize each Link
 - Execute Chain: execute each Link
 - Finalize Chain: finalize each Link; if setting is true export datastore and configurations for each intermediate chain
- finalize():** Finalizes execution

Print ()

Print process-manager summary

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

__init__ ()

Initialize process-manager singleton instance

The init function takes no arguments. Chains are added with the “add_chain” method.

add_chain (input_chain)

Add a chain to the process manager

Add a chain to be run by the process manager. A check is performed that a chain with this name does not already exist.

Parameters `input_chain` (*str or Chain*) – (name of) the chain to be added

Raises `RuntimeError` – if chain with same name already exists

Returns the chain that has been added

Return type *Chain*

static check_io_config (io_conf)

Check I/O config and set name/version if not specified

Parameters `io_conf` (*dict*) – I/O config to check

execute (chain)

Execute a particular chain

Execution of a chain comprises:

- Initialize Chain:**
 - Instantiates internal variables

- Initialize each Link

•**Execute Chain:**

- Execute each Link

•**Finalize Chain:**

- Finalize each Link
- If setting is true, export datastore and configurations for each intermediate chain

Parameters `chain` – The chain to execute

Returns status code of execution attempt

Return type *StatusCode*

execute_all ()

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*)

Find the chain with the given name

Parameters **name** (*str*) – The name to search for

Returns the found chain

Return type *Chain*

Raises **RuntimeError** – if chain name not found

get_chain_idx (*name*)

Find index of the chain with given name

Parameters **name** (*str*) – The name of the chain to search for

Returns the index of the chain

Return type `int`

Raises **Exception** – if chain name not 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

has_chain (*name*)
 Check if chain exists for this name

Parameters **name** (*str*) – the name of the chain to check

Returns boolian of search result

Return type bool

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

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_chain (*name*)
 Remove chain with specified name

Remove specified chain. If chain is not found, print warning.

Parameters **name** (*str*) – Name of the chain to remove

remove_chains ()

Remove all configured chains

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

Remove specified process service

Parameters

- **serv_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.

service (*serv_spec*)

Get or register process service

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

Returns registered instance

Return type *ProcessService*

eskapade.core.process_services module

class `eskapade.core.process_services.ConfigObject`

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

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.
>>> import logging
>>> logging.basicConfig(level=logging.DEBUG)
```

Obtain the ConfigObject from any location as follows:

```
>>> from eskapade import ProcessManager, ConfigObject
>>> proc_mgr = ProcessManager()
>>> settings = proc_mgr.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).

Print ()

Print a summary of the settings

__init__ ()

Initialize ConfigObject instance

add_macros (macro_paths)

Add configuration macros for Eskapade run

io_base_dirs ()

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.
>>> import logging
>>> logging.basicConfig(level=logging.DEBUG)
```

Obtain the global datastore from any location as follows:

```
>>> from eskapade import ProcessManager, DataStore
>>> proc_mgr = ProcessManager()
>>> ds = proc_mgr.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: `eskapade.mixins.LoggingMixin`

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`

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_elements module**class** `eskapade.core.run_elements.Chain` (*name*)Bases: `eskapade.mixins.LoggingMixin`, `eskapade.mixins.TimerMixin`

Chain of links

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’.

Chains are added to the processManager (PM) thusly:

```
>>> proc_mgr = ProcessManager()
>>> overview = proc_mgr.add_chain('Overview')
```

And Links are added to a chain as follows:

```
>>> # add a link to the chain
>>> from eskapade import analysis
>>> overview.add_link(analysis.ReadToDf(path='foo.csv', key='foo'))
```

A Chain has an initialize, execute, and finalize method.

- `initialize()`: Initialise internal chain variables Then initialize the links in the chain.
- `execute()`: Execute the links in the chain
- `finalize()`: Finalize the links in the chain If configured, then store datastore and configuration

__init__ (*name*)

Initialize the new chain with certain name

Parameters **name** (*str*) – name of the chain**add_link** (*obj*)

Add link as a pre-built object

Parameters **obj** – The link to add**Returns** the link just added**Return type** *Link***execute** ()

Initialize, execute, finalize links in the chain

Returns status code of execution attempt**Return type** *StatusCode***finalize** ()

Finalize the chain and the links in the chain

Returns status code of finalization attempt**Return type** *StatusCode***get_link** (*name*)

Find the link with the given name

Parameters **name** (*str*) – The name of the link to search for**Returns** the link just found

Return type *Link*

Raises **RuntimeError** – if link name not found

has_link (*name*)

Check if link with name exists in this chain

Parameters **name** (*str*) – The name of the link to search for

Returns boolean answer

Return type bool

initialize ()

Initialize internal variables and links

Returns status code of initialization attempt

Return type *StatusCode*

remove_link (*mod*)

Remove link from this chain

Parameters **mod** – Link of the chain to remove

remove_links ()

Remove all links in the chain

class `eskapade.core.run_elements.Link` (*name*)

Bases: `eskapade.mixins.ArgumentsMixin`, `eskapade.mixins.LoggingMixin`, `eskapade.mixins.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:

```
>>> # add a chain first
>>> proc_mgr = ProcessManager()
>>> overview = proc_mgr.add_chain('Overview')
>>>
>>> # add a link to the chain
>>> from eskapade import analysis
>>> reader = analysis.ReadToDf(name='CsvReader', key='foo')
>>> reader.path = 'foo.csv'
>>> overview.add_link(reader)
```

A link has an `initialize`, `execute`, and `finalize` method. `execute()` is the central function that executes the algorithm code. `initialize()` and `finalize()` are supporting functions before and after `execute()`.

- `initialize()`: Initialise internal link variables, if needed

```
>>> # done here by hand, but normally done in the process manager
>>> pds.initialize()
```

- `execute()`: Execute the algorithm


```
>>> # by hand
>>> pds.execute()
```

- `finalize()`: Finalize the link, if needed

```
>>> # by hand
>>> pds.finalize()
```

See the documentation of these methods, or example links for more info.

Parameters `name` – The name of the link

`__init__` (*name*)

Initialize link instance

`execute` ()

Execute the link

This function is supposed to be overloaded by the actual link.

Returns status code of execute attempt

Return type *StatusCode*

`execute_link` ()

Execute the link

Returns status code of execute attempt

Return type *StatusCode*

`finalize` ()

Finalizing the link

This function is supposed to be overloaded by the actual link.

Returns status code of finalize attempt

Return type *StatusCode*

`finalize_link` ()

Finalizing the link

Returns status code of finalize attempt

Return type *StatusCode*

`initialize` ()

Initialize the link

This function is supposed to be overloaded by the actual link.

Returns status code of initialize attempt

Return type *StatusCode*

`initialize_link` ()

Initialize the link

Returns status code of initialize attempt

Return 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.ifInputMissing = statuscode*

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

Return type (*StatusCode*,list)

name

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.ifOutputExists = 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.run_utils module

`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

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

Bases: `eskapade.core.run_elements.Link`

Asserts that specified item(s) exists in the datastore

`__init__ (**kwargs)`
Store the configuration of link AssertInDs

Parameters

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

`execute ()`
Execute AssertInDs

eskapade.core_ops.links.break_link module

class `eskapade.core_ops.links.break_link.Break (**kwargs)`
Bases: `eskapade.core.run_elements.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 Break instance

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

`execute ()`
Execute Break

eskapade.core_ops.links.ds_object_deleter module

class `eskapade.core_ops.links.ds_object_deleter.DsObjectDeleter (**kwargs)`
Bases: `eskapade.core.run_elements.Link`

Delete objects from data store

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

`__init__ (**kwargs)`
Initialize DsObjectDeleter instance

Parameters

- **name** (*str*) – name of link
- **deletionKeys** (*list*) – keys to clear. Overwrites clearAll to false.
- **deletionClasses** (*list*) – delete object(s) by class type.
- **keepOnly** (*lst*) – keys to keep. Overwrites clearAll to false.
- **clearAll** (*bool*) – clear all key-value pairs in the datastore. Default is true.

`execute ()`
Execute DsObjectDeleter

`initialize ()`
Initialize DsObjectDeleter

eskapade.core_ops.links.ds_to_ds module

class `eskapade.core_ops.links.ds_to_ds.DsToDs` (**kwargs)

Bases: `eskapade.core.run_elements.Link`

Moves or copies an object in the datastore.

`__init__` (**kwargs)

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

Parameters

- **name** (*str*) – name of link
- **readKey** (*str*) – key of data to read from data store
- **storeKey** (*str*) – key of data to store in data store
- **move** (*bool*) – move readKey item to storeKey. Default is true.
- **copy** (*bool*) – if True the readKey key, value pair will not be deleted. Default is false.
- **remove** (*bool*) – if True the item corresponding to readKey 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 DsToDs

`initialize` ()

Initialize DsToDs

eskapade.core_ops.links.event_looper module

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

Bases: `eskapade.core.run_elements.Link`

EventLooper algorithm processes input lines and reprints them

`__init__` (**kwargs)

EventLooper processes input lines and reprints or stores them.

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

Parameters

- **name** (*str*) – name of link
- **filename** (*str*) – file name where the strings are located (txt or similar). Default is None. (optional)
- **storeKey** (*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

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

Bases: `eskapade.core.run_elements.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 HelloWorld

eskapade.core_ops.links.ipython_embed module

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

Bases: `eskapade.core.run_elements.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*)

Init of link IPythonEmbed

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

execute ()

Execute IPythonEmbed

eskapade.core_ops.links.line_printer module

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

Bases: `eskapade.core.run_elements.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
- **readKey** (*str*) – key of input data to read from data store

execute ()

Execute LinePrinter

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 LinePrinter

eskapade.core_ops.links.print_ds module

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

Bases: `eskapade.core.run_elements.Link`

Print the content of the datastore

__init__ (**kwargs)

Print overview of the datastore in current state

Parameters

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

execute ()

Execute PrintDs

eskapade.core_ops.links.repeat_chain module

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

Bases: `eskapade.core.run_elements.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
- **listenTo** (*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 RepeatChain

initialize ()

Initialize RepeatChain

eskapade.core_ops.links.skip_chain_if_empty module

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

Bases: `eskapade.core.run_elements.Link`

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

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

`__init__` (**kwargs)

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

Parameters

- **name** (*str*) – name of link
- **collectionSet** (*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
- **checkAtInitialize** (*bool*) – perform dataset empty is check at initialize. Default is true.
- **checkAtExecute** (*bool*) – perform dataset empty is check at initialize. Default is false.

checkCollectionSet ()

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 SkipChainIfEmpty

initialize ()

Initialize SkipChainIfEmpty

eskapade.core_ops.links.to_ds_dict module

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

Bases: `eskapade.core.run_elements.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.

dostorage (*ds*)

perform storage in datastore

Function makes a distinction between dicts and any other object.

execute ()

Execute ToDsDict

initialize ()

Initialize ToDsDict

Module contents

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

Bases: `eskapade.core.run_elements.Link`

Asserts that specified item(s) exists in the datastore

__init__ (**kwargs)

Store the configuration of link AssertInDs

Parameters

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

execute ()

Execute AssertInDs

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

Bases: `eskapade.core.run_elements.Link`

Delete objects from data store

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

__init__ (**kwargs)

Initialize DsObjectDeleter instance

Parameters

- **name** (*str*) – name of link
- **deletionKeys** (*list*) – keys to clear. Overwrites clearAll to false.
- **deletionClasses** (*list*) – delete object(s) by class type.
- **keepOnly** (*list*) – keys to keep. Overwrites clearAll to false.
- **clearAll** (*bool*) – clear all key-value pairs in the datastore. Default is true.

execute ()

Execute DsObjectDeleter

initialize()
Initialize DsObjectDeleter

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

Bases: `eskapade.core.run_elements.Link`

Moves or copies an object in the datastore.

__init__ (**kwargs)
Link to move, copy, or remove an object in the datastore.

Parameters

- **name** (*str*) – name of link
- **readKey** (*str*) – key of data to read from data store
- **storeKey** (*str*) – key of data to store in data store
- **move** (*bool*) – move readKey item to storeKey. Default is true.
- **copy** (*bool*) – if True the readKey key, value pair will not be deleted. Default is false.
- **remove** (*bool*) – if True the item corresponding to readKey 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 DsToDs

initialize()
Initialize DsToDs

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

Bases: `eskapade.core.run_elements.Link`

EventLooper algorithm processes input lines and reprints them

__init__ (**kwargs)
EventLooper processes input lines and reprints or stores them.

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

Parameters

- **name** (*str*) – name of link
- **filename** (*str*) – file name where the strings are located (txt or similar). Default is None. (optional)
- **storeKey** (*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.LinePrinter` (**kwargs)
Bases: `eskapade.core.run_elements.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
- **readKey** (*str*) – key of input data to read from data store

execute()
Execute LinePrinter

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 LinePrinter

class `eskapade.core_ops.links.HelloWorld` (**kwargs)
Bases: `eskapade.core.run_elements.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 HelloWorld

class `eskapade.core_ops.links.PrintDs` (**kwargs)
Bases: `eskapade.core.run_elements.Link`

Print the content of the datastore

__init__ (**kwargs)
Print overview of the datastore in current state

Parameters

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

execute()
Execute PrintDs

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

Bases: `eskapade.core.run_elements.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
- **listenTo** (*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 RepeatChain

initialize ()

Initialize RepeatChain

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

Bases: `eskapade.core.run_elements.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)

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

Parameters

- **name** (*str*) – name of link
- **collectionSet** (*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
- **checkAtInitialize** (*bool*) – perform dataset empty is check at initialize. Default is true.
- **checkAtExecute** (*bool*) – perform dataset empty is check at initialize. Default is false.

checkCollectionSet ()

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 SkipChainIfEmpty

initialize ()
 Initialize SkipChainIfEmpty

class `eskapade.core_ops.links.ToDsDict (**kwargs)`
 Bases: `eskapade.core.run_elements.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.

dostorage (ds)
 perform storage in datastore
 Function makes a distinction between dicts and any other object.

execute ()
 Execute ToDsDict

initialize ()
 Initialize ToDsDict

class `eskapade.core_ops.links.IPythonEmbed (**kwargs)`
 Bases: `eskapade.core.run_elements.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)**
 Init of link IPythonEmbed

Parameters name (*str*) – name of link

execute ()
 Execute IPythonEmbed

class `eskapade.core_ops.links.Break (**kwargs)`
 Bases: `eskapade.core.run_elements.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 Break instance

Parameters name (*str*) – name of link

```
execute ()
    Execute Break
```

Module contents

eskapade.data_quality package

Subpackages

eskapade.data_quality.links package

Submodules

eskapade.data_quality.links.fix_pandas_dataframe module

```
class eskapade.data_quality.links.fix_pandas_dataframe.FixPandasDataFrame (**kwargs)
    Bases: eskapade.core.run_elements.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 FixPandasDataFrame 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 FixPandasDataFrame

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 FixPandasDataFrame
```

```
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.run_elements.Link
```

Fix dirty Pandas dataframe with inconsistent datatypes

Default settings perform the following clean-up steps on an input dataframe:

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

Boolean columns with contamination get converted to string columns by default. Optionally, they can be converted to integer columns as well.

The `FixPandasDataFrame` link can be used in a dataframe loop, in which case any data type assessed per column in the first dataframe iteration will be used for the next dataframes as well.

The default settings should work pretty well in many circumstances, but can be configured pretty flexibly. Optionally:

- Instead of dynamically assessed, the data type can also be imposed per column
- All nans in a column can be converted to a value consistent with the data type of that column. E.g. for integer columns, `nan -> -999`
- An alternative nan can be set per column and datatype
- Modifications can be applied inplace, i.e. directly to the input dataframe

```
__init__ (**kwargs)
    Initialize FixPandasDataFrame 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 FixPandasDataFrame

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 FixPandasDataFrame
```

eskapade.data_quality.tests package

Subpackages

eskapade.data_quality.tests.integration package

Submodules

eskapade.data_quality.tests.integration.test_tutorial_macros module

Module contents

Module contents

Submodules

eskapade.data_quality.dq_helper module

```
eskapade.data_quality.dq_helper.bool_to_int(val)  
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.root_analysis` package

Subpackages

`eskapade.root_analysis.decorators` package

Submodules

`eskapade.root_analysis.decorators.histograms` module

`eskapade.root_analysis.decorators.histograms.bin_centers(self, q)`

Get bin centers of histogram axis

Parameters `q` (`int`) – axis index, should be 0 (=x), 1 (=y), or 2(=z)

Returns numpy array of bin centers

Return type numpy.array

eskapade.root_analysis.decorators.histograms.**bin_edges** (*self*, *q*)
Get bin edges of histogram axis

Parameters *q* (*int*) – axis index, should be 0 (=x), 1 (=y), or 2(=z)

Returns numpy array of bin edges

Return type numpy.array

eskapade.root_analysis.decorators.histograms.**bin_entries** (*self*)
Get bin values

Returns numpy array of bin entries

Return type numpy.array

eskapade.root_analysis.decorators.histograms.**bin_entries_2dgrid** (*self*)
Get 2-D grid of bin entries

Returns 2-d grid of bin entries

eskapade.root_analysis.decorators.histograms.**bin_labels** (*self*, *q*)
Get bin labels of histogram axis

Parameters *q* (*int*) – axis index, should be 0 (=x), 1 (=y), or 2(=z)

Returns numpy array of bin labels

Return type numpy.array

eskapade.root_analysis.decorators.histograms.**bin_range** (*self*, *q*)
Get bin range of histogram axis

Parameters *q* (*int*) – axis index, should be 0 (=x), 1 (=y), or 2(=z)

Returns tuple of bin range

Return type tuple

eskapade.root_analysis.decorators.histograms.**bin_vals** (*self*)
Get NumPy-style histogram values

Returns comma-separated np array of bin entries, np array of bin edges

Return type numpy.array, numpy.array

eskapade.root_analysis.decorators.histograms.**datatype**
Data type of histogram variable

Return data type of the variable represented by the histogram. If not already set, will determine datatype automatically.

Returns data type

Return type type or list(type)

eskapade.root_analysis.decorators.histograms.**high**
High edge of last bin

Returns up edge of first bin

Return type float

eskapade.root_analysis.decorators.histograms.**low**
Low edge of first bin

Returns low edge of first bin

Return type float

`eskapade.root_analysis.decorators.histograms.n_bins`
Number of bins of x axis

Returns number of bins for x axis

Return type int

`eskapade.root_analysis.decorators.histograms.n_dim`
Histogram dimension

`eskapade.root_analysis.decorators.histograms.xy_ranges_grid(self)`
Get x and y ranges and x, y grid

Returns x and y ranges and x,y grid

eskapade.root_analysis.decorators.roofit module

`eskapade.root_analysis.decorators.roofit.coll_iter(coll)`
Iterate over items in RooAbsCollection

`eskapade.root_analysis.decorators.roofit.data_set_iter(self)`
Iterate over events in RooDataSet

`eskapade.root_analysis.decorators.roofit.ws_contains(ws, key)`
Check if RooWorkspace contains an object

`eskapade.root_analysis.decorators.roofit.ws_put(ws, *args)`
Put object in RooWorkspace

`eskapade.root_analysis.decorators.roofit.ws_setitem(ws, key, value)`
Put object in RooWorkspace dict-style

Module contents

eskapade.root_analysis.links package

Submodules

eskapade.root_analysis.links.add_propagated_error_to_roodataset module

class `eskapade.root_analysis.links.add_propagated_error_to_roodataset.AddPropagatedErrorToRooDataSet`
Bases: `eskapade.core.run_elements.Link`

Evaluates errors on a provided roofit function

The evaluated errors are added as a new column to the input dataset. Optionally, the evaluated function values are added as a column to the input dataset as well.

`__init__` (***kwargs*)
Initialize AddPropagatedErrorToRooDataSet instance

Parameters

- **name** (*str*) – name of link

- **data** (*str*) – key of input data to read from data store or workspace
- **function** (*str*) – key of rooabsreal function to read from data store or workspace
- **fit_result** (*str*) – key of roofitresult object to read from data store or workspace, used as input for error propagation
- **from_ws** (*bool*) – if true, try to retrieve data, function, and fit_result from workspace instead of datastore. Default is false.
- **function_error_name** (*str*) – column name assigned to propagated errors that are appended to data
- **add_function_to_data** (*bool*) – add column of the function values to the data. Default is true

execute ()

Execute AddPropagatedErrorToRooDataSet

initialize ()

Initialize AddPropagatedErrorToRooDataSet

eskapade.root_analysis.links.convert_dataframe_2_roodataset module

class `eskapade.root_analysis.links.convert_dataframe_2_roodataset.ConvertDataFrame2RooDataSet`

Bases: `eskapade.core.run_elements.Link`

Convert Pandas dataframe into a RooFit dataset

By default all observables of the dataframe are interpreted as continuous (not category), except for boolean and numpy category variables. Other category observable first need to be converted ('factorized') to numeric values, eg. using the link record_factorizer. These other category variables can be picked up by setting: `map_to_factorized`, which is dictionary to map columns to factorized ones. `map_to_factorized` is a dict of dicts, ie. one dict for each column.

RooDataHistFiller stores a roodatahist object, a rooargset containing all corresponding roofit observables and roocategories. Also stored, under key `sk_map_to_original`, is a dictionary to map all factorized columns back to original.

For each observable one can set the number of bins, and min and max values. The total number of bins in the roodatahist may not exceed `n_max_total_bins`.

`__init__` (**kwargs)

Initialize ConvertDataFrame2RooDataSet instance

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **read_key_vars** (*str*) – key of input rooargset of observables from data store (optional)
- **columns** (*list*) – list of columns to pick up from dataset. Default is all columns. (optional)
- **ignore_columns** (*list*) – list of columns to ignore from dataset. (optional)
- **store_key** (*str*) – key of output roodataset to store in data store (optional)
- **store_key_vars** (*str*) – key of output rooargset of observables to store in data store. (optional)

- **into_ws** (*bool*) – if true, store in workspace, not datastore. Default is True
- **rm_original** (*bool*) – if true, remove original histogram. Default is False
- **map_to_factorized** (*dict*) – dictionary or key to dictionary to map back columns to factorized ones. `map_to_factorized` is a dict of dicts, one dict for each column. (optional)
- **sk_map_to_original** (*str*) – store key of dictionary to map factorized columns to original. Default is 'key' + '_' + store_key + '_to_original'. (optional)
- **n_max_total_bins** (*int*) – max number of bins in roodatahist. Default is 1e6.
- **store_index** (*bool*) – If true, copy df's index to rds. Default is true.
- **create_keys_pdf** (*str*) – if set, create keys pdf from rds with this name and add to ds or workspace (optional)

execute ()
Execute ConvertDataFrame2RooDataSet

initialize ()
Initialize ConvertDataFrame2RooDataSet

eskapade.root_analysis.links.convert_roodataset_2_dataframe module

class `eskapade.root_analysis.links.convert_roodataset_2_dataframe.ConvertRooDataSet2DataFrame`
Bases: `eskapade.core.run_elements.Link`

Convert an input RooFit dataset into a Pandas dataframe

Input roodataset can be picked up from either datastore or rooworkspace. The output dataframe is stored in the datastore.

__init__ (**kwargs)
Initialize ConvertRooDataSet2DataFrame instance

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input roodataset to read from data store or workspace
- **store_key** (*str*) – key of output data to store in data store (optional)
- **from_ws** (*bool*) – if true, pick up input roodataset from workspace, not datastore. Default is false.
- **rm_original** (*bool*) – if true, input roodataset is removed from ds/ws. Default is false.

execute ()
Execute ConvertRooDataSet2DataFrame

initialize ()
Initialize ConvertRooDataSet2DataFrame

eskapade.root_analysis.links.convert_roodataset_2_roodatahist module

class `eskapade.root_analysis.links.convert_roodataset_2_roodatahist.ConvertRooDataSet2RooData`
Bases: `eskapade.core.run_elements.Link`

Convert input RooFit dataset into a Pandas dataframe

Input RooDataSet can be picked up from either the data store or the workspace. The output data frame is stored in the data store.

`__init__` (**kwargs)
Initialize ConvertRooDataSet2RooDataHist instance

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input roodataset to read from data store or workspace
- **store_key** (*str*) – key of output data to store in data store (optional)
- **from_ws** (*bool*) – if true, pick up input roodataset from workspace instead of data store (default is False)
- **rm_original** (*bool*) – if true, input roodataset is removed from ds/ws (default is False)
- **columns** (*list*) – columns to pick up from input RooDataSet
- **binning_name** (*str*) – name of binning configuration with which to construct RooDataHist

`execute` ()
Execute ConvertRooDataSet2RooDataHist

`initialize` ()
Initialize ConvertRooDataSet2RooDataHist

eskapade.root_analysis.links.convert_root_hist_2_roodatahist module

class `eskapade.root_analysis.links.convert_root_hist_2_roodatahist.ConvertRootHist2RooDataHist`
Bases: `eskapade.core.run_elements.Link`

Convert a ROOT histogram into a RooFit histogram

Input histograms can have up to three dimensions. RooFit observables are deduced from the histogram axes. By default all observables are interpreted as continuous.

ConvertRootHist2RooDataHist stores a roodatahist object, a rooarglist containing all corresponding roofit observables. Optionally, a roohistpdf is created from the roodatahist object and stored as well.

`__init__` (**kwargs)
Initialize ConvertRootHist2RooDataHist instance

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – histogram to pick up from datastore (or, if set, from histogram dict)
- **hist_dict_key** (*str*) – histograms dictionary from data store. If set, the histogram is read from this dict (optional)
- **store_key** (*str*) – key of roodatahist (optional)
- **into_ws** (*bool*) – if true, store in workspace, not datastore. Default is True.
- **rm_original** (*bool*) – if true, remove original histogram. Default is False.
- **create_hist_pdf** (*str*) – if set, create keys pdf from roodatahist with this name and add to ds or workspace

```
execute ()  
    Execute ConvertRootHist2RooDataHist  
initialize ()  
    Initialize ConvertRootHist2RooDataHist
```

eskapade.root_analysis.links.convert_root_hist_2_roodataset module

```
class eskapade.root_analysis.links.convert_root_hist_2_roodataset.ConvertRootHist2RooDataSet  
    Bases: eskapade.core.run_elements.Link
```

Convert a ROOT histogram into a RooFit dataset

Input histograms can have up to three dimensions. RooFit observables are deduced from the histogram axes. By default all observables are interpreted as continuous.

ConvertRootHist2RooDataSet stores a roodataset object, a rooargset containing all corresponding roofit observables. Optionally, a keys-pdf is created from the roodataset object and stored as well.

```
__init__ (**kwargs)  
    Initialize ConvertRootHist2RooDataSet instance
```

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – histogram to pick up from datastore (or, if set, from histogram dict)
- **hist_dict_key** (*str*) – histograms dictionary from data store. If set, histogram is read from this dict (optional)
- **store_key** (*str*) – key of roodatahist (optional)
- **store_key_vars** (*str*) – key of output rooargset of observables to store in data store (optional)
- **into_ws** (*bool*) – if true, store in workspace, not datastore. Default is True.
- **rm_original** (*bool*) – if true, remove original histogram. Default is False.
- **create_keys_pdf** (*str*) – if set, create keys pdf from rds with this name and add to ds or workspace

```
execute ()  
    Execute ConvertRootHist2RooDataSet  
initialize ()  
    Initialize ConvertRootHist2RooDataSet
```

eskapade.root_analysis.links.print_ws module

```
class eskapade.root_analysis.links.print_ws.PrintWs (**kwargs)  
    Bases: eskapade.core.run_elements.Link
```

Print the contents of the RooFit workspace

```
__init__ (**kwargs)  
    Initialize PrintWs instance
```

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

execute ()
Execute PrintWs

eskapade.root_analysis.links.read_from_root_file module

eskapade.root_analysis.links.roodatahist_filler module

class `eskapade.root_analysis.links.roodatahist_filler.RooDataHistFiller (**kwargs)`
Bases: `eskapade.core.run_elements.Link`

Fill a RooFit histogram with columns from a Pandas dataframe

Histograms can have any number of dimensions. Only numeric observables are picked up. By default all observables are interpreted as continuous (not category), except for boolean and numpy category variables. Other category observable first need to be converted ('factorized') to numeric values, eg. using the `link_record_factorizer`. These other category variables can be picked up by setting: `map_to_factorized`, which is dictionary to map columns to factorized ones. `map_to_factorized` is a dict of dicts, ie. one dict for each column.

`RooDataHistFiller` stores a `roodatahist` object, a `rooargset` containing all corresponding roofit observables and `roocategories`, and a `rooargset` containing only the observables that are `roocategories`. Also stored, under key `sk_map_to_original`, is a dictionary to map all factorized columns back to original.

For each observable one can set the number of bins, and min and max values. The total number of bins in the `roodatahist` may not exceed `n_max_total_bins`.

The `roodatahist` histogram can be filled iteratively, while looping over multiple dataframes.

__init__ (kwargs)**
Initialize `RooDataHistFiller` instance

Parameters

- **name** (*str*) – name of link
- **read_key** (*str*) – key of input data to read from data store
- **columns** (*list*) – list of columns to pick up from dataset. Default is all columns. (optional)
- **ignore_columns** (*list*) – list of columns to ignore from dataset. (optional)
- **store_key** (*str*) – key of output roodataset to store in data store. (optional)
- **store_key_vars** (*str*) – key of output rooargset of all observables to store in data store. (optional)
- **store_key_cats** (*str*) – key of output rooargset of category observables to store in data store. (optional)
- **store_at_finalize** (*bool*) – if true, store in workspace at `finalize()`, not at `execute()`. (optional)
- **into_ws** (*bool*) – if true, store in workspace, not datastore. Default is False. (optional)
- **rm_original** (*bool*) – if true, remove original dataframe. Default is False. (optional)
- **map_to_factorized** (*dict*) – dictionary or key to dictionary to map columns to factorized ones. `map_to_factorized` is a dict of dicts, one dict for each column. (optional)
- **sk_map_to_original** (*str*) – store key of dictionary to map factorized columns to original. Default is 'key' + '_' + `store_key` + '_to_original'. (optional)

- **var_number_of_bins** (*dict*) – number of bins for histogram of certain variable (optional)
- **var_min_value** (*dict*) – min value for histogram of certain variable (optional)
- **var_max_value** (*dict*) – max value for histogram of certain variable (optional)
- **n_max_total_bins** (*int*) – max number of bins in roodatahist. Default is 1e6. (optional)
- **create_hist_pdf** (*str*) – if filled, create hist pdf from rdh with this name and add to datastore or workspace. (optional)

do_storage()

Storage of the created RooDataHist object

execute()

Execute RooDataHistFiller

Fill a roodatahist object with a pandas dataframe. It is possible to fill the roodatahist iteratively, in a loop over dataframes.

There are 5 steps to the code:

1. basic checks of the dataframe
2. convert the dataframe to a roodataset
3. instantiate a roodatahist object
4. fill the roodatahist object with the roodataset
5. store the roodatahist. optionally, at the storage stage a pdf can be created of the roodatahist as well.

finalize()

Finalize RooDataHistFiller

initialize()

Initialize RooDataHistFiller

eskapade.root_analysis.links.roofit_percentile_binning module

eskapade.root_analysis.links.root_hist_filler module

class `eskapade.root_analysis.links.root_hist_filler.RootHistFiller` (**kwargs)

Bases: `eskapade.analysis.histogram_filling.HistogramFillerBase`

Create ROOT histograms from columns in Pandas dataframe

Histograms can be up to 3 dimensions. The data type for the histogram can be automatically assessed from the column(s), or set as input. For each histogram (axis) one can set the label, number of bins, min and max values, and data type. String based columns are not allowed, for now.

The histograms are created at initialize() and can be filled iteratively, while looping over multiple dataframes.

__init__ (**kwargs)

Initialize RootHistFiller instance

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 in data store. Default is ‘hist’
- **columns** (*list*) – required list of column (variables) from input data to turn into histograms. Eg. [‘x’,‘y’, [‘y’,‘z’], [‘x’,‘y’,‘z’]]. A list of columns is interpreted as a multi-dimensional histogram. Up to 3 dimensions is allowed.
- **pair_up_columns** (*list*) – required list of column (variables) to pair up for 2-dim histograms
- **var_label** (*dict*) – title for histogram of certain variable (optional)
- **var_number_of_bins** (*dict*) – number of bins for histogram of certain variable (optional)
- **var_min_value** (*dict*) – min value for histogram of certain variable (optional)
- **var_max_value** (*dict*) – max value for histogram of certain variable (optional)
- **var_dtype** (*dict*) – impose data type of variable (optional)
- **weight** (*str*) – name of weight column (optional)

categorize_columns (*df*)

Categorize columns of dataframe by data type

Parameters *df* – input (pandas) data frame

construct_empty_hist (*columns*)

Create an (empty) histogram of right type

Create a multi-dim histogram by iterating through the columns.

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

Returns created ROOT histogram

Return type ROOT.TH1

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)

initialize ()

Initialize RootHistFiller

eskapade.root_analysis.links.trunc_exp_fit module

class `eskapade.root_analysis.links.trunc_exp_fit.TruncExpFit` (**kwargs)

Bases: `eskapade.core.run_elements.Link`

Fit truncated exponential PDF to data

Fit an exponential PDF in a range with a variable upper bound to data. That is, the PDF is truncated at a different value for each record in the data. A dataset with values for the PDF exponential variable and the upper bound must be provided.

Optionally an additional dataset with range upper bounds may be provided. This will be considered to be a dataset of samples, which may or may not have a corresponding event in the first dataset. In this case an estimate of the number of events without any range upper bounds will be provided, given the samples in the

upper-bound dataset. This result is also visualized in a plot of the PDFs with and without range upper bounds and histograms of the corresponding event data.

`__init__` (**kwargs)

Initialize the TruncExpFit instance

Parameters

- **name** (*str*) – name of link instance
- **read_key** (*str*) – data-store key of input data
- **max_var_data_key** (*str*) – data-store key of dataset with range upper-bound values
- **model_name** (*str*) – name of truncated-exponential model to use

`execute` ()

Execute TruncExpFit

`initialize` ()

Initialize the TruncExpFit execution

`eskapade.root_analysis.links.trunc_exp_fit.est_norm_ratio` (*norm_full*, *norm_data*,
fit_result)

Estimate total number of events without range bounds

`eskapade.root_analysis.links.trunc_exp_fit.make_plots` (*data*, *model*, *full_norm_ratio*,
plots_path)

Make plots of data and model

eskapade.root_analysis.links.trunc_exp_gen module

class `eskapade.root_analysis.links.trunc_exp_gen.TruncExpGen` (**kwargs)

Bases: `eskapade.core.run_elements.Link`

Generate with truncated exponential PDF

Generate data with an exponential PDF in a range with a variable upper bound. That is, the PDF is truncated at a different value for each record in the data.

A dataset with upper bounds must be provided. Each bound is considered to be a sample, for which an event may or not be generated. The fraction of samples for which an event is generated can be set. For each event, a value is generated with the exponential PDF truncated at the range upper bound for the corresponding sample.

`__init__` (**kwargs)

Initialize the TruncExpGen instance

Parameters

- **name** (*str*) – name of link instance
- **store_key** (*str*) – data-store key of generated data
- **max_var_data_key** (*str*) – data-store key of dataset with range upper-bound values
- **model_name** (*str*) – name of truncated-exponential model to use

`execute` ()

Execute TruncExpGen

`initialize` ()

Initialize the TruncExpGen execution

`eskapade.root_analysis.links.trunc_exp_gen.gen_max_var_data(model)`

Generate range upper-bound data

`eskapade.root_analysis.links.trunc_exp_gen.sel_max_var_data(model, max_var_data, event_frac)`

Select upper-bound data with PDF integral values

`eskapade.root_analysis.links.uncorrelation_hypothesis_tester` module

`eskapade.root_analysis.links.ws_utils` module

class `eskapade.root_analysis.links.ws_utils.WsUtils(**kwargs)`

Bases: `eskapade.core.run_elements.Link`

Apply standard operations to object in the RooFit workspace

Operations include: - moving object to and from the datastore/workspace - execute rooworkspace factory commands - simulation, use the function: `add_simulate()` - fitting, use the function: `add_fit()` - plotting, use the function `add_plot()`

`__init__` (***kwargs*)

Initialize WsUtils instance

Note: - For simulations use the function: `add_simulate()` - For fitting use the function: `add_fit()` - For plotting use the function: `add_plot()`

Parameters

- **name** (*str*) – name of link
- **copy_into_ws** (*list*) – key of input data to read from data store, to be inserted in rooworkspace
- **copy_into_ds** (*list*) – key of input data to read from rooworkspace, to be inserted in data store
- **rm_original** (*bool*) – if true, objects inserted in ws/ds are removed from ds/ws. Default is false.
- **rm_from_ws** (*list*) – rm keys from rooworkspace
- **factory** (*list*) – list of commands passed to workspace factory at `execute()`
- **apply** (*list*) – list of functions to pass workspace through at `execute()`
- **results_path** (*str*) – output path of plot (optional)
- **pages_key** (*str*) – data store key of existing report pages

`add_fit` (**args, **kwargs*)

Add fit task

Stash args and kwargs, to be executed by `do_fit()` during `execute()`

Parameters

- **args** – positional arguments passed on to `do_fit()`
- **kwargs** – key word arguments passed on to `do_fit()`

add_plot (**args, **kwargs*)

Add plotting task

Stash args and kwargs, to be executed by do_fit() during execute()

Parameters

- **args** – positional arguments passed on to do_plot()
- **kwargs** – key word arguments passed on to do_plot()

add_simulate (**args, **kwargs*)

Add simulation task

Stash args and kwargs, to be executed by do_simulate() during execute()

Parameters

- **args** – positional arguments passed on to do_simulate()
- **kwargs** – key word arguments passed on to do_simulate()

do_fit (*ds, ws, pdf, data, key='', replace_existing=True, into_ws=False, *args, **kwargs*)

Fit PDF to data set

Parameters

- **ds** – input data store, from which to retrieve pdf and dataset to fit
- **ws** (*ROOT.RooWorkspace*) – input workspace, from which to retrieve pdf and dataset to fit
- **pdf** – input pdf used for fitting, can be a key to look up or RooAbsPdf
- **data** – input dataset used for fitting, can be a key to look up or RooAbsData
- **key** (*str*) – key under which to store the fit result object
- **into_ws** (*bool*) – if true, store simulated data in workspace, not the datastore
- **args** – all other positional arguments are passed on to the roofit fit function.
- **kwargs** – all other key word arguments are passed on to the roofit fit function.

do_plot (*ds, ws, obs, data=None, pdf=None, func=None, data_args=(), pdf_args=(), func_args=(), data_kwargs={}, pdf_kwargs={}, func_kwargs={}, key='', into_ws=False, output_file=None, bins=40, logy=False, miny=0, plot_range=None*)

Make a plot of data and/or a pdf, or of a function

Either a dataset, pdf, or function needs to be provided as input for plotting.

Parameters

- **ds** – input data store, from which to retrieve pdf and dataset to fit
- **ws** (*ROOT.RooWorkspace*) – input workspace, from which to retrieve pdf and dataset to fit
- **data** – input dataset used for plotting, can be a key to look up or RooAbsData (optional)
- **pdf** – input pdf used for plotting, can be a key to look up or RooAbsPdf (optional)
- **func** – input function used for plotting, can be a key to look up or RooAbsReal (optional)
- **data_args** – positional arguments passed on to the plotting of the data. (optional)
- **data_kwargs** – key word arguments passed on to the plotting of the data. (optional)
- **pdf_args** – positional arguments passed on to the plotting of the pdf. (optional)

- **pdf_kwargs** – key word arguments passed on to the plotting of the pdf. (optional)
- **func_args** – positional arguments passed on to the plotting of the function. (optional)
- **func_kwargs** – key word arguments passed on to the plotting of the function. (optional)
- **key** (*str*) – key under which to store the plot frame (=RooPlot). If key exists in ds/workspace, plot in the existing frame. (optional)
- **into_ws** (*bool*) – if true, store simulated data in workspace, not the datastore
- **output_file** (*str*) – if set, store plot with this file name (optional)
- **bins** (*int*) – number of bins in the plot. default is 40. (optional)
- **logy** (*bool*) – if true, set y-axis to log scale (optional)
- **miny** (*float*) – set minimum value of y-axis to miny value (optional)
- **plot_range** (*tuple*) – specify x-axis plot range as (min, max) (optional)

do_simulate (*ds, ws, pdf, num, obs, key='', into_ws=False, *args, **kwargs*)
 Simulate data based on input PDF

Parameters

- **ds** – input data store, from which to retrieve pdf and observables
- **ws** (*ROOT.RooWorkspace*) – input workspace, from which to retrieve pdf and observables
- **pdf** – input pdf used for simulation, can be a key to look up or RooAbsPdf
- **num** (*int*) – number of records to generate
- **obs** – input observables used for simulation, can be a key to look up or a RooArgSet
- **key** (*str*) – key under which to store the simulated data
- **into_ws** (*bool*) – if true, store simulated data in workspace, not the datastore
- **args** – all other positional arguments are passed on to the roofit generate function.
- **kwargs** – all other key word arguments are passed on to the roofit generate function.

execute ()

Execute WsUtils

Operations are executed in this order:

- 1.put objects from the datastore into rooworkspace
- 2.execute rooworkspace factory commands
- 3.pass the workspace to (a list of) functions, to execute bits of (workspace) code
- 4.simulate data from a pdf
- 5.fit a pdf to a dataset
- 6.make a plot of a dataset, pdf, or function
- 7.move objects from the workspace to the datastore

finalize ()

Finalize WsUtils

initialize ()

Initialize WsUtils

Module contents

Submodules

eskapade.root_analysis.data_conversion module

eskapade.root_analysis.roofit_manager module

class `eskapade.root_analysis.roofit_manager.RooFitManager`

Bases: `eskapade.core.process_services.ProcessService`

Process service for managing RooFit operations

`__init__` ()

Initialize RooFit manager instance

model (*name*, *model_cls=None*, ***kwargs*)

Get RooFit model

Return the RooFit model with the specified name. Create the model if it does not yet exist. Arguments and key-word arguments are passed on to the model when it is initialized.

Parameters

- **name** (*str*) – name of the model
- **model_cls** – model class; must inherit from `RooFitModel`

set_var_vals (*vals*)

Set values of workspace variables

Parameters **vals** (*dict*) – values and errors to set: {name1: (val1, err1), name2: (val2, err2), ...}

ws

RooFit workspace for Eskapade run

eskapade.root_analysis.roofit_models module

class `eskapade.root_analysis.roofit_models.LinearRegression` (*ws*, *name=''*,
fit_intercept=True,
minimizer='Minuit2',
strategy=2)

Bases: `eskapade.root_analysis.roofit_models.RooFitModel`

Least-squares linear regression model

`__init__` (*ws*, *name=''*, *fit_intercept=True*, *minimizer='Minuit2'*, *strategy=2*)

Initialize LinearRegression instance

Parameters

- **ws** (*ROOT.RooWorkspace*) – RooFit workspace
- **fit_intercept** (*bool*) – build model with variable intercept parameter
- **minimizer** (*str*) – minimization tool for fit (minimizer type for `ROOT.RooMinimizer`)
- **strategy** (*int*) – fit strategy (strategy for `ROOT.RooMinimizer`)

build_model ()
Build prediction function and corresponding χ^2

coef_
Return fitted values of coefficients

fit (*X*, *y*)
Fit model to specified data

Parameters

- **X** – samples of features
- **y** – targets

fit_dataset ()
Fit model to pre-loaded data

intercept_
Return fitted intercept value

load_data (*X*, *y*)
Load target and feature data

predict (*X*)
Predict values for given samples of features

score (*X*, *y*)
Calculate R^2 coefficient for specified samples

class `eskapade.root_analysis.roofit_models.RooFitModel` (*ws*, *name*='',
load_libesroofit=False)

Bases: object

Base class for RooFit models

__init__ (*ws*, *name*='', *load_libesroofit*=False)
Initialize model instance

Parameters

- **ws** (*ROOT.RooWorkspace*) – RooFit workspace
- **name** (*str*) – name of model
- **load_libesroofit** (*bool*) – load Eskapade RooFit library upon initialization (default is False)

is_built
Is model built in workspace

name
Name of model

pdf
Model PDF

pdf_name
Name of model PDF

ws
Workspace to create model in

```
class eskapade.root_analysis.roofit_models.TruncExponential (ws, name='',
                                                           var_range=None,
                                                           var=None,
                                                           max_var=None,
                                                           exp=None,
                                                           fracs=None)
```

Bases: `eskapade.root_analysis.roofit_models.RooFitModel`

Exponential model with variable range upper bound

```
__init__ (ws, name='', var_range=None, var=None, max_var=None, exp=None, fracs=None)
Initialize TruncExponential instance
```

Parameters

- **ws** (`ROOT.RooWorkspace`) – RooFit workspace
- **name** (`str`) – name of model
- **var_range** (`tuple`) – range of PDF variable: (min, max)
- **var** (`tuple`) – PDF variable: (name, value)
- **max_var** (`tuple`) – variable upper limit of PDF range: (name, value)
- **exp** (`list`) – list of exponential-rate parameters: [(name0, value0), (name1, value1), ...]
- **fracs** (`list`) – list of exponential-fraction parameters: [(name0, value0), ...]

all_vars_set

Set containing PDF variables

build_model()

Build truncated exponential model in workspace

create_norm (`data=None, range_min=None, range_max=None`)

Create PDF normalization-integral objects

Calculate two objects representing normalization integrals of the PDF function. The first integral is calculated in the range determined by the specified minimum value and the range upper bounds in the provided dataset. The second integral is calculated in the specified range.

Parameters

- **data** (`ROOT.RooDataSet`) – dataset with upper bounds of range
- **range_min** (`float`) – lower bound of integration range
- **range_max** (`float`) – upper bound of integration range

Returns normalization integrals

Return type tuple

max_var

Range upper bound of PDF variable

max_var_set

Set containing range upper bound of PDF variable

var

PDF variable

var_set

Set containing PDF variable

eskapade.root_analysis.roofit_utils module

```
eskapade.root_analysis.roofit_utils.create_roofit_opts (opts_type='', create_linked_list=True,
**kwargs)
```

Build list of options for RooFit functions

Additional keyword arguments are appended to default options list for specified type. The keys must be names of functions in the RooFit namespace that return a RooCmdArg, while the values are the arguments of the specified function:

```
>>> opts = create_roofit_opts('data_plot', LineColor=ROOT.kRed, Asymmetry=cat,
↳Binning=(10, 0., 10.))
>>> data.plotOn(frame, opts)
```

Parameters `opts_type` (*str*) – options type {'', 'fit', 'pdf_plot', 'data_plot', 'obs_frame'}

Returns collection of RooFit options

Return type dict or ROOT.RooLinkedList

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

Load Eskapade RooFit library

```
eskapade.root_analysis.roofit_utils.set_rf_log_level (level)
```

Set RooFit log level

eskapade.root_analysis.style module

Module contents

eskapade.visualization package

Subpackages

eskapade.visualization.links package

Submodules

eskapade.visualization.links.correlation_summary module

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

Bases: `eskapade.core.run_elements.Link`

Create a heatmap of correlations between dataframe variables

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

Initialize CorrelationSummary 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 CorrelationSummary

finalize ()
Finalize CorrelationSummary

initialize ()
Initialize CorrelationSummary

eskapade.visualization.links.df_boxplot module

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

Bases: `eskapade.core.run_elements.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 the DfBoxplot link

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 DfBoxplot

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 DfBoxplot

initialize()
Initialize DfBoxplot link

eskapade.visualization.links.df_summary module

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

Bases: `eskapade.core.run_elements.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 the DfSummary link

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 DfSummary

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 DfSummary

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 DfSummary 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.DfSummary` (**kwargs)

Bases: `eskapade.core.run_elements.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 the DfSummary link

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 DfSummary

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 DfSummary

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 DfSummary 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

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

Bases: `eskapade.core.run_elements.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 the DfBoxplot link

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 DfBoxplot

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 DfBoxplot

initialize ()

Initialize DfBoxplot link

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

Bases: `eskapade.core.run_elements.Link`

Create a heatmap of correlations between dataframe variables

__init__ (**kwargs)

Initialize CorrelationSummary 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 CorrelationSummary

finalize ()

Finalize CorrelationSummary

initialize ()

Initialize CorrelationSummary

Submodules

`eskapade.visualization.vis_utils` module

`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 data frame

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.exceptions module

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

`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.mixins module

class `eskapade.mixins.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.mixins.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.mixins.LoggingMixin`

Bases: `object`

Mixin base class for logging

classmethod `log()`
Get logger of the module of this class

classmethod `set_log_level(level)`
Set logging level for the module of this class

class `eskapade.mixins.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.utils module

`eskapade.utils.build_cxx_library(lib_key='', accept_existing=False)`
Build Eskapade C++ library

Parameters

- **lib_key** (*str*) – key of the library to build (build all if empty)
- **accept_existing** (*bool*) – accept existing library if build fails

`eskapade.utils.collect_python_modules()`
Function to collect Eskapade Python modules

`eskapade.utils.get_dir_path(key)`

Function to retrieve Eskapade specific directory path

Parameters `key` (*str*) – Eskapade specific project directory key

Returns directory path

Return type `str`

`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.get_file_path(key)`

Function to retrieve Eskapade specific directory file path

Parameters `key` (*str*) – Eskapade specific project file key

Returns file path

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

Module contents

Appendices

MacOS

This Eskapade installation guide is written using **MacOS Sierra**.

To install Eskapade on MacOS there are basically four challenges to overcome certain versioning issues:

- Getting an isolated Python 3.5.2 environment,
- Getting ROOT 6.08/06 to work with Python 3.5.2 (ROOT conflicts with Anaconda),
- Getting Python packages similar to those in Anaconda,
- Getting Spark 2.1.0 set up (the homebrew version does not cooperate).

Therefore we are building a KaveToolBox for MacOS and re-use the Vagrant code.

Setup Python 3

First install a clean Python 3.5.2 (same as Anaconda):

```
env PYTHON_CONFIGURE_OPTS="--enable-framework" pyenv install -v 3.5.2
```

Note: Using homebrew will give Python 3.6.x which is incompatible with Spark 2.1.0

Then we need to create an isolated Python environment which we call 'eskapade':

```
brew install virtualenv
pyenv-sh-virtualenvwrapper
source /usr/local/bin/virtualenvwrapper.sh
mkvirtualenv -p .pyenv/versions/3.5.2/bin/python eskapade
workon eskapade
```

Installing ROOT 6

Clone ROOT from the git repository:

```
git clone http://root.cern.ch/git/root.git
cd root
git tag -l
git checkout -b v6-08-06 v6-08-06
mkdir ~/root
cd ~/root
```

Next, apply some patches to make it work properly with Python 3.5. Luckily that was already sorted out by the Eskapade team for the Vagrant/Docker installation:

```
for patchfile in $(ls ${ESKAPADE}/vagrant/dev/root/patches/*.patch); do \
  patch -p1 -i "${patchfile}" \
done
```

Next we compile it with some additional flags to ensure it has the desired functionality. Note the CXX stuff (this is important):

```
cmake /Users/${USER}/git/root/. \
-DPYTHON_EXECUTABLE=/Users/${USER}/.pyenv/versions/3.5.2/bin/python \
-DPYTHON_INCLUDE_DIR=/Users/${USER}/.pyenv/versions/3.5.2/include/python3.5m \
-DPYTHON_LIBRARY=/Users/${USER}/.pyenv/versions/3.5.2/lib/libpython3.5m.dylib \
-Dfail-on-missing=ON \
-Dcxx14=ON -Droot7=ON -Dshared=ON -Dsoversion=ON -Dthread=ON -Dfortran=ON -Dpython=ON \
↪ -Dcling=ON -Dx11=ON -Dssl=ON \
-Dxml=ON -Dfft3=ON -Dbuiltin_fft3=OFF -Dmathmore=ON -Dminuit2=ON -Droofit=ON -
↪ Dtmva=ON -Dopengl=ON -Dgviz=ON \
-Dalien=OFF -Dbonjour=OFF -Dcastor=OFF -Dchirp=OFF -Ddavix=OFF -Ddcache=OFF -
↪ Dfitsio=OFF -Dgfal=OFF -Dhdfs=OFF \
-Dkrb5=OFF -Dldap=OFF -Dmonalisa=OFF -Dmysql=OFF -Dodbc=OFF -Doracle=OFF -Dpgsql=OFF -
↪ Dpythia6=OFF -Dpythia8=OFF \
-Dsqlite=OFF -Drfio=OFF -Dxrootd=OFF \
cmake --build . -- -j4
```

When this is done, we need to setup the ROOTSYS environment, so we can install some additional Python packages:

```
source ~/root/bin/thisroot.sh
```

And finally install some Python packages for ROOT bindings:

```
pip install rootpy root-numpy root_pandas
```

Getting Python packages

Then we want to install the same packages as there are in a proper KaveToolBox environment to avoid version conflicts and random issues. We can use a requirements file, obtained through a `pip freeze` on the Vagrant/Docker installation that works.

To install the requirements run:

```
pip install -r requirements.txt
```

Input for the requirements.txt file is the following:

```
alabaster==0.7.8
amqp==2.1.4
appdirs==1.4.3
appnope==0.1.0
argcomplete==1.0.0
arrow==0.10.0
astroid==1.4.9
astropy==1.2.1
autopep8==1.3.1
Babel==2.3.3
backports.shutil-get-terminal-size==1.0.0
beautifulsoup4==4.4.1
billiard==3.5.0.2
binaryornot==0.4.0
bitarray==0.8.1
blaze==0.10.1
bokeh==0.12.0
boto==2.40.0
Bottleneck==1.1.0
branca==0.2.0
bson==0.4.6
cairocffi==0.8.0
CairoSVG==2.0.2
celery==4.0.2
cffi==1.6.0
chardet==2.3.0
cheroot==5.4.0
CherryPy==10.2.1
chest==0.2.3
click==6.6
cloudpickle==0.2.1
clyent==1.2.1
colorama==0.3.7
configobj==5.0.6
contextlib2==0.5.3
cookiecutter==1.5.1
coverage==4.3.4
cryptography==1.4
cssselect==1.0.1
```

```
cycler==0.10.0
Cython==0.24
cytoolz==0.8.0
dask==0.10.0
datashape==0.5.2
decorator==4.0.10
Delorean==0.6.0
descartes==1.1.0
dill==0.2.5
Django==1.10.5
django-filter==1.0.2
djangorestframework==3.6.2
docutils==0.12
entrypoints==0.2.2
et-xmlfile==1.0.1
fastcache==1.0.2
Flask==0.11.1
Flask-Cors==2.1.2
folium==0.3.0
future==0.16.0
gevent==1.1.1
gnureadline==6.3.3
greenlet==0.4.10
h5py==2.6.0
HeapDict==1.0.0
humanize==0.5.1
idna==2.1
imagesize==0.7.1
ipykernel==4.3.1
ipython==4.2.0
ipython-genutils==0.1.0
ipywidgets==4.1.1
isort==4.2.5
itsdangerous==0.24
JayDeBeApi==1.1.1
jdcal==1.2
jedi==0.9.0
Jinja2==2.8
jinja2-time==0.2.0
JPype1==0.6.2
jsonschema==2.5.1
jupyter==1.0.0
jupyter-client==4.3.0
jupyter-console==4.1.1
jupyter-core==4.1.0
kombu==4.0.2
lazy-object-proxy==1.2.2
loket==0.2.0
lxml==3.6.0
Markdown==2.6.8
MarkupSafe==0.23
matplotlib==1.5.1
mccabe==0.6.1
mistune==0.7.2
mock==2.0.0
modernize==0.5
mpld3==0.3
mpmath==0.19
```

```
multipledispatch==0.4.8
names==0.3.0
nbconvert==4.2.0
nbformat==4.0.1
nbpresent==3.0.0
networkx==1.11
nltk==3.2.1
nose==1.3.7
notebook==4.2.1
numpy==1.11.1
odo==0.5.0
packaging==16.8
pandas==0.18.1
patsy==0.4.1
pbr==2.0.0
pexpect==4.0.1
pickleshare==0.7.2
Pillow==3.2.0
portend==1.8
poyo==0.4.1
prompt-toolkit==1.0.14
psutil==4.3.0
ptyprocess==0.5.1
py4j==0.10.4
pyasn1==0.1.9
pycodestyle==2.3.1
pyparser==2.14
Pygments==2.1.3
pymongo==3.4.0
pyparsing==2.1.4
python-dateutil==2.5.3
pytz==2016.4
PyYAML==3.11
pyzmq==15.2.0
qtconsole==4.2.1
requests==2.13.0
root-numpy==4.7.2
root-pandas==0.1.1
rootpy==0.9.0
scikit-learn==0.18.1
scipy==0.19.0
seaborn==0.7.1
simplegeneric==0.8.1
six==1.10.0
sklearn==0.0
snowballstemmer==1.2.1
sortedcontainers==1.5.7
Sphinx==1.5.3
sphinx-rtd-theme==0.2.4
SQLAlchemy==1.0.13
statsmodels==0.8.0
tabulate==0.7.7
tempora==1.6.1
terminado==0.6
tinycss==0.4
toolz==0.8.0
tornado==4.3
traitlets==4.2.1
```

```
tzlocal==1.3
vine==1.1.3
wcmwidth==0.1.7
Werkzeug==0.11.10
whichcraft==0.4.0
wrapt==1.10.10
```

Setting up Spark 2.1.0

Now download Spark from apache, extract it, and compile it:

```
wget "http://archive.apache.org/dist/spark/spark-2.1.0/spark-2.1.0.tgz"
tar -xzf "spark-2.1.0.tgz"
cd spark-2.1.0
mvn -DskipTests clean package
```

Ensure it has the py4j package:

```
pip install py4j
```

Add docker containers to hosts

Add the following aliases to the localhost line in /etc/hosts, so it looks like:

```
127.0.0.1 localhost es-service es-mongo es-proxy
```

This will ensure you can reach the docker containers via the port forwards from the container to the docker host (i.e. localhost).

Cleaning the environment

Everytime you want to have a clean Eskapade environment run the following:

```
# --- setup Python
source /usr/local/bin/virtualenvwrapper.sh
workon eskapade

# --- setup ROOT
source ~/root/bin/thisroot.sh

# --- setup Spark
export SPARK_HOME=$HOME/spark-2.1.0
export PYTHONPATH=$SPARK_HOME/python/:$PYTHONPATH
export PYSARK_SUBMIT_ARGS="--master local[4] --num-executors 1 --executor-cores 4 --
↪executor-memory 4g pyspark-shell"

# --- setup Eskapade
cd ~/git/gitlab-nl/decision-engine
source ./eskapade/setup.sh
source ./analyticsengine/setup.sh
```

To automate this you can put it in a 'setup_eskapade.sh' script, but at the time of writing we have not done this yet.

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 real-time streaming (micro batch) processing with Apache Spark.

Requirements

A default 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](#).

NB: not all combinations of Spark and Python versions work properly together.

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 using three different methods:

1. environment variables
2. an Eskapade macro (preferred)
3. an Eskapade link

This is demonstrated in the following tutorial macro:

```
$ run_escaped tutorials/esk601_spark_configuration.py
```

The methods are described in the sections below. For a description of 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.

Environment variables

Configuration for Spark jobs can be set through the `PYSPARK_SUBMIT_ARGS` environment variable, e.g.:

```
PYSPARK_SUBMIT_ARGS=--master local[4] --num-executors 1 --executor-cores 4 --executor-
↪memory 4g pyspark-shell
```

The Spark session is then started directly from a macro with specified settings.

```
sm = proc_mgr.service(SparkManager)
sm.spark_session
sm.spark_context.setLogLevel('INFO')
```

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.

Configuration for Spark jobs can be set through the `SparkConf` class that holds a list of key/value pairs with configuration settings, e.g.:

```
conf = pyspark.conf.SparkConf()
conf.setAppName(settings['analysisName'])
conf.setAll([('spark.driver.host', '127.0.0.1'), ('spark.master', 'local[2]')])
```

The Spark session is then started directly from a macro with specified settings.

```
sm = proc_mgr.service(SparkManager)
sm.spark_conf = conf
sm.spark_session
sm.spark_context.setLogLevel('INFO')
```

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.:

```
conf_link = spark_analysis.SparkConfigurator(name='SparkConfigurator')
conf_link.sparkConf = [('spark.master', 'local[3]')]
conf_link.setLogLevel = 'INFO'
proc_mgr.add_chain('Config').add_link(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:

```
proc_mgr.service(SparkManager).spark_context.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.

Indices and tables

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

e

- eskapade, 3
- eskapade.analysis, 73
- eskapade.analysis.datetime, 60
- eskapade.analysis.histogram, 61
- eskapade.analysis.histogram_filling, 68
- eskapade.analysis.links, 51
- eskapade.analysis.links.apply_func_to_df, 41
- eskapade.analysis.links.apply_selection_to_df, 41
- eskapade.analysis.links.assign_random_class, 42
- eskapade.analysis.links.basic_generator, 42
- eskapade.analysis.links.df_concatenator, 43
- eskapade.analysis.links.df_merger, 43
- eskapade.analysis.links.histogrammar_filler, 44
- eskapade.analysis.links.random_sample_splitter, 45
- eskapade.analysis.links.read_to_df, 46
- eskapade.analysis.links.record_factorizer, 47
- eskapade.analysis.links.record_vectorizer, 48
- eskapade.analysis.links.value_counter, 49
- eskapade.analysis.links.write_from_df, 50
- eskapade.analysis.statistics, 71
- eskapade.core, 86
- eskapade.core.definitions, 74
- eskapade.core.execution, 75
- eskapade.core.persistence, 75
- eskapade.core.process_manager, 76
- eskapade.core.process_services, 80
- eskapade.core.run_elements, 83
- eskapade.core.run_utils, 86
- eskapade.core_ops, 97
- eskapade.core_ops.links, 92
- eskapade.core_ops.links.assert_in_ds, 86
- eskapade.core_ops.links.break_link, 87
- eskapade.core_ops.links.ds_object_deleter, 87
- eskapade.core_ops.links.ds_to_ds, 88
- eskapade.core_ops.links.event_looper, 88
- eskapade.core_ops.links.hello_world, 89
- eskapade.core_ops.links.ipython_embed, 89
- eskapade.core_ops.links.line_printer, 89
- eskapade.core_ops.links.print_ds, 90
- eskapade.core_ops.links.repeat_chain, 90
- eskapade.core_ops.links.skip_chain_if_empty, 91
- eskapade.core_ops.links.to_ds_dict, 91
- eskapade.data_quality, 102
- eskapade.data_quality.dq_helper, 101
- eskapade.data_quality.links, 99
- eskapade.data_quality.links.fix_pandas_dataframe, 97
- eskapade.exceptions, 129
- eskapade.helpers, 130
- eskapade.mixins, 131
- eskapade.root_analysis.decorators, 104
- eskapade.root_analysis.decorators.histograms, 102
- eskapade.root_analysis.decorators.roofit, 104
- eskapade.root_analysis.links.add_propagated_error_t, 104
- eskapade.root_analysis.links.convert_dataframe_2_r, 105
- eskapade.root_analysis.links.convert_roodataset_2_c

106
eskapade.root_analysis.links.convert_roodataset_2_roodatahist,
106
eskapade.root_analysis.links.convert_root_hist_2_roodatahist,
107
eskapade.root_analysis.links.convert_root_hist_2_roodataset,
108
eskapade.root_analysis.links.print_ws,
108
eskapade.root_analysis.links.roodatahist_filler,
109
eskapade.root_analysis.links.root_hist_filler,
110
eskapade.root_analysis.links.trunc_exp_fit,
111
eskapade.root_analysis.links.trunc_exp_gen,
112
eskapade.root_analysis.links.ws_utils,
113
eskapade.root_analysis.roofit_manager,
116
eskapade.root_analysis.roofit_models,
116
eskapade.root_analysis.roofit_utils, 119
eskapade.root_analysis.style, 119
eskapade.utils, 132
eskapade.version, 133
eskapade.visualization, 128
eskapade.visualization.links, 123
eskapade.visualization.links.correlation_summary,
119
eskapade.visualization.links.df_boxplot,
120
eskapade.visualization.links.df_summary,
121
eskapade.visualization.vis_utils, 126

Symbols

- `__init__()` (eskapade.analysis.datetime.FreqTimePeriod method), 60
`__init__()` (eskapade.analysis.datetime.TimePeriod method), 60
`__init__()` (eskapade.analysis.datetime.UniformTsTimePeriod method), 61
`__init__()` (eskapade.analysis.histogram.BinningUtil method), 61
`__init__()` (eskapade.analysis.histogram.Histogram method), 63
`__init__()` (eskapade.analysis.histogram.ValueCounts method), 66
`__init__()` (eskapade.analysis.histogram_filling.HistogramFillerBase method), 68
`__init__()` (eskapade.analysis.links.ApplyFuncToDf method), 51
`__init__()` (eskapade.analysis.links.ApplySelectionToDf method), 51
`__init__()` (eskapade.analysis.links.AssignRandomClass method), 55
`__init__()` (eskapade.analysis.links.BasicGenerator method), 52
`__init__()` (eskapade.analysis.links.DfConcatenator method), 52
`__init__()` (eskapade.analysis.links.DfMerger method), 53
`__init__()` (eskapade.analysis.links.HistogrammarFiller method), 59
`__init__()` (eskapade.analysis.links.RandomSampleSplitter method), 56
`__init__()` (eskapade.analysis.links.ReadToDf method), 54
`__init__()` (eskapade.analysis.links.RecordFactorizer method), 56
`__init__()` (eskapade.analysis.links.RecordVectorizer method), 54
`__init__()` (eskapade.analysis.links.ValueCounter method), 57
`__init__()` (eskapade.analysis.links.WriteFromDf method), 55
`__init__()` (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 41
`__init__()` (eskapade.analysis.links.apply_selection_to_df.ApplySelectionToDf method), 41
`__init__()` (eskapade.analysis.links.assign_random_class.AssignRandomClass method), 42
`__init__()` (eskapade.analysis.links.basic_generator.BasicGenerator method), 43
`__init__()` (eskapade.analysis.links.df_concatenator.DfConcatenator method), 43
`__init__()` (eskapade.analysis.links.df_merger.DfMerger method), 43
`__init__()` (eskapade.analysis.links.histogrammar_filler.HistogrammarFiller method), 44
`__init__()` (eskapade.analysis.links.random_sample_splitter.RandomSampleSplitter method), 45
`__init__()` (eskapade.analysis.links.read_to_df.ReadToDf method), 46
`__init__()` (eskapade.analysis.links.record_factorizer.RecordFactorizer method), 47
`__init__()` (eskapade.analysis.links.record_vectorizer.RecordVectorizer method), 48
`__init__()` (eskapade.analysis.links.value_counter.ValueCounter method), 49
`__init__()` (eskapade.analysis.links.write_from_df.WriteFromDf method), 50
`__init__()` (eskapade.analysis.statistics.ArrayStats method), 72
`__init__()` (eskapade.analysis.statistics.GroupByStats method), 73
`__init__()` (eskapade.core.definitions.RandomSeeds method), 74
`__init__()` (eskapade.core.persistence.IoConfig method), 75
`__init__()` (eskapade.core.process_manager.ProcessManager method), 77
`__init__()` (eskapade.core.process_services.ConfigObject method), 81

- `__init__()` (eskapade.core.process_services.ProcessService method), 82
- `__init__()` (eskapade.core.run_elements.Chain method), 83
- `__init__()` (eskapade.core.run_elements.Link method), 85
- `__init__()` (eskapade.core_ops.links.AssertInDs method), 92
- `__init__()` (eskapade.core_ops.links.Break method), 96
- `__init__()` (eskapade.core_ops.links.DsObjectDeleter method), 92
- `__init__()` (eskapade.core_ops.links.DsToDs method), 93
- `__init__()` (eskapade.core_ops.links.EventLooper method), 93
- `__init__()` (eskapade.core_ops.links.HelloWorld method), 94
- `__init__()` (eskapade.core_ops.links.IPythonEmbed method), 96
- `__init__()` (eskapade.core_ops.links.LinePrinter method), 94
- `__init__()` (eskapade.core_ops.links.PrintDs method), 94
- `__init__()` (eskapade.core_ops.links.RepeatChain method), 95
- `__init__()` (eskapade.core_ops.links.SkipChainIfEmpty method), 95
- `__init__()` (eskapade.core_ops.links.ToDsDict method), 96
- `__init__()` (eskapade.core_ops.links.assert_in_ds.AssertInDs method), 86
- `__init__()` (eskapade.core_ops.links.break_link.Break method), 87
- `__init__()` (eskapade.core_ops.links.ds_object_deleter.DsObjectDeleter method), 87
- `__init__()` (eskapade.core_ops.links.ds_to_ds.DsToDs method), 88
- `__init__()` (eskapade.core_ops.links.event_looper.EventLooper method), 88
- `__init__()` (eskapade.core_ops.links.hello_world.HelloWorld method), 89
- `__init__()` (eskapade.core_ops.links.ipython_embed.IPythonEmbed method), 89
- `__init__()` (eskapade.core_ops.links.line_printer.LinePrinter method), 89
- `__init__()` (eskapade.core_ops.links.print_ds.PrintDs method), 90
- `__init__()` (eskapade.core_ops.links.repeat_chain.RepeatChain method), 90
- `__init__()` (eskapade.core_ops.links.skip_chain_if_empty.SkipChainIfEmpty method), 91
- `__init__()` (eskapade.core_ops.links.to_ds_dict.ToDsDict method), 91
- `__init__()` (eskapade.data_quality.links.FixPandasDataFrame method), 99
- `__init__()` (eskapade.data_quality.links.fix_pandas_dataframe.FixPandasDataFrame method), 97
- `__init__()` (eskapade.exceptions.MissingPackageError method), 129
- `__init__()` (eskapade.exceptions.MissingPy4jError method), 129
- `__init__()` (eskapade.exceptions.MissingRooFitError method), 129
- `__init__()` (eskapade.exceptions.MissingRooStatsError method), 129
- `__init__()` (eskapade.exceptions.MissingRootError method), 129
- `__init__()` (eskapade.exceptions.MissingSparkError method), 130
- `__init__()` (eskapade.mixins.ConfigMixin method), 131
- `__init__()` (eskapade.mixins.TimerMixin method), 132
- `__init__()` (eskapade.root_analysis.links.add_propagated_error_to_roodataset method), 104
- `__init__()` (eskapade.root_analysis.links.convert_dataframe_2_roodataset.C method), 105
- `__init__()` (eskapade.root_analysis.links.convert_roodataset_2_dataframe.C method), 106
- `__init__()` (eskapade.root_analysis.links.convert_roodataset_2_roodatahist.C method), 107
- `__init__()` (eskapade.root_analysis.links.convert_root_hist_2_roodatahist.C method), 107
- `__init__()` (eskapade.root_analysis.links.convert_root_hist_2_roodataset.Co method), 108
- `__init__()` (eskapade.root_analysis.links.print_ws.PrintWs method), 108
- `__init__()` (eskapade.root_analysis.links.roodatahist_filler.RooDataHistFiller method), 109
- `__init__()` (eskapade.root_analysis.links.root_hist_filler.RootHistFiller method), 110
- `__init__()` (eskapade.root_analysis.links.trunc_exp_fit.TruncExpFit method), 112
- `__init__()` (eskapade.root_analysis.links.trunc_exp_gen.TruncExpGen method), 112
- `__init__()` (eskapade.root_analysis.links.ws_utils.WsUtils method), 113
- `__init__()` (eskapade.root_analysis.roofit_manager.RooFitManager method), 116
- `__init__()` (eskapade.root_analysis.roofit_models.LinearRegression method), 116
- `__init__()` (eskapade.root_analysis.roofit_models.RooFitModel method), 117
- `__init__()` (eskapade.root_analysis.roofit_models.TruncExponential method), 118
- `__init__()` (eskapade.visualization.links.CorrelationSummary method), 126
- `__init__()` (eskapade.visualization.links.DfBoxplot method), 125
- `__init__()` (eskapade.visualization.links.DfSummary method), 123
- `__init__()` (eskapade.visualization.links.correlation_summary.CorrelationSu method), 119

- `__init__()` (eskapade.visualization.links.df_boxplot.DfBoxplot method), 120
- `__init__()` (eskapade.visualization.links.df_summary.DfSummary method), 121
- ## A
- `add_chain()` (eskapade.core.process_manager.ProcessManager method), 77
- `add_fit()` (eskapade.root_analysis.links.ws_utils.WsUtils method), 113
- `add_link()` (eskapade.core.run_elements.Chain method), 83
- `add_macros()` (eskapade.core.process_services.ConfigObject method), 81
- `add_plot()` (eskapade.root_analysis.links.ws_utils.WsUtils method), 113
- `add_simulate()` (eskapade.root_analysis.links.ws_utils.WsUtils method), 114
- `addApplyFunc()` (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 41
- `addApplyFunc()` (eskapade.analysis.links.ApplyFuncToDf method), 51
- `AddPropagatedErrorToRooDataSet` (class in eskapade.root_analysis.links.add_propagated_error_to_root_dataset), 104
- `all_vars_set` (eskapade.root_analysis.roofit_models.TruncatedExponential attribute), 118
- `apply_transform_funcs()` (in module eskapade.helpers), 130
- `ApplyFuncToDf` (class in eskapade.analysis.links), 51
- `ApplyFuncToDf` (class in eskapade.analysis.links.apply_func_to_df), 41
- `ApplySelectionToDf` (class in eskapade.analysis.links), 51
- `ApplySelectionToDf` (class in eskapade.analysis.links.apply_selection_to_df), 41
- `ArgumentsMixin` (class in eskapade.mixins), 131
- `ArrayStats` (class in eskapade.analysis.statistics), 71
- `assert_data_type()` (eskapade.visualization.links.df_summary.DfSummary method), 121
- `assert_data_type()` (eskapade.visualization.links.DfSummary method), 123
- `assert_dataframe()` (eskapade.analysis.histogram_filling.HistogramFillerBase method), 69
- `AssertInDs` (class in eskapade.core_ops.links), 92
- `AssertInDs` (class in eskapade.core_ops.links.assert_in_ds), 86
- ## B
- `AssignRandomClass` (class in eskapade.analysis.links), 55
- `AssignRandomClass` (class in eskapade.analysis.links.assign_random_class), 42
- `BasicGenerator` (class in eskapade.analysis.links), 52
- `BasicGenerator` (class in eskapade.analysis.links.basic_generator), 42
- `bin_centers()` (eskapade.analysis.histogram.Histogram method), 63
- `bin_centers()` (in module eskapade.root_analysis.decorators.histograms), 102
- `bin_edges()` (eskapade.analysis.histogram.Histogram method), 63
- `bin_edges()` (in module eskapade.root_analysis.decorators.histograms), 103
- `bin_entries()` (eskapade.analysis.histogram.Histogram method), 63
- `bin_entries()` (in module eskapade.root_analysis.decorators.histograms), 103
- `bin_entries_2dgrid()` (in module eskapade.root_analysis.decorators.histograms), 103
- `bin_labels()` (eskapade.analysis.histogram.Histogram method), 63
- `bin_labels()` (in module eskapade.root_analysis.decorators.histograms), 103
- `bin_range()` (in module eskapade.root_analysis.decorators.histograms), 103
- `bin_specs` (eskapade.analysis.histogram.BinningUtil attribute), 61
- `bin_vals()` (in module eskapade.root_analysis.decorators.histograms), 103
- `BinningUtil` (class in eskapade.analysis.histogram), 61
- `bool_to_int()` (in module eskapade.data_quality.dq_helper), 101
- `bool_to_str()` (in module eskapade.data_quality.dq_helper), 101
- `box_plot()` (in module eskapade.visualization.vis_utils), 126
- `Break` (class in eskapade.core_ops.links), 96
- `Break` (class in eskapade.core_ops.links.break_link), 87
- `build_cxx_library()` (in module eskapade.utils), 132
- `build_model()` (eskapade.root_analysis.roofit_models.LinearRegression method), 116

- build_model() (eskapade.root_analysis.roofit_models.TruncExponential method), 118
- C**
- categorize_columns() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 69
- categorize_columns() (eskapade.root_analysis.links.root_hist_filler.RootHistFiller method), 111
- Chain (class in eskapade.core.run_elements), 83
- check_arg_callable() (eskapade.mixins.ArgumentsMixin method), 131
- check_arg_iters() (eskapade.mixins.ArgumentsMixin method), 131
- check_arg_opts() (eskapade.mixins.ArgumentsMixin method), 131
- check_arg_types() (eskapade.mixins.ArgumentsMixin method), 131
- check_arg_vals() (eskapade.mixins.ArgumentsMixin method), 131
- check_extra_kwargs() (eskapade.mixins.ArgumentsMixin method), 131
- check_io_config() (eskapade.core.process_manager.ProcessManager static method), 77
- check_nan() (in module eskapade.data_quality.dq_helper), 101
- check_required_args() (eskapade.mixins.ArgumentsMixin method), 131
- checkCollectionSet() (eskapade.core_ops.links.skip_chain_if_empty.SkipChainIfEmpty method), 91
- checkCollectionSet() (eskapade.core_ops.links.SkipChainIfEmpty method), 95
- cleanup_string() (in module eskapade.data_quality.dq_helper), 101
- coef_ (eskapade.root_analysis.roofit_models.LinearRegression fit attribute), 117
- coll_iter() (in module eskapade.root_analysis.decorators.roofit), 104
- collect_python_modules() (in module eskapade.utils), 132
- combine_hists() (eskapade.analysis.histogram.Histogram class method), 64
- config_path (eskapade.mixins.ConfigMixin attribute), 131
- ConfigMixin (class in eskapade.mixins), 131
- ConfigObject (class in eskapade.core.process_services), 80
- construct_empty_hist() (eskapade.analysis.links.histogrammar_filler.HistogrammarFiller method), 45
- construct_empty_hist() (eskapade.analysis.links.HistogrammarFiller method), 60
- construct_empty_hist() (eskapade.root_analysis.links.root_hist_filler.RootHistFiller method), 111
- convert() (in module eskapade.data_quality.dq_helper), 101
- ConvertDataFrame2RooDataSet (class in eskapade.root_analysis.links.convert_dataframe_2_roodataset), 105
- ConvertRooDataSet2DataFrame (class in eskapade.root_analysis.links.convert_roodataset_2_dataframe), 106
- ConvertRooDataSet2RooDataHist (class in eskapade.root_analysis.links.convert_roodataset_2_roodatahist), 106
- ConvertRootHist2RooDataHist (class in eskapade.root_analysis.links.convert_root_hist_2_roodatahist), 107
- ConvertRootHist2RooDataSet (class in eskapade.root_analysis.links.convert_root_hist_2_roodataset), 108
- copy() (eskapade.analysis.histogram.Histogram method), 64
- CorrelationSummary (class in eskapade.visualization.links), 126
- CorrelationSummary (class in eskapade.visualization.links.correlation_summary), 119
- count() (eskapade.analysis.histogram.ValueCounts method), 67
- counts (eskapade.analysis.histogram.ValueCounts attribute), 67
- create() (eskapade.core.process_services.ProcessService class method), 82
- create_arg_parser() (in module eskapade.core.run_utils), 86
- create_dir() (in module eskapade.core.persistence), 75
- create_mpv_stat() (eskapade.analysis.statistics.ArrayStats method), 72
- create_norm() (eskapade.root_analysis.roofit_models.TruncExponential method), 118
- create_roofit_opts() (in module eskapade.root_analysis.roofit_utils), 119
- create_stats() (eskapade.analysis.statistics.ArrayStats method), 72
- create_sub_counts() (eskapade.analysis.histogram.ValueCounts method), 67
- D**
- data_set_iter() (in module eskapade

- pade.root_analysis.decorators.roofit), 104
- DataStore (class in eskapade.core.process_services), 81
- datatype (eskapade.analysis.histogram.Histogram attribute), 64
- datatype (in module eskapade.root_analysis.decorators.histograms), 103
- delete_smallstat() (in module eskapade.visualization.vis_utils), 127
- determine_preferred_dtype() (in module eskapade.data_quality.links.fix_pandas_dataframe), 99
- DfBoxplot (class in eskapade.visualization.links), 125
- DfBoxplot (class in eskapade.visualization.links.df_boxplot), 120
- DfConcatenator (class in eskapade.analysis.links), 52
- DfConcatenator (class in eskapade.analysis.links.df_concatenator), 43
- DfMerger (class in eskapade.analysis.links), 53
- DfMerger (class in eskapade.analysis.links.df_merger), 43
- DfSummary (class in eskapade.visualization.links), 123
- DfSummary (class in eskapade.visualization.links.df_summary), 121
- do_fit() (eskapade.root_analysis.links.ws_utils.WsUtils method), 114
- do_plot() (eskapade.root_analysis.links.ws_utils.WsUtils method), 114
- do_simulate() (eskapade.root_analysis.links.ws_utils.WsUtils method), 115
- do_storage() (eskapade.root_analysis.links.roodatahist_filler.RooDataHistFiller method), 110
- dostorage() (eskapade.core_ops.links.to_ds_dict.ToDsDict method), 92
- dostorage() (eskapade.core_ops.links.ToDsDict method), 96
- drop_inconsistent_keys() (eskapade.analysis.links.value_counter.ValueCounter method), 50
- drop_inconsistent_keys() (eskapade.analysis.links.ValueCounter method), 58
- drop_requested_keys() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 69
- DsObjectDeleter (class in eskapade.core_ops.links), 92
- DsObjectDeleter (class in eskapade.core_ops.links.ds_object_deleter), 87
- DsToDs (class in eskapade.core_ops.links), 93
- DsToDs (class in eskapade.core_ops.links.ds_to_ds), 88
- dt_string() (eskapade.analysis.datetime.FreqTimePeriod method), 60
- ## E
- eskapade (module), 1, 133
- eskapade.analysis (module), 73
- eskapade.analysis.datetime (module), 60
- eskapade.analysis.histogram (module), 61
- eskapade.analysis.histogram_filling (module), 68
- eskapade.analysis.links (module), 51
- eskapade.analysis.links.apply_func_to_df (module), 41
- eskapade.analysis.links.apply_selection_to_df (module), 41
- eskapade.analysis.links.assign_random_class (module), 42
- eskapade.analysis.links.basic_generator (module), 42
- eskapade.analysis.links.df_concatenator (module), 43
- eskapade.analysis.links.df_merger (module), 43
- eskapade.analysis.links.histogrammer_filler (module), 44
- eskapade.analysis.links.random_sample_splitter (module), 45
- eskapade.analysis.links.read_to_df (module), 46
- eskapade.analysis.links.record_factorizer (module), 47
- eskapade.analysis.links.record_vectorizer (module), 48
- eskapade.analysis.links.value_counter (module), 49
- eskapade.analysis.links.write_from_df (module), 50
- eskapade.analysis.statistics (module), 71
- eskapade.core (module), 86
- eskapade.core.definitions (module), 74
- eskapade.core.execution (module), 75
- eskapade.core.persistence (module), 75
- eskapade.core.process_manager (module), 76
- eskapade.core.process_services (module), 80
- eskapade.core.run_elements (module), 83
- eskapade.core.run_utils (module), 86
- eskapade.core_ops (module), 97
- eskapade.core_ops.links (module), 92
- eskapade.core_ops.links.assert_in_ds (module), 86
- eskapade.core_ops.links.break_link (module), 87
- eskapade.core_ops.links.ds_object_deleter (module), 87
- eskapade.core_ops.links.ds_to_ds (module), 88
- eskapade.core_ops.links.event_looper (module), 88
- eskapade.core_ops.links.hello_world (module), 89
- eskapade.core_ops.links.ipython_embed (module), 89
- eskapade.core_ops.links.line_printer (module), 89
- eskapade.core_ops.links.print_ds (module), 90
- eskapade.core_ops.links.repeat_chain (module), 90
- eskapade.core_ops.links.skip_chain_if_empty (module), 91
- eskapade.core_ops.links.to_ds_dict (module), 91
- eskapade.data_quality (module), 102
- eskapade.data_quality.dq_helper (module), 101
- eskapade.data_quality.links (module), 99
- eskapade.data_quality.links.fix_pandas_dataframe (module), 97
- eskapade.exceptions (module), 129
- eskapade.helpers (module), 130

- eskapade.mixins (module), 131
- eskapade.root_analysis.decorators (module), 104
- eskapade.root_analysis.decorators.histograms (module), 102
- eskapade.root_analysis.decorators.roofit (module), 104
- eskapade.root_analysis.links.add_propagated_error_to_roodataset (module), 104
- eskapade.root_analysis.links.convert_dataframe_2_roodataset (module), 105
- eskapade.root_analysis.links.convert_roodataset_2_dataframe (module), 106
- eskapade.root_analysis.links.convert_roodataset_2_roodatahist (module), 106
- eskapade.root_analysis.links.convert_root_hist_2_roodatahist (module), 107
- eskapade.root_analysis.links.convert_root_hist_2_roodataset (module), 108
- eskapade.root_analysis.links.print_ws (module), 108
- eskapade.root_analysis.links.roodatahist_filler (module), 109
- eskapade.root_analysis.links.root_hist_filler (module), 110
- eskapade.root_analysis.links.trunc_exp_fit (module), 111
- eskapade.root_analysis.links.trunc_exp_gen (module), 112
- eskapade.root_analysis.links.ws_utils (module), 113
- eskapade.root_analysis.roofit_manager (module), 116
- eskapade.root_analysis.roofit_models (module), 116
- eskapade.root_analysis.roofit_utils (module), 119
- eskapade.root_analysis.style (module), 119
- eskapade.utils (module), 132
- eskapade.version (module), 133
- eskapade.visualization (module), 128
- eskapade.visualization.links (module), 123
- eskapade.visualization.links.correlation_summary (module), 119
- eskapade.visualization.links.df_boxplot (module), 120
- eskapade.visualization.links.df_summary (module), 121
- eskapade.visualization.vis_utils (module), 126
- est_norm_ratio() (in module eskapade.root_analysis.links.trunc_exp_fit), 112
- EventLooper (class in eskapade.core_ops.links), 93
- EventLooper (class in eskapade.core_ops.links.event_looper), 88
- execute() (eskapade.analysis.histogram_filling.HistogramFilling method), 69
- execute() (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 41
- execute() (eskapade.analysis.links.apply_selection_to_df.ApplySelectionToDf method), 42
- execute() (eskapade.analysis.links.ApplyFuncToDf method), 51
- execute() (eskapade.analysis.links.ApplySelectionToDf method), 52
- execute() (eskapade.analysis.links.assign_random_class.AssignRandomClass method), 42
- execute() (eskapade.analysis.links.AssignRandomClass method), 56
- execute() (eskapade.analysis.links.basic_generator.BasicGenerator method), 43
- execute() (eskapade.analysis.links.BasicGenerator method), 52
- execute() (eskapade.analysis.links.df_concatenator.DfConcatenator method), 43
- execute() (eskapade.analysis.links.df_merger.DfMerger method), 44
- execute() (eskapade.analysis.links.DfConcatenator method), 53
- execute() (eskapade.analysis.links.DfMerger method), 53
- execute() (eskapade.analysis.links.random_sample_splitter.RandomSampleSplitter method), 46
- execute() (eskapade.analysis.links.RandomSampleSplitter method), 56
- execute() (eskapade.analysis.links.read_to_df.ReadToDf method), 46
- execute() (eskapade.analysis.links.ReadToDf method), 54
- execute() (eskapade.analysis.links.record_factorizer.RecordFactorizer method), 47
- execute() (eskapade.analysis.links.record_vectorizer.RecordVectorizer method), 48
- execute() (eskapade.analysis.links.RecordFactorizer method), 57
- execute() (eskapade.analysis.links.RecordVectorizer method), 55
- execute() (eskapade.analysis.links.write_from_df.WriteFromDf method), 51
- execute() (eskapade.analysis.links.WriteFromDf method), 55
- execute() (eskapade.core.process_manager.ProcessManager method), 77
- execute() (eskapade.core.run_elements.Chain method), 83
- execute() (eskapade.core.run_elements.Link method), 85
- execute() (eskapade.core_ops.links.assert_in_ds.AssertInDs method), 87
- execute() (eskapade.core_ops.links.AssertInDs method), 92
- execute() (eskapade.core_ops.links.Break method), 96
- execute() (eskapade.core_ops.links.break_link.Break method), 87
- execute() (eskapade.core_ops.links.ds_object_deleter.DsObjectDeleter method), 87
- execute() (eskapade.core_ops.links.ds_to_ds.DsToDs method), 88
- execute() (eskapade.core_ops.links.DsObjectDeleter method), 92
- execute() (eskapade.core_ops.links.DsToDs method), 93
- execute() (eskapade.core_ops.links.event_looper.EventLooper

- method), 88
- execute() (eskapade.core_ops.links.EventLooper method), 93
- execute() (eskapade.core_ops.links.hello_world.HelloWorld method), 89
- execute() (eskapade.core_ops.links.HelloWorld method), 94
- execute() (eskapade.core_ops.links.ipython_embed.IPythonEmbed method), 89
- execute() (eskapade.core_ops.links.IPythonEmbed method), 96
- execute() (eskapade.core_ops.links.line_printer.LinePrinter method), 90
- execute() (eskapade.core_ops.links.LinePrinter method), 94
- execute() (eskapade.core_ops.links.print_ds.PrintDs method), 90
- execute() (eskapade.core_ops.links.PrintDs method), 94
- execute() (eskapade.core_ops.links.repeat_chain.RepeatChain method), 90
- execute() (eskapade.core_ops.links.RepeatChain method), 95
- execute() (eskapade.core_ops.links.skip_chain_if_empty.SkipChainIfEmpty method), 91
- execute() (eskapade.core_ops.links.SkipChainIfEmpty method), 95
- execute() (eskapade.core_ops.links.to_ds_dict.ToDsDict method), 92
- execute() (eskapade.core_ops.links.ToDsDict method), 96
- execute() (eskapade.data_quality.links.fix_pandas_dataframe.FixPandasDataFrame method), 98
- execute() (eskapade.data_quality.links.FixPandasDataFrame method), 100
- execute() (eskapade.root_analysis.links.add_propagated_error_to_roodataset.AddPropagatedErrorToRooDataSet method), 105
- execute() (eskapade.root_analysis.links.convert_dataframe_2_roodataset.ConvertDataFrame2RooDataSet method), 106
- execute() (eskapade.root_analysis.links.convert_roodataset_2_dataframe.ConvertRooDataSet2DataFrame method), 106
- execute() (eskapade.root_analysis.links.convert_roodataset_2_roodatahist.ConvertRooDataSet2RooDataHist method), 107
- execute() (eskapade.root_analysis.links.convert_root_hist_2_roodatahist.ConvertRootHist2RooDataHist method), 107
- execute() (eskapade.root_analysis.links.convert_root_hist_2_roodataset.ConvertRootHist2RooDataSet method), 108
- execute() (eskapade.root_analysis.links.print_ws.PrintWs method), 108
- execute() (eskapade.root_analysis.links.roodatahist_filler.RooDataHistFiller method), 110
- execute() (eskapade.root_analysis.links.trunc_exp_fit.TruncExpFit method), 112
- execute() (eskapade.root_analysis.links.trunc_exp_gen.TruncExpGen method), 112
- execute() (eskapade.root_analysis.links.ws_utils.WsUtils method), 115
- execute() (eskapade.visualization.links.correlation_summary.CorrelationSummary method), 120
- execute() (eskapade.visualization.links.CorrelationSummary method), 126
- execute() (eskapade.visualization.links.df_boxplot.DfBoxplot method), 120
- execute() (eskapade.visualization.links.df_summary.DfSummary method), 121
- execute() (eskapade.visualization.links.DfBoxplot method), 126
- execute() (eskapade.visualization.links.DfSummary method), 124
- execute_all() (eskapade.core.process_manager.ProcessManager method), 78
- execute_link() (eskapade.core.run_elements.Link method), 85
- execute_macro() (eskapade.core.process_manager.ProcessManager method), 78

F

- Failure (eskapade.core.definitions.StatusCode attribute), 74
- fill_histogram() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 70
- fill_histogram() (eskapade.analysis.links.histogrammar_filler.HistogrammarFiller method), 45
- fill_histogram() (eskapade.analysis.links.HistogrammarFiller method), 90
- fill_histogram() (eskapade.analysis.links.value_counter.ValueCounter method), 50
- fill_histogram() (eskapade.analysis.links.ValueCounter method), 50
- fill_histogram() (eskapade.root_analysis.links.root_hist_filler.RootHistFiller method), 111
- finalize() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 70
- finalize() (eskapade.core.process_manager.ProcessManager method), 93
- finalize() (eskapade.core.run_elements.Chain method), 85
- finalize() (eskapade.core.run_elements.Link method), 85
- finalize() (eskapade.core_ops.links.EventLooper method), 89
- finalize() (eskapade.core_ops.links.EventLooper method), 94
- finalize() (eskapade.root_analysis.links.roodatahist_filler.RooDataHistFiller method), 110
- finalize() (eskapade.root_analysis.links.ws_utils.WsUtils method), 115
- finalize() (eskapade.visualization.links.correlation_summary.CorrelationSummary method), 120

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

finalize() (eskapade.visualization.links.df_boxplot.DfBoxplot method), 120

finalize() (eskapade.visualization.links.df_summary.DfSummary method), 122

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

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

finalize_link() (eskapade.core.run_elements.Link method), 85

finish() (eskapade.core.process_services.ProcessService method), 82

fit() (eskapade.root_analysis.roofit_models.LinearRegression method), 117

fit_dataset() (eskapade.root_analysis.roofit_models.LinearRegression method), 117

FixPandasDataFrame (class in eskapade.data_quality.links), 99

FixPandasDataFrame (class in eskapade.data_quality.links.fix_pandas_dataframe), 97

freq (eskapade.analysis.datetime.FreqTimePeriod attribute), 60

FreqTimePeriod (class in eskapade.analysis.datetime), 60

G

gen_max_var_data() (in module eskapade.root_analysis.links.trunc_exp_gen), 112

get_all_columns() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 70

get_all_columns() (eskapade.visualization.links.df_summary.DfSummary method), 122

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

get_bin_center() (eskapade.analysis.histogram.BinningUtil method), 62

get_bin_count() (eskapade.analysis.histogram.Histogram method), 64

get_bin_edges() (eskapade.analysis.histogram.BinningUtil method), 62

get_bin_edges_range() (eskapade.analysis.histogram.BinningUtil method), 62

get_bin_labels() (eskapade.analysis.histogram.Histogram method), 64

get_bin_range() (eskapade.analysis.histogram.Histogram method), 64

get_bin_vals() (eskapade.analysis.histogram.Histogram method), 64

get_chain() (eskapade.core.process_manager.ProcessManager method), 78

get_chain_idx() (eskapade.core.process_manager.ProcessManager method), 78

get_col_props() (eskapade.analysis.statistics.ArrayStats method), 72

get_col_props() (in module eskapade.analysis.statistics), 73

get_config() (eskapade.mixins.ConfigMixin method), 131

get_data_type() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 70

get_dir_path() (in module eskapade.utils), 132

get_env_var() (in module eskapade.utils), 133

get_file_path() (in module eskapade.utils), 133

get_histogram() (eskapade.analysis.histogram.Histogram method), 65

get_latex_table() (eskapade.analysis.statistics.ArrayStats method), 72

get_latex_table() (eskapade.analysis.statistics.GroupByStats method), 73

get_left_bin_edge() (eskapade.analysis.histogram.BinningUtil method), 62

get_length() (eskapade.visualization.links.df_summary.DfSummary method), 122

get_length() (eskapade.visualization.links.DfSummary method), 124

get_link() (eskapade.core.run_elements.Chain method), 83

get_nonone_bin_centers() (eskapade.analysis.histogram.Histogram method), 65

get_nonone_bin_counts() (eskapade.analysis.histogram.Histogram method), 65

get_nonone_bin_edges() (eskapade.analysis.histogram.Histogram method), 65

get_nonone_bin_range() (eskapade.analysis.histogram.Histogram method), 65

get_print_stats() (eskapade.analysis.statistics.ArrayStats method), 72

get_right_bin_edge() (eskapade.analysis.histogram.BinningUtil method), 62

get_sample() (eskapade.visualization.links.df_summary.DfSummary method), 122

get_sample() (eskapade.visualization.links.DfSummary method), 124

get_service_tree() (eska-

pade.core.process_manager.ProcessManager initialize() (eskapade.analysis.links.assign_random_class.AssignRandomClass method), 78 method), 42
 get_services() (eskapade.core.process_manager.ProcessManager initialize() (eskapade.analysis.links.AssignRandomClass method), 79 method), 56
 get_uniform_bin_edges() (eskapade.analysis.histogram.Histogram initialize() (eskapade.analysis.links.basic_generator.BasicGenerator method), 65 method), 43
 get_values() (eskapade.analysis.histogram.ValueCounts initialize() (eskapade.analysis.links.BasicGenerator method), 67 method), 52
 get_x_label() (eskapade.analysis.statistics.ArrayStats initialize() (eskapade.analysis.links.df_concatenator.DfConcatenator method), 72 method), 43
 groupbyapply() (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf initialize() (eskapade.analysis.links.DfMerger method), 41 method), 44
 groupbyapply() (eskapade.analysis.links.ApplyFuncToDf initialize() (eskapade.analysis.links.DfConcatenator method), 51 method), 53
 GroupByStats (class in eskapade.analysis.statistics), 73 initialize() (eskapade.analysis.links.DfMerger method), 53
 initialize() (eskapade.analysis.links.random_sample_splitter.RandomSampleSplitter method), 46
 initialize() (eskapade.analysis.links.RandomSampleSplitter method), 56
 has_chain() (eskapade.core.process_manager.ProcessManager initialize() (eskapade.analysis.links.read_to_df.ReadToDf method), 79 method), 46
 has_link() (eskapade.core.run_elements.Chain method), 84 initialize() (eskapade.analysis.links.ReadToDf method), 54
 HelloWorld (class in eskapade.core_ops.links), 94 initialize() (eskapade.analysis.links.record_factorizer.RecordFactorizer method), 47
 HelloWorld (class in eskapade.core_ops.links.hello_world), 89 initialize() (eskapade.analysis.links.record_vectorizer.RecordVectorizer method), 103 method), 48
 high (in module eskapade.root_analysis.decorators.histograms), 103 initialize() (eskapade.analysis.links.RecordFactorizer method), 57
 Histogram (class in eskapade.analysis.histogram), 62 initialize() (eskapade.analysis.links.RecordVectorizer method), 55
 HistogramFillerBase (class in eskapade.analysis.histogram_filling), 68 initialize() (eskapade.analysis.links.value_counter.ValueCounter method), 50
 HistogrammarFiller (class in eskapade.analysis.links), 59 initialize() (eskapade.analysis.links.ValueCounter method), 58
 HistogrammarFiller (class in eskapade.analysis.links.histogrammar_filler), 44 initialize() (eskapade.analysis.links.write_from_df.WriteFromDf method), 51
 initialize() (eskapade.analysis.links.WriteFromDf method), 55
 initialize() (eskapade.core.process_manager.ProcessManager method), 79
 initialize() (eskapade.core.run_elements.Chain method), 84
 initialize() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 70 initialize() (eskapade.core.run_elements.Link method), 85
 initialize() (eskapade.analysis.links.apply_func_to_df.ApplyFuncToDf method), 41 initialize() (eskapade.core_ops.links.ds_object_deleter.DsObjectDeleter method), 87
 initialize() (eskapade.analysis.links.apply_selection_to_df.ApplySelectionToDf method), 42 initialize() (eskapade.core_ops.links.ds_to_ds.DsToDs method), 88
 initialize() (eskapade.analysis.links.ApplyFuncToDf method), 51 initialize() (eskapade.core_ops.links.DsObjectDeleter method), 92
 initialize() (eskapade.analysis.links.ApplySelectionToDf method), 52 initialize() (eskapade.core_ops.links.DsToDs method), 93
 initialize() (eskapade.core_ops.links.event_looper.EventLooper method), 93

- load_data() (eskapade.root_analysis.roofit_models.LinearRegression method), 117
- load_libesroofit() (in module eskapade.root_analysis.roofit_utils), 119
- log() (eskapade.mixins.LoggingMixin class method), 132
- LoggingMixin (class in eskapade.mixins), 132
- low (in module eskapade.root_analysis.decorators.histograms), 103
- ## M
- make_histogram() (eskapade.analysis.statistics.ArrayStats method), 72
- make_plots() (in module eskapade.root_analysis.links.trunc_exp_fit), 112
- max_var (eskapade.root_analysis.roofit_models.TruncExponential attribute), 118
- max_var_set (eskapade.root_analysis.roofit_models.TruncExponential attribute), 118
- MissingPackageError, 129
- MissingPy4jError, 129
- MissingRooFitError, 129
- MissingRooStatsError, 129
- MissingRootError, 129
- MissingSparkError, 130
- model() (eskapade.root_analysis.roofit_manager.RooFitManager method), 116
- ## N
- n_bins (eskapade.analysis.histogram.Histogram attribute), 65
- n_bins (in module eskapade.root_analysis.decorators.histograms), 104
- n_dim (eskapade.analysis.histogram.Histogram attribute), 65
- n_dim (in module eskapade.root_analysis.decorators.histograms), 104
- name (eskapade.core.run_elements.Link attribute), 86
- name (eskapade.root_analysis.roofit_models.RooFitModel attribute), 117
- nononecounts (eskapade.analysis.histogram.ValueCounts attribute), 67
- num_bins (eskapade.analysis.histogram.Histogram attribute), 65
- num_bins (eskapade.analysis.histogram.ValueCounts attribute), 67
- num_nonone_bins (eskapade.analysis.histogram.ValueCounts attribute), 68
- obj_repr() (in module eskapade.helpers), 130
- offset (eskapade.analysis.datetime.UniformTsTimePeriod attribute), 61
- only_bool() (in module eskapade.analysis.histogram_filling), 70
- only_float() (in module eskapade.analysis.histogram_filling), 70
- only_int() (in module eskapade.analysis.histogram_filling), 70
- only_str() (in module eskapade.analysis.histogram_filling), 71
- ## P
- pandasReader() (in module eskapade.analysis.links.read_to_df), 47
- pandasWriter() (in module eskapade.analysis.links.write_from_df), 51
- parse_date_time() (eskapade.analysis.datetime.TimePeriod class method), 60
- parse_time_period() (eskapade.analysis.datetime.TimePeriod class method), 60
- pdf (eskapade.root_analysis.roofit_models.RooFitModel attribute), 117
- pdf_name (eskapade.root_analysis.roofit_models.RooFitModel attribute), 117
- period (eskapade.analysis.datetime.UniformTsTimePeriod attribute), 61
- period_index() (eskapade.analysis.datetime.FreqTimePeriod method), 60
- period_index() (eskapade.analysis.datetime.TimePeriod method), 61
- period_index() (eskapade.analysis.datetime.UniformTsTimePeriod method), 61
- persist (eskapade.core.process_services.ProcessServiceMeta attribute), 82
- persist_in_file() (eskapade.core.process_services.ProcessService method), 82
- persist_services() (eskapade.core.process_manager.ProcessManager method), 79
- plot_2d_histogram() (in module eskapade.visualization.vis_utils), 127
- plot_correlation_matrix() (in module eskapade.visualization.vis_utils), 128
- plot_histogram() (in module eskapade.visualization.vis_utils), 128
- predict() (eskapade.root_analysis.roofit_models.LinearRegression method), 117
- Print() (eskapade.core.process_manager.ProcessManager method), 77

Print() (eskapade.core.process_services.ConfigObject method), 81

Print() (eskapade.core.process_services.DataStore method), 82

print_chains() (eskapade.core.process_manager.ProcessManager method), 79

print_services() (eskapade.core.process_manager.ProcessManager method), 79

PrintDs (class in eskapade.core_ops.links), 94

PrintDs (class in eskapade.core_ops.links.print_ds), 90

PrintWs (class in eskapade.root_analysis.links.print_ws), 108

process_1d_histogram() (eskapade.visualization.links.df_summary.DfSummary method), 122

process_1d_histogram() (eskapade.visualization.links.DfSummary method), 124

process_2d_histogram() (eskapade.visualization.links.df_summary.DfSummary method), 122

process_2d_histogram() (eskapade.visualization.links.DfSummary method), 124

process_and_store() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 70

process_and_store() (eskapade.analysis.links.value_counter.ValueCounter method), 50

process_and_store() (eskapade.analysis.links.ValueCounter method), 58

process_columns() (eskapade.analysis.histogram_filling.HistogramFillerBase method), 70

process_columns() (eskapade.analysis.links.value_counter.ValueCounter method), 50

process_columns() (eskapade.analysis.links.ValueCounter method), 58

process_counts() (eskapade.analysis.histogram.ValueCounts method), 68

process_nan_histogram() (eskapade.visualization.links.df_summary.DfSummary method), 122

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

process_sample() (eskapade.visualization.links.df_summary.DfSummary method), 123

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

process_series() (eskapade.visualization.links.df_summary.DfSummary method), 123

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

process_transform_funcs() (in module eskapade.helpers), 130

ProcessManager (class in eskapade.core.process_manager), 76

ProcessService (class in eskapade.core.process_services), 82

ProcessServiceMeta (class in eskapade.core.process_services), 82

R

RandomSampleSplitter (class in eskapade.analysis.links), 56

RandomSampleSplitter (class in eskapade.analysis.links.random_sample_splitter), 45

RandomSeeds (class in eskapade.core.definitions), 74

ReadToDf (class in eskapade.analysis.links), 53

ReadToDf (class in eskapade.analysis.links.read_to_df), 46

record_file_number() (in module eskapade.core.persistence), 76

record_vectorizer() (in module eskapade.analysis.links.record_vectorizer), 48

RecordFactorizer (class in eskapade.analysis.links), 56

RecordFactorizer (class in eskapade.analysis.links.record_factorizer), 47

RecordVectorizer (class in eskapade.analysis.links), 54

RecordVectorizer (class in eskapade.analysis.links.record_vectorizer), 48

Recoverable (eskapade.core.definitions.StatusCode attribute), 74

remove_all_services() (eskapade.core.process_manager.ProcessManager method), 79

remove_chain() (eskapade.core.process_manager.ProcessManager method), 79

remove_chains() (eskapade.core.process_manager.ProcessManager method), 80

remove_keys_of_inconsistent_type() (eskapade.analysis.histogram.Histogram method), 66

remove_keys_of_inconsistent_type() (eskapade.analysis.histogram.ValueCounts method), 68

remove_link() (eskapade.core.run_elements.Chain method), 84

- remove_links() (eskapade.core.run_elements.Chain method), 84
- remove_service() (eskapade.core.process_manager.ProcessManager method), 80
- RepeatChain (class in eskapade.core_ops.links), 94
- RepeatChain (class in eskapade.core_ops.links.repeat_chain), 90
- RepeatChain (eskapade.core.definitions.StatusCode attribute), 74
- repl_whites() (in module eskapade.core.persistence), 76
- reset() (eskapade.core.process_manager.ProcessManager method), 80
- reset_config() (eskapade.mixins.ConfigMixin method), 131
- reset_eskapade() (in module eskapade.core.execution), 75
- RooDataHistFiller (class in eskapade.root_analysis.links.roodatahist_filler), 109
- RooFitManager (class in eskapade.root_analysis.roofit_manager), 116
- RooFitModel (class in eskapade.root_analysis.roofit_models), 117
- RootHistFiller (class in eskapade.root_analysis.links.root_hist_filler), 110
- run_eskapade() (in module eskapade.core.execution), 75
- ## S
- score() (eskapade.root_analysis.roofit_models.LinearRegression method), 117
- sel_max_var_data() (in module eskapade.root_analysis.links.trunc_exp_gen), 113
- service() (eskapade.core.process_manager.ProcessManager method), 80
- set_begin_end_chain_opt() (in module eskapade.core.definitions), 75
- set_custom_user_vars() (in module eskapade.core.definitions), 75
- set_log_level() (eskapade.mixins.LoggingMixin class method), 132
- set_log_level_opt() (in module eskapade.core.definitions), 75
- set_matplotlib_backend() (in module eskapade.utils), 133
- set_opt_var() (in module eskapade.core.definitions), 75
- set_rf_log_level() (in module eskapade.root_analysis.roofit_utils), 119
- set_seeds() (in module eskapade.core.definitions), 75
- set_single_chain_opt() (in module eskapade.core.definitions), 75
- set_user_opts() (eskapade.core.process_services.ConfigObject method), 81
- set_var_vals() (eskapade.root_analysis.roofit_manager.RooFitManager method), 116
- setChunkSize() (eskapade.analysis.links.read_to_df.ReadToDf method), 47
- setChunkSize() (eskapade.analysis.links.ReadToDf method), 54
- simulate() (eskapade.analysis.histogram.Histogram method), 66
- skew (eskapade.analysis.histogram.ValueCounts attribute), 68
- SkipChain (eskapade.core.definitions.StatusCode attribute), 74
- SkipChainIfEmpty (class in eskapade.core_ops.links), 95
- SkipChainIfEmpty (class in eskapade.core_ops.links.skip_chain_if_empty), 91
- start_timer() (eskapade.mixins.TimerMixin method), 132
- StatusCode (class in eskapade.core.definitions), 74
- stop_timer() (eskapade.mixins.TimerMixin method), 132
- store() (eskapade.core.run_elements.Link method), 86
- Success (eskapade.core.definitions.StatusCode attribute), 74
- sum_counts (eskapade.analysis.histogram.ValueCounts attribute), 68
- sum_data_length() (eskapade.analysis.links.read_to_df.ReadToDf method), 47
- sum_data_length() (eskapade.analysis.links.ReadToDf method), 54
- sum_nonone_counts (eskapade.analysis.histogram.ValueCounts attribute), 68
- summary() (eskapade.core.run_elements.Link method), 86
- surface() (eskapade.analysis.histogram.Histogram method), 66
- ## T
- TimePeriod (class in eskapade.analysis.datetime), 60
- TimerMixin (class in eskapade.mixins), 132
- to_date_time() (in module eskapade.data_quality.dq_helper), 102
- to_float() (in module eskapade.data_quality.dq_helper), 102
- to_int() (in module eskapade.data_quality.dq_helper), 102
- to_normalized() (eskapade.analysis.histogram.Histogram method), 66
- to_ns() (in module eskapade.analysis.histogram_filling), 71
- to_str() (in module eskapade.analysis.histogram_filling), 71
- to_str() (in module eskapade.data_quality.dq_helper), 102
- ToDsDict (class in eskapade.core_ops.links), 96

ToDsDict (class in `eskapade.core_ops.links.to_ds_dict`), 91

`total_time()` (`eskapade.mixins.TimerMixin` method), 132

`truncated_bin_edges()` (`eskapade.analysis.histogram.BinningUtil` method), 62

`TruncExpFit` (class in `eskapade.root_analysis.links.trunc_exp_fit`), 111

`TruncExpGen` (class in `eskapade.root_analysis.links.trunc_exp_gen`), 112

`TruncExponential` (class in `eskapade.root_analysis.roofit_models`), 117

U

`Undefined` (`eskapade.core.definitions.StatusCode` attribute), 74

`UniformTsTimePeriod` (class in `eskapade.analysis.datetime`), 61

V

`value_to_bin_center()` (in module `eskapade.analysis.histogram_filling`), 71

`value_to_bin_index()` (in module `eskapade.analysis.histogram_filling`), 71

`value_to_bin_label()` (`eskapade.analysis.histogram.BinningUtil` method), 62

`ValueCounter` (class in `eskapade.analysis.links`), 57

`ValueCounter` (class in `eskapade.analysis.links.value_counter`), 49

`ValueCounts` (class in `eskapade.analysis.histogram`), 66

`var` (`eskapade.root_analysis.roofit_models.TruncExponential` attribute), 118

`var_set` (`eskapade.root_analysis.roofit_models.TruncExponential` attribute), 118

`variable` (`eskapade.analysis.histogram.Histogram` attribute), 66

W

`weighted_quantile()` (in module `eskapade.analysis.statistics`), 73

`WriteFromDf` (class in `eskapade.analysis.links`), 55

`WriteFromDf` (class in `eskapade.analysis.links.write_from_df`), 50

`ws` (`eskapade.root_analysis.roofit_manager.RooFitManager` attribute), 116

`ws` (`eskapade.root_analysis.roofit_models.RooFitModel` attribute), 117

`ws_contains()` (in module `eskapade.root_analysis.decorators.roofit`), 104

`ws_put()` (in module `eskapade.root_analysis.decorators.roofit`), 104

`ws_setitem()` (in module `eskapade.root_analysis.decorators.roofit`), 104

`WsUtils` (class in `eskapade.root_analysis.links.ws_utils`), 113

X

`xy_ranges_grid()` (in module `eskapade.root_analysis.decorators.histograms`), 104