

---

# **kedro-mflow**

***Release 0.2.0***

**Yolan Honoré-Rougé**

**Jul 18, 2020**



# CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Introduction	1
1.1.1	What is Kedro?	1
1.1.2	What is Mlflow?	1
1.1.3	A brief comparison between Kedro and Mlflow	2
1.1.3.1	Configuration and prototyping: Kedro 1 - 0 Mlflow	2
1.1.3.2	Versioning: Kedro 1 - 1 Mlflow	2
1.1.3.3	Model packaging and service: Kedro 1 - 2 Mlflow	3
1.1.3.4	Conclusion: Use Kedro and add Mlflow for machine learning projects	3
1.2	Motivation	3
1.2.1	When should I use kedro-mlflow?	3
1.2.2	Why should I use kedro-mlflow ?	3
1.2.2.1	Benchmark of existing solutions	3
1.2.2.2	Enforcing Kedro principles	4
1.3	Installation	4
1.3.1	Pre-requisites	4
1.3.2	Installation guide	5
1.3.3	Check the installation	5
1.3.4	Available commands	5
<b>2</b>	<b>Introduction</b>	<b>7</b>
2.1	Example project	7
2.1.1	Check your installation	7
2.1.2	Install the toy project	7
2.1.2.1	Installation with <code>kedro&gt;=0.16.3</code>	7
2.1.2.2	Installation with <code>kedro&gt;=0.16.0, &lt;=0.16.2</code>	8
2.2	Install dependencies	9
2.3	First step	9
2.3.1	Initialize kedro-mlflow	9
2.3.2	Run the pipeline	11
2.3.3	Open the UI	14
2.3.3.1	Parameters versioning	16
2.3.3.2	Journal information	18
2.3.3.3	Artifacts	18
<b>3</b>	<b>Introduction</b>	<b>19</b>
3.1	Scope	19
3.1.1	In the scope of the tutorial	19
3.1.2	Out of scope of the tutorial	19
3.2	Setup your Kedro project	20

3.2.1	Create a kedro project . . . . .	20
3.2.2	Update the template of your kedro project . . . . .	20
3.2.3	Automatic template update (recommended) . . . . .	20
3.2.3.1	Default situation . . . . .	20
3.2.3.2	Special case: what happens if you have a custom <code>run.py</code> ? . . . . .	20
3.2.4	Manual update . . . . .	21
3.3	Configure mlflow inside your project . . . . .	22
3.3.1	Context: mlflow tracking under the hood . . . . .	22
3.3.2	The <code>mlflow.yml</code> file . . . . .	22
3.4	Parameters versioning . . . . .	23
3.4.1	Automatic parameters versioning . . . . .	23
3.4.2	How does <code>MlflowNodeHook</code> operates under the hood? . . . . .	23
3.4.3	Frequently Asked Questions . . . . .	23
3.4.3.1	Will parameters be recorded if the pipeline fails during execution? . . . . .	23
3.4.3.2	How are parameters detected by the plugin? . . . . .	24
3.4.3.3	How can I register a parameter if I use a <code>TemplatedConfigLoader</code> ? . . . . .	24
3.5	Opening the UI . . . . .	24
3.5.1	The mlflow user interface . . . . .	24
3.5.2	The kedro-mlflow helper . . . . .	24
<b>4</b>	<b>Indices and tables</b>	<b>25</b>

## INTRODUCTION

### 1.1 Introduction

#### 1.1.1 What is Kedro?

Kedro is a python package which facilitates the prototyping of data pipelines. It aims at implementing software engineering best practices (separation between I/O and compute, abstraction, templating...). It is specifically useful for machine learning projects since it provides within the same interface both interactive objects for the exploration phase and *Command Line Interface* (CLI) and configuration files for the production phase. This makes the transition from exploration to production as smooth as possible.

For more details, see [Kedro's official documentation](#).

#### 1.1.2 What is Mlflow?

Mlflow is a library which helps managing the lifecycle of machine learning models. Mlflow provides 4 modules:

- **Mlflow Tracking:** This module focuses on experiment versioning. The goal is to store all the objects needed to reproduce any code execution. This includes code through version control, but also parameters and artifacts (i.e objects fitted on data like encoders, binarizers...). These elements vary wildly during machine learning experimentation phase. Mlflow also enable to track metrics to evaluate runs, and provides a *User Interface* (UI) to browse the different runs and compare them.
- **Mlflow Projects:** This module provides a configuration files and CLI to enable reproducible execution of pipelines in production phase.
- **Mlflow Models:** This module defines a standard way for packaging machine learning models, and provides built-in ways to serve registered models. Such standardization enable to serve these models across a wide range of tools.
- **Mlflow Model Registry:** This module aims at monitoring deployed models. The registry manages the transition between different versions of the same model (when the dataset is retrained on new data, or when parameters are updated) while it is in production.

For more details, see [Mlflow's official documentation](#).

### 1.1.3 A brief comparison between Kedro and Mlflow

While Kedro and Mlflow do not compete in the same field, they provide some overlapping functionalities. Mlflow is specifically dedicated to machine learning and its lifecycle management, while Kedro focusing on data pipeline development. Below chart compare the different functionalities:

We can draw the following conclusions from the chart, discussed hereafter.

#### 1.1.3.1 Configuration and prototyping: Kedro 1 - 0 Mlflow

Mlflow and Kedro are essentially overlapping on the way they offer a dedicated configuration files for running the pipeline from CLI. However:

- Mlflow provides a single configuration file (the `MLProject`) where all elements are declared (data, parameters and pipelines). Its goal is mainly to enable CLI execution of the project, but it is not very flexible. In my opinion, this file is **production oriented** and is not really intended to use for exploration.
- Kedro offers a bunch of files (`catalog.yml`, `parameters.yml`, `pipeline.py`) and their associated abstraction (`AbstractDataSet`, `DataCatalog`, `Pipeline` and `node` objects). Kedro is much more opinionated: each object has a dedicated place (and only one!) in the template. This makes the framework both **exploration and production oriented**. The downside is that it could make the learning curve a bit sharper since a newcomer has to learn all Kedro specifications. It also provides a `kedro-viz` plugin to visualize the DAG interactively, which is particularly handy in medium-to-big projects.

#### 1.1.3.2 Versioning: Kedro 1 - 1 Mlflow

The Kedro [Journal](#) aims at [reproducibility](#), but is not focused on machine learning. The Journal keeps track of two elements:

- the CLI arguments , including *on the fly* parameters. This makes the command used to run the pipeline fully reproducible.
- the `AbstractVersionedDataSet` for which versioning is activated. It consists in copying the data whom `versioned` argument is `True` when the `save` method of the `AbstractVersionedDataSet` is called. This approach suffers from two main drawbacks:
  - the configuration is assumed immutable (including parameters), which is not realistic in machine learning projects where they are very volatile. To fix this, the `git sha` has been recently added to the `Journal`, but it has still some bugs in my experience (including the fact that the current `git sha` is logged even if the pipeline is ran with uncommitted change, which prevents reproducibility). This is still recent and will likely evolve in the future.
  - there is no support for browsing old runs, which prevents [cleaning the database with old and unused datasets](#), compare runs between each other...

On the other hand, Mlflow:

- distinguishes between artifacts (i.e. any data file), metrics (integers that may evolve over time) and parameters. The logging is very straightforward since there is a one-liner function for logging the desired type. This separation makes further manipulation easier.
- offers a way to configure the logging in a database through the `mlflow_tracking_uri` parameter. This database-like logging comes with easy [querying of different runs through a client](#) (for instance “find the most recent run with a metric at least above a given threshold” is immediate with Mlflow but hacky in Kedro).
- [comes with a User Interface \(UI\)](#) which enable to browse / filter / sort the runs, display graphs of the metrics, render plots... This make the run management much easier than in Kedro.

- has a command to reproduce exactly the run from a given `git sha`, which is not possible in `Kedro`.

### 1.1.3.3 Model packaging and service: Kedro 1 - 2 Mlflow

`Kedro` offers a way to package the code to make the pipelines callable, but does not manage specifically machine learning models.

`Mlflow` offers a way to store machine learning models with a given “flavor”, which is the minimal amount of information necessary to use the model for prediction:

- a configuration file
- all the artifacts, i.e. the necessary data for the model to run (including encoder, binarizer...)
- a loader
- a conda configuration through an `environment.yml` file

When a stored model meets these requirements, `Mlflow` provides built-in tools to serve the model (as an API or for batch prediction) on many machine learning tools (Microsoft Azure ML, Amazon Sagemaker, Apache SparkUDF) and locally.

### 1.1.3.4 Conclusion: Use Kedro and add Mlflow for machine learning projects

In my opinion, `Kedro`’s will to enforce software engineering best practice makes it really useful for machine learning teams. It is extremely well documented and the support is excellent, which makes it very user friendly even for people with no CS background. However, it lacks some machine learning-specific functionalities (better versioning, model service), and it is where `Mlflow` fills the gap.

## 1.2 Motivation

### 1.2.1 When should I use kedro-mlflow?

Basically, you should use `kedro-mlflow` in **any `Kedro` project which involves machine learning** / deep learning. As stated in the [introduction](#), `Kedro`’s current versioning (as of version 0.16.1) is not sufficient for machine learning projects: it lacks a UI and a run management system. Besides, the `KedroPipelineModel` ability to serve a `kedro` pipeline as an API or a batch in one line of code is a great addition for collaboration and transition to production.

If you do not use `Kedro` or if you do pure data manipulation which do not involve machine learning, this plugin is not what you are seeking for ;)

### 1.2.2 Why should I use kedro-mlflow ?

#### 1.2.2.1 Benchmark of existing solutions

This paragraph gives a (quick) overview of existing solutions for `mlflow` integration inside `Kedro` projects.

`Mlflow` is very simple to add to any existing code. It is a 2-step process:

- add `log_{XXX}` (either param, artifact, metric or model) functions where they are needed inside the code
- add a `MLProject` at the root of the project to enable CLI execution. This file must contain all the possible execution steps (like the `pipeline.py` in a `kedro` project).

Including mlflow inside a `kedro` project is consequently very easy: the logging functions can be added in the code, and the `MLProject` is very simple and is composed almost only of the `kedro run` command. You can find examples of such implementation:

- the [medium paper](#) by QuantumBlack employees.
- the associated [github repo](#)
- other examples can be found on Github, but AFAIK all of them follow the very same principles.

### 1.2.2.2 Enforcing Kedro principles

Above implementations have the advantage of being very straightforward and *mlflow compliant*, but it breaks several Kedro principles:

- the `MLFLOW_TRACKING_URI` which registers the database where runs are logged is declared inside the code instead of a configuration file, which **hinders portability across environments** and makes transition to production more difficult
- the logging of different elements can be put in many places in the Kedro template (in the code of any function involved in a node, in a Hook, in the `ProjectContext`, in a `transformer...`). This is not compliant with the Kedro template where any object has a dedicated location. We want to avoid the logging to occur anywhere because:
  - it is **very error-prone** (one can forget to log one parameter)
  - it is **hard to modify** (if you want to remove / add / modify an mlflow action you have to find it in the code)
  - it **prevents reuse** (re-usable function must not contain mlflow specific code unrelated to their functional specificities, only their execution must be tracked).

`kedro-mlflow` enforces these best practices while implementing a clear interface for each mlflow action in Kedro template. Below chart maps the mlflow action to perform with the Python API provided by `kedro-mlflow` and the location in Kedro template where the action should be performed.

In the current version (`kedro_mlflow=0.2.0`), `kedro-mlflow` do not provide interface to log metrics, set tags or log models outside a Kedro Pipeline. These decisions are subject to debate and design decisions (for instance, metrics are often updated in a loop during each epoch / training iteration and it does not always make sense to register the metric between computation steps, e.g. as a an I/O operation after a node run).

***Note:** the version 0.2.0 does not need any `MLProject` file to use mlflow inside your Kedro project. As seen in the introduction, this file overlaps with Kedro configuration files.*

## 1.3 Installation

### 1.3.1 Pre-requisites

I strongly recommend to use `conda` (a package manager) to create an environment in order to avoid version conflicts between packages.

I also recommend to read [Kedro installation guide](#) to set up your Kedro project.



### 1.3.2 Installation guide

The plugin is compatible with `kedro>=0.16.0`. Since Kedro tries to enforce backward compatibility, it will very likely remain compatible with further versions.

First, install Kedro from PyPI and ensure you have a `0.16.0` version:

```
pip install --upgrade "kedro>=0.16.0,<0.17.0"
```

Second, and since the `kedro-mlflow` plugin is not on PyPi yet, you must install it from sources:

```
pip install git+https://github.com/Galileo-Galilei/kedro-mlflow.git
```

You may want to install the develop branch which has unreleased features:

```
pip install git+https://github.com/Galileo-Galilei/kedro-mlflow.git@develop
```

### 1.3.3 Check the installation

Type `kedro info` in a terminal to check the installation. If it has succeeded, you should see the following ascii art:

```

_ | _ _ _ _ _ | _ _ _ _
| | / / _ \ / _ | ' _ / _ \
| | < _ / ( | | | | ( ) |
|_| \ \ _ _ | \ _ _ | \ _ _ /
v0.16.2

kedro allows teams to create analytics
projects. It is developed as part of
the Kedro initiative at QuantumBlack.

Installed plugins:
kedro_mlflow: 0.2.0 (hooks:global,project)
```

The version `0.2.0` of the plugin is installed and has both global and project commands.

That's it! You are now ready to go!

### 1.3.4 Available commands

With the `kedro mlflow -h` command outside of a kedro project, you now see the following output:

```
Usage: kedro mlflow [OPTIONS] COMMAND [ARGS]...

    Use mlflow-specific commands inside kedro project.

Options:
  -h, --help  Show this message and exit.

Commands:
  new  Create a new kedro project with updated template.
```



## INTRODUCTION

### 2.1 Example project

#### 2.1.1 Check your installation

Create a conda environment and `kedro-mlflow` (this will automatically install `kedro>=0.16.0`).

```
conda create -n km_example python=3.6.8 --yes
conda activate km_example
pip install kedro-mlflow
```

#### 2.1.2 Install the toy project

For this end to end example, we will use the `kedro starter` of with the `iris dataset`.

We use this project because:

- it covers most of the common use cases
- it is compatible with older version of `Kedro` so newcomers are used to it
- it is maintained by `Kedro` maintainers and therefore enforces some best practices.

##### 2.1.2.1 Installation with `kedro>=0.16.3`

The default starter is now called “`pandas-iris`”. In a new console, enter:

```
kedro new --starter=pandas-iris
```

Answer `Kedro Mlflow Example`, `km-example` and `km_example` to the three setup questions of a new `kedro` project:

```
Project Name:
=====
Please enter a human readable name for your new project.
Spaces and punctuation are allowed.
[New Kedro Project]: Kedro Mlflow Example

Repository Name:
=====
Please enter a directory name for your new project repository.
Alphanumeric characters, hyphens and underscores are allowed.
```

(continues on next page)

(continued from previous page)

```

Lowercase is recommended.
[kedro-mlflow-example]: km-example

Python Package Name:
=====
Please enter a valid Python package name for your project package.
Alphanumeric characters and underscores are allowed.
Lowercase is recommended. Package name must start with a letter or underscore.
[kedro_mlflow_example]: km_example

```

### 2.1.2.2 Installation with `kedro>=0.16.0, <=0.16.2`

With older versions of Kedro, the starter option is not available, but this `kedro new` provides an “Include example” question. Answer `y` to this question to get the same starter as above. In a new console, enter:

```
kedro new
```

Answer Kedro Mlflow Example, `km-example`, `km_example` and `y` to the four setup questions of a new kedro project:

```

Project Name:
=====
Please enter a human readable name for your new project.
Spaces and punctuation are allowed.
[New Kedro Project]: Kedro Mlflow Example

Repository Name:
=====
Please enter a directory name for your new project repository.
Alphanumeric characters, hyphens and underscores are allowed.
Lowercase is recommended.
[kedro-mlflow-example]: km-example

Python Package Name:
=====
Please enter a valid Python package name for your project package.
Alphanumeric characters and underscores are allowed.
Lowercase is recommended. Package name must start with a letter or underscore.
[kedro_mlflow_example]: km_example

Generate Example Pipeline:
=====
Do you want to generate an example pipeline in your project?
Good for first-time users. (default=N)
[y/N]: y

```

## 2.2 Install dependencies

Move to the project directory:

```
cd km-example
```

Install the project dependencies:

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

**Warning:** Do not use `kedro install` commands does not seem to install the packages in your activated environment.

## 2.3 First step

### 2.3.1 Initialize kedro-mlflow

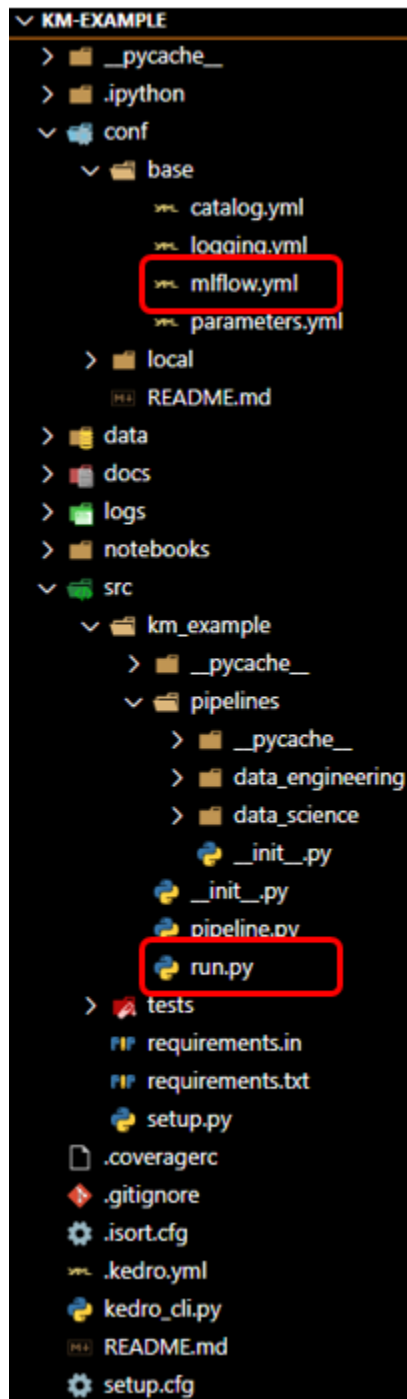
Run

```
kedro mlflow init
```

You have the following message:

```
'conf/base/mlflow.yml' successfully updated.  
'run.py' successfully updated
```

The `conf/base` folder is updated:



If you have configured your own mlflow server, you can specify the tracking uri in the `mlflow.yml` (replace the highlighted line below:):

```

mlflow.yml X
conf > base > mlflow.yml
1  # GLOBAL CONFIGURATION -----
2
3  # `mlflow_tracking_uri` is the path where the runs will be recorded.
4  # For more informations, see https://www.mlflow.org/docs/latest/tracking.html#where-runs-are-recorded
5  # kedro-mlflow accepts relative path from the project root.
6  # For instance, default `mlruns` will create a `mlruns` folder
7  # at the root of the project
8  mlflow_tracking_uri: mlruns
9
10
11 # EXPERIMENT-RELATED PARAMETERS -----
12
13 # `name` is the name of the experiment (~subfolder
14 # where the runs are recorded). Change the name to
15 # switch between different experiments
16 experiment:
17   name: km_example
18   create: True # if the specified `name` does not exists, should it be created?
19
20
21 # RUN-RELATED PARAMETERS -----
22
23 run:
24   id: null # if `id` is None, a new run will be created
25   name: null # if `name` is None, pipeline name will be used for the run name
26   nested: True # if `nested` is False, you won't be able to launch sub-runs inside your nodes
27
28 # UI-RELATED PARAMETERS -----
29
30 ui:
31   port: null # the port to use for the ui. Find a free port if null.
32   host: null # the host to use for the ui. Default to "localhost" if null.
33

```

## 2.3.2 Run the pipeline

Open a new command and launch

```
kedro run
```

If the pipeline executes properly, you should see the following log:

```

2020-07-13 21:29:24,939 - kedro.versioning.journal - WARNING - Unable to git describe_
↳ path/to/km-example
2020-07-13 21:29:25,401 - kedro.io.data_catalog - INFO - Loading data from `example_
↳ iris_data` (CSVDataSet)...
2020-07-13 21:29:25,562 - kedro.io.data_catalog - INFO - Loading data from_
↳ `params:example_test_data_ratio` (MemoryDataSet)...
2020-07-13 21:29:25,969 - kedro.pipeline.node - INFO - Running node: split_
↳ data([example_iris_data,params:example_test_data_ratio]) -> [example_test_x,example_
↳ test_y,example_train_x,example_train_y]
2020-07-13 21:29:26,053 - kedro.io.data_catalog - INFO - Saving data to `example_
↳ train_x` (MemoryDataSet)...
2020-07-13 21:29:26,368 - kedro.io.data_catalog - INFO - Saving data to `example_
↳ train_y` (MemoryDataSet)...
2020-07-13 21:29:26,484 - kedro.io.data_catalog - INFO - Saving data to `example_test_
↳ x` (MemoryDataSet)...

```

(continues on next page)

(continued from previous page)

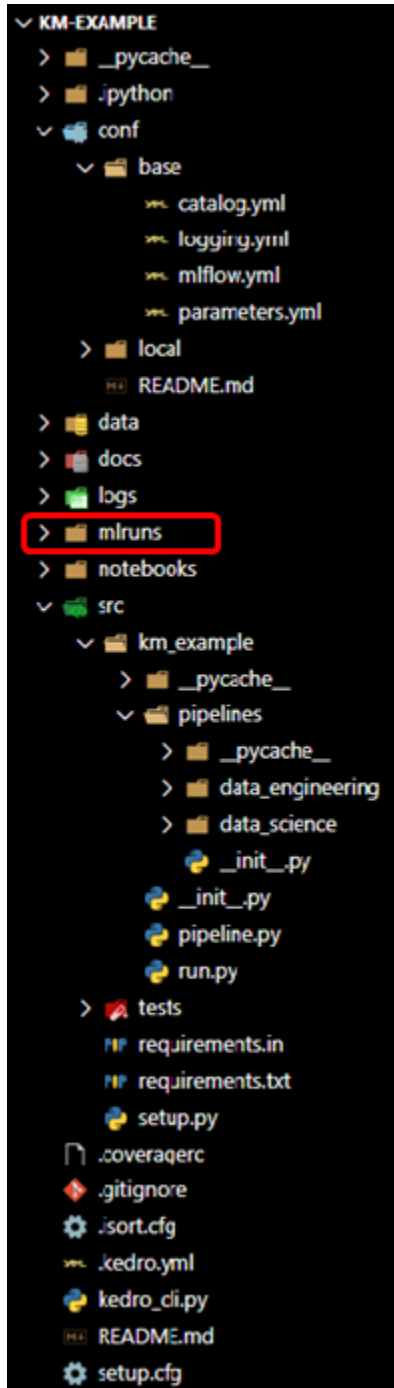
```

2020-07-13 21:29:26,486 - kedro.io.data_catalog - INFO - Saving data to `example_test_
↳y` (MemoryDataSet)...
2020-07-13 21:29:26,610 - kedro.runner.sequential_runner - INFO - Completed 1 out of 4
↳tasks
2020-07-13 21:29:26,850 - kedro.io.data_catalog - INFO - Loading data from `example_
↳train_x` (MemoryDataSet)...
2020-07-13 21:29:26,851 - kedro.io.data_catalog - INFO - Loading data from `example_
↳train_y` (MemoryDataSet)...
2020-07-13 21:29:26,965 - kedro.io.data_catalog - INFO - Loading data from
↳`parameters` (MemoryDataSet)...
2020-07-13 21:29:26,972 - kedro.pipeline.node - INFO - Running node: train_
↳model([example_train_x,example_train_y,parameters]) -> [example_model]
2020-07-13 21:29:27,756 - kedro.io.data_catalog - INFO - Saving data to `example_
↳model` (MemoryDataSet)...
2020-07-13 21:29:27,763 - kedro.runner.sequential_runner - INFO - Completed 2 out of 4
↳tasks
2020-07-13 21:29:28,141 - kedro.io.data_catalog - INFO - Loading data from `example_
↳model` (MemoryDataSet)...
2020-07-13 21:29:28,161 - kedro.io.data_catalog - INFO - Loading data from `example_
↳test_x` (MemoryDataSet)...
2020-07-13 21:29:28,670 - kedro.pipeline.node - INFO - Running node: predict([example_
↳model,example_test_x]) -> [example_predictions]
2020-07-13 21:29:29,002 - kedro.io.data_catalog - INFO - Saving data to `example_
↳predictions` (MemoryDataSet)...
2020-07-13 21:29:29,248 - kedro.runner.sequential_runner - INFO - Completed 3 out of 4
↳tasks
2020-07-13 21:29:29,433 - kedro.io.data_catalog - INFO - Loading data from `example_
↳predictions` (MemoryDataSet)...
2020-07-13 21:29:29,730 - kedro.io.data_catalog - INFO - Loading data from `example_
↳test_y` (MemoryDataSet)...
2020-07-13 21:29:29,911 - kedro.pipeline.node - INFO - Running node: report_
↳accuracy([example_predictions,example_test_y]) -> None
2020-07-13 21:29:30,056 - km_example.pipelines.data_science.nodes - INFO - Model
↳accuracy on test set: 100.00%
2020-07-13 21:29:30,214 - kedro.runner.sequential_runner - INFO - Completed 4 out of 4
↳tasks
2020-07-13 21:29:30,372 - kedro.runner.sequential_runner - INFO - Pipeline execution
↳completed successfully.

```

Since we have kept the default value of the `mlflow.yml`, the tracking uri (the place where runs are recorded) is a local `mlruns` folder which has just been created with the execution:





### 2.3.3 Open the UI

Launch the ui:

```
kedro mlflow ui
```

And open the following address in your favorite browser

<http://localhost:5000/>

The screenshot shows the mlflow web interface. On the left sidebar, under 'Experiments', 'km\_example' is selected. A red box highlights this selection, with a red text annotation: 'The name of the experiment in "mlflow.yml"'. The main panel displays details for the 'km\_example' experiment (ID: 1). It shows the artifact location as 'file:///C:/Local/path/to/m-example/mlruns/1'. Below this, there's a search bar with filters: 'metrics.rmse < 1 and params.model = "tree" and tags.mlflow.source.type = "LOCAL"'. A table of runs is shown below, with one run highlighted by a red box. This run has a green status icon and a timestamp of '2020-07-13 21:29:24'. A red text annotation 'Last run executed' points to this row. The table columns include Start Time, Run Name, User (You), Source (Python path), Version (0.2), Parameters (example\_test\_data\_ref, parameters), Env (env), Extra Params ({}), and From Inputs ({}).

Click now on the last run executed, you will land on this page:

**km\_example > Run 9128c4c15e2c438db27749561f543c97** ▾

Date: 2020-07-13 21:29:24

Source: 

\\km\_example\Scripts\kedro

Duration: 5.6s

Status: FINISHED

▼ Notes [🔗](#)

None

## ▼ Parameters

Name	Value
example_test_data_ratio	0.2
parameters	{'example_test_data_ratio': 0.2, 'example_num_train_iter': 10000, 'example_learning_rate': 0.01}

## ▼ Metrics

Name	Value
------	-------

## ▼ Tags

Name	Value	Actions
env	local	<a href="#">🔗</a> <a href="#">🗑️</a>
extra_params	{}	<a href="#">🔗</a> <a href="#">🗑️</a>
from_inputs	[]	<a href="#">🔗</a> <a href="#">🗑️</a>
from_nodes	[]	<a href="#">🔗</a> <a href="#">🗑️</a>
git_sha	None	<a href="#">🔗</a> <a href="#">🗑️</a>
kedro_command	kedro run	<a href="#">🔗</a> <a href="#">🗑️</a>
kedro_version	0.16.3	<a href="#">🔗</a> <a href="#">🗑️</a>
load_versions	{}	<a href="#">🔗</a> <a href="#">🗑️</a>
node_names	()	<a href="#">🔗</a> <a href="#">🗑️</a>
pipeline_name	None	<a href="#">🔗</a> <a href="#">🗑️</a>
project_path	\\km-example	<a href="#">🔗</a> <a href="#">🗑️</a>
run_id	2020-07-13T19:29:20.514Z	<a href="#">🔗</a> <a href="#">🗑️</a>
tags	()	<a href="#">🔗</a> <a href="#">🗑️</a>
to_nodes	[]	<a href="#">🔗</a> <a href="#">🗑️</a>

## Add Tag

<input type="text" value="Name"/>	<input type="text" value="Value"/>	<input type="button" value="Add"/>
-----------------------------------	------------------------------------	------------------------------------

## ▼ Artifacts

No Artifacts Recorded

Use the log artifact APIs to store file outputs from MLflow runs.

### 2.3.3.1 Parameters versioning

Note that the parameters have been recorded *automagically*. Here, two parameters format are used:

1. The parameter `example_test_data_ratio`, which is called in the `pipeline.py` file with the `params :` prefix
2. the dictionary of all parameters in `parameters.yml` which is a reserved key word in Kedro. Note that **this is bad practice** because you cannot know which parameters are really used inside the function called. Another problem is that it can generate too long parameters names and lead to mlflow errors.

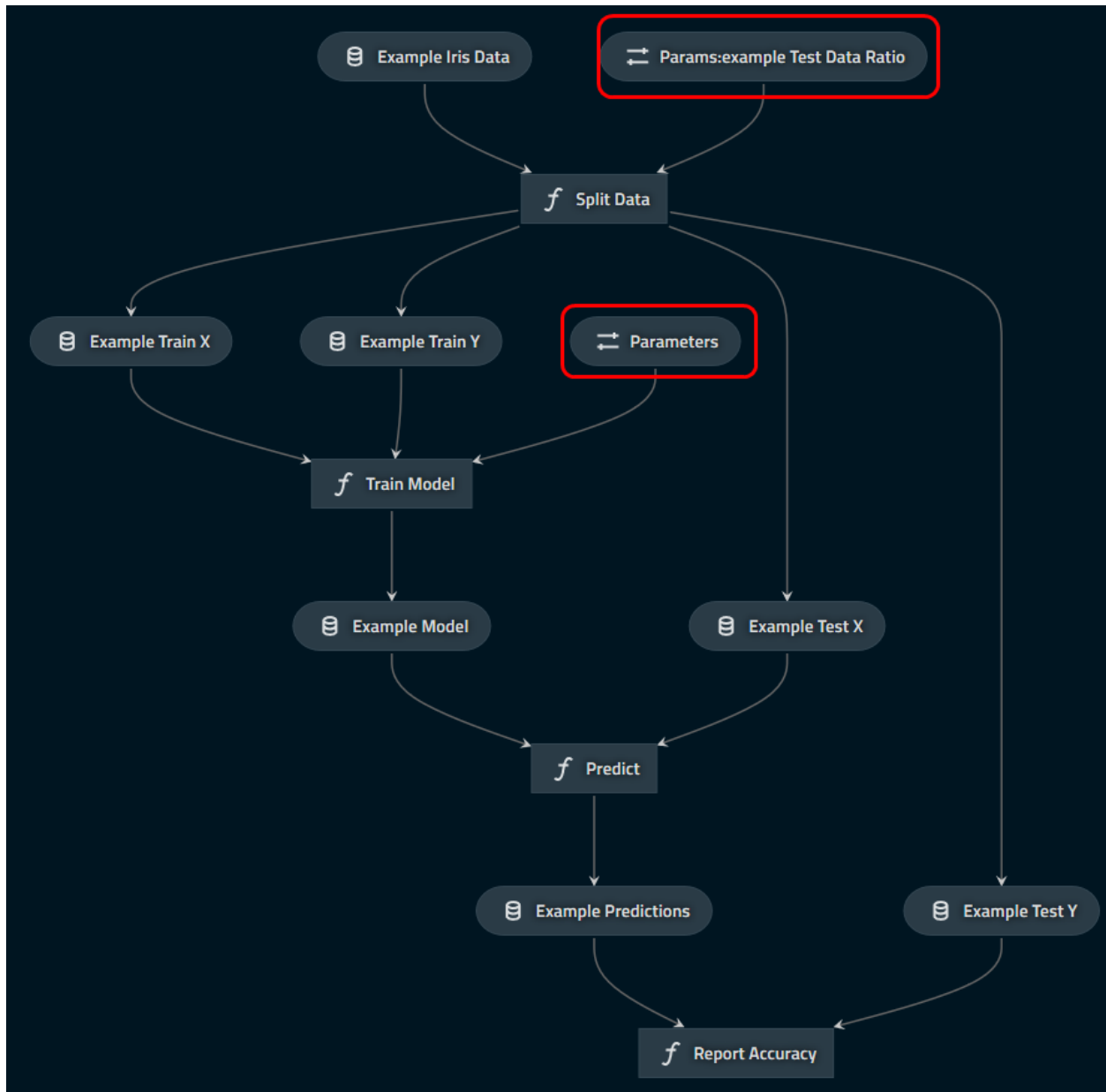
You can see that these are effectively the registered parameters in the pipeline with the `kedro-viz` plugin:

```
pip install kedro-viz
kedro viz
```

Open your browser at the following adress:

```
http://localhost:4141/
```

You should see the following graph:



which indicates clearly which parameters are logged (in the red boxes with the “parameter” icon).

### 2.3.3.2 Journal information

The informations provided by the `Kedro's Journal` are also recorded as `tags` in the `mlflow ui` in order to make reproducible. In particular, the exact command used for running the pipeline and the `kedro` version used are stored.

### 2.3.3.3 Artifacts

With this run, artifacts are empty. This is expected: `mlflow` does not know what it should log and it will not log all your data by default. However, you want to save your model (at least) or your run is likely useless!

First, open the `catalog.yml` file which should like this:

```
# This is a data set used by the "Hello World" example pipeline provided with the project
# template. Please feel free to remove it once you remove the example pipeline.

example_iris_data:
  type: pandas.CSVDataSet
  filepath: data/01_raw/iris.csv
```

And persist the model as a pickle with the `MlflowDataSet` class:

```
# This is a data set used by the "Hello World" example pipeline provided with the project
# template. Please feel free to remove it once you remove the example pipeline.

example_iris_data:
  type: pandas.CSVDataSet
  filepath: data/01_raw/iris.csv

example_model:
  type: kedro_mlflow.io.MlflowDataSet
  data_set:
    type: pickle.PickleDataSet
    filepath: data/06_models/trained_model.pkl
```

Reopen the ui, select the last run and see that the file was uploaded:

#### ▼ Artifacts



This works for any type of file (including images with `MatplotlibWriter`) and the UI even offers a preview for `png` and `csv`, which is really convenient to compare runs.

*Note: `Mlflow` offers specific logging for machine learning models that should be better suited for your use case, but is not supported yet in `kedro-mlflow==0.2.0`*

## INTRODUCTION

### 3.1 Scope

#### 3.1.1 In the scope of the tutorial

This tutorial addresses the following items:

1. How to include `kedro-mlflow` capabilities in a Kedro project:
  1. create a new kedro project with updated template
  2. update an existing kedro project
2. Configure mlflow inside a Kedro project
3. Version and track objects during execution with mlflow:
  1. Version parameters inside a Kedro project
  2. Version data inside a Kedro project
  3. **(COMING in 0.3.0)** Version machine learning models inside a Kedro project
  4. **(COMING in 0.3.0)** Version metrics inside a Kedro project
  5. Open mlflow ui with project configuratio

This is a step by step tutorial and it is recommended to read the different chapters above order, but not mandatory.

#### 3.1.2 Out of scope of the tutorial

Some advanced capabilities are adressed in the [advanced use section](#):

- saving a kedro pipeline as a mlflow model and serve it
- **(COMING in 0.3.0)** launching a Kedro project directly with mlflow trthough the `MLProject` file.

## 3.2 Setup your Kedro project

### 3.2.1 Create a kedro project

This plugins must be used in an existing kedro project. If you do not have a kedro project yet, you can create it with `kedro new` command. [See the kedro docs for a tutorial](#).

For this tutorial and if you do not have a real-world project, I strongly suggest that you accept to include the proposed example to make a demo of this plugin out of the box.

### 3.2.2 Update the template of your kedro project

In order to use the `kedro-mlflow` plugin, you need to perform 2 actions:

1. Create an `mlflow.yml` file for [configuring mlflow in a dedicated file](#).
2. Update the `src/PYTHON_PACKAGE/run.py` to add the [necessary hooks](#) to the project context. The `MlflowPipelineHook` manages the configuration and registers the `PipelineML`, while the `MlflowNodeHook` autolog the parameters.

### 3.2.3 Automatic template update (recommended)

#### 3.2.3.1 Default situation

The first and recommended possibility to setup this context is to use a [dedicated command line](#) offered by the plugin. Position yourself with at the root (i.e. the folder with the `.kedro.yml` file)

```
$ cd path/to/your/project
```

Run the init command :

```
$ kedro mlflow init
```

*Note : If the warning "You have not updated your template yet. This is mandatory to use 'kedro-mlflow' plugin. Please run the following command before you can access to other commands : '\$ kedro mlflow init' is raised, this is a bug to be corrected and you can safely ignore it. If you have never modified your `run.py` manually, it should run smoothly and you should get the following message:*

```
'conf/base/mlflow.yml' successfully updated.  
'run.py' successfully updated
```

#### 3.2.3.2 Special case: what happens if you have a custom `run.py` ?

You may have modified the `run.py` manually since the creation of the project. This may happen in the following situations:

- you have added hooks (of another plugin for instance)
- you have modified the `ConfigLoader`, for instance to us a `TemplatedConfigLoader` to make your configuration dynamic and link the files with one another



- you have modified the `get_pipelines` functions to implement specific logic ... These are advanced features of Kedro and if you have made such modifications they are very likely conscious; however some other plugins may have modified this file without any warning.

Whatever the reason is, if you `run.py` was modified since the project creation, the *previous process* will return the following warning message:

```
You have modified your 'run.py' since project creation.
In order to use kedro-mlflow, you must either:
- set up your run.py with the following instructions :
INSERT_DOC_URL
- call the following command:
$ kedro mlflow init --force
```

In this situation, the `mlflow.yml` is still created, but the `run.py` is left unchanged to avoid messing up with your own changes. You can still erase your `run.py` and replace it with the one of the plugin with below command.

```
kedro mlflow init --force
```

**USE AT YOUR OWN RISK: This will erase definitely all the modifications you made to your own `run.py` with no possible recovery.** In consequence, this is not the recommended way to setup the project if you have a custom `run.py`. The best way to continue the setup is to *set up the hooks manually*.

### 3.2.4 Manual update

The `MlflowPipelineHook` and `MlflowNodeHook` hooks need to be registered in the `run.py` file. The kedro documentation explain since tail [how to register a hook](#).

Your `run.py` should look like the following code snippet :

```
from kedro_mlflow.framework.hooks import MlflowNodeHook, MlflowPipelineHook
from YOUR_PYTHON_PACKAGE.pipeline import create_pipelines

class ProjectContext(KedroContext):
    """Users can override the remaining methods from the parent class here,
    or create new ones (e.g. as required by plugins)
    """

    project_name = "YOUR PROJECT NAME"
    project_version = "0.16.2"
    hooks = (
        MlflowNodeHook(flatten_dict_params=False),
        MlflowPipelineHook(model_name="YOUR_PYTHON_PACKAGE",
                           conda_env="src/requirements.txt")
    ) # <-- the new lines to add
```

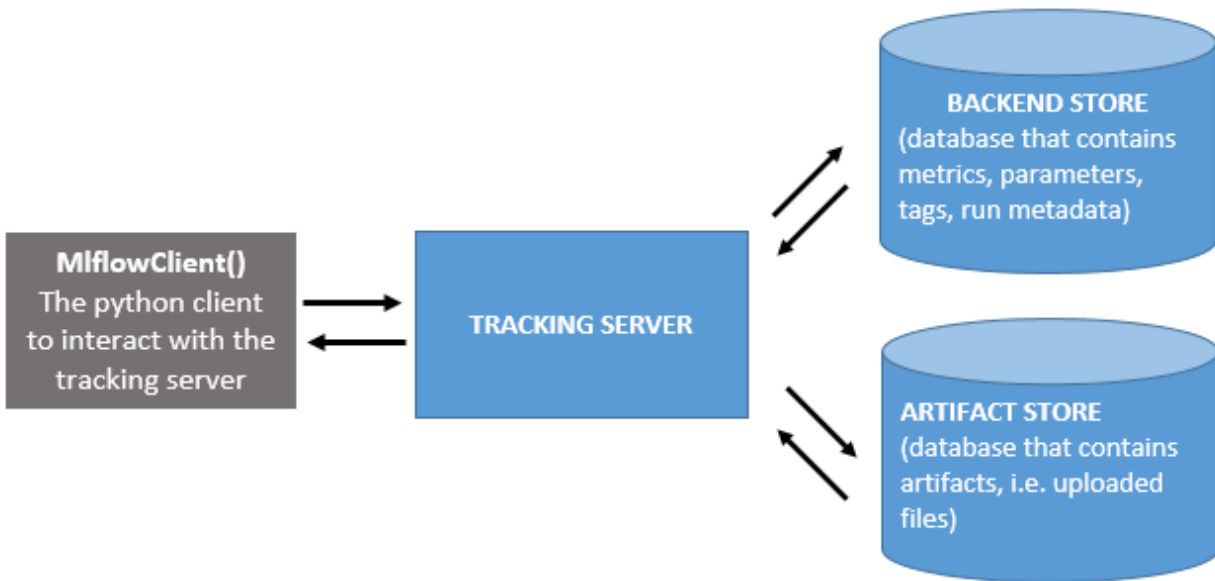
Pay attention to the following elements:

- if you have other hooks (custom, from other plugins...), you can just add them to the hooks tuple
- you **must register both hooks** for the plugin to work
- the hooks are highly parametrizable, you can find a [detailed description of their parameters here](#).

## 3.3 Configure mlflow inside your project

### 3.3.1 Context: mlflow tracking under the hood

Mlflow is composed of four modules which are described in the [introduction section](#). The main module is “tracking”. The goal of this module is to keep track of every varying parameters across different code execution (parameters, metrics and artifacts). The following schema describes how this module operates under the hood:



Basically, this schema shows that mlflow separates WHERE the artifacts are logged from HOW they are logged inside your code. You need to setup your mlflow tracking server separately from your code, and then each logging will send a request to the tracking server to store the elements you want to track in the appropriate location. The advantage of such a setup are numerous:

- once the mlflow tracking server is setup, there is single parameter to set before logging which is the tracking server uri. This makes configuration very easy in your project.
- since the different storage locations are well identified, it is easy to define custom solutions for each of them. They can be [database](#) or [even local folders](#).

The rationale behind the separation of the backend store and the artifacts store is that artifacts can be very big and are duplicated across runs, so they need a special management with extensible storage. This is typically [cloud storage like AWS S3 or Azure Blob storage](#).

### 3.3.2 The `mlflow.yml` file

kedro-mlflow needs the tracking uri of your mlflow tracking server to operate properly. The `mlflow.yml` file must have the `mlflow_tracking_uri` key with a [valid mlflow\\_tracking\\_uri associated value](#). The `mlflow.yml` default have this keys set to `mlruns`. This will create a `mlruns` folder locally at the root of your kedro project and enable you to use the plugin without any setup of a mlflow tracking server.

```
mlflow_tracking_uri: mlruns
```

This is the only mandatory key in the `mlflow.yml` file, but there are many others that provides fine-grained control on your mlflow setup. Please see the [mlflow.yml description](#) for further details.

## 3.4 Parameters versioning

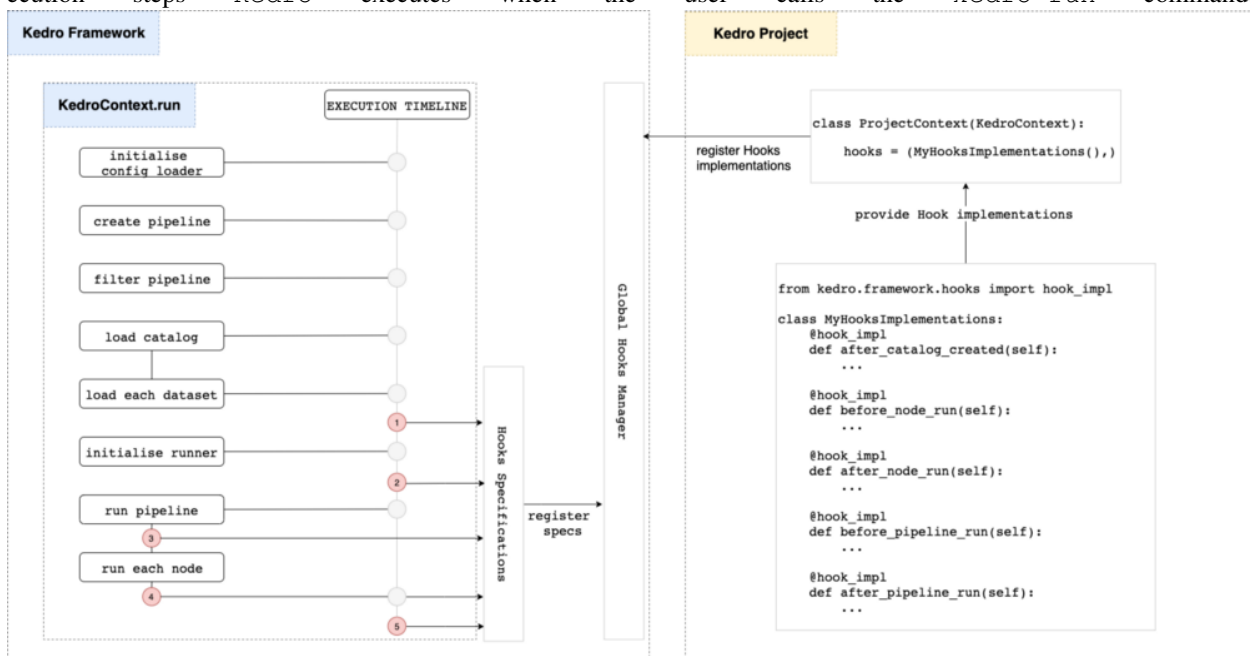
### 3.4.1 Automatic parameters versioning

Parameters versioning is automatic when the `MlflowNodeHook` is added to the hook list of the `ProjectContext`. In `kedro-mlflow==0.2.0`, this hook has a parameter called `flatten_dict_params` which enables to log as distinct parameters the (key, value) pairs of a ``Dict`` parameter.

You **do not need any additional configuration** to benefit from parameters versioning.

### 3.4.2 How does `MlflowNodeHook` operates under the hood?

The `medium post` which introduces hooks explains in detail the different execution steps Kedro executes when the user calls the `kedro run` command.



The `MlflowNodeHook` registers the parameters before each node (entry point number 3 on above picture) by calling ``mlflow.log\_parameter(param\_name, param\_value)`` on each parameters of the node.

### 3.4.3 Frequently Asked Questions

#### 3.4.3.1 Will parameters be recorded if the pipeline fails during execution?

The parameters are registered node by node (and not in a single batch at the beginning of the execution). If the pipeline fails in the middle of its execution, the **parameters of the nodes who have been run will be recorded, but not the parameters of non executed nodes.**

### 3.4.3.2 How are parameters detected by the plugin?

The hook **detects parameters through their prefix `params:` or the value `parameters`**. These are the **reserved keywords** used by Kedro to define parameters in the `pipeline.py` file(s).

### 3.4.3.3 How can I register a parameter if I use a `TemplatedConfigLoader`?

If you use a `TemplatedConfigLoader` to enable dynamic parameters construction at runtime or dependency between configuration files, and if we assume your `src/<project-name>/run.py` file looks like:

```
from kedro.config import TemplatedConfigLoader # new import
from datetime import date

class ProjectContext(KedroContext):
    def _create_config_loader(self, conf_paths: Iterable[str]) ->
    ↪TemplatedConfigLoader:
        return TemplatedConfigLoader(
            conf_paths,
            globals_pattern="*globals.yml", # read the globals dictionary from
    ↪project config
            globals_dict={ # extra keys to add to the globals dictionary, take
    ↪precedence over globals_pattern
                execution_date: date.today()
            },
        )
```

Then you need to add this entry in your `conf/<env>/parameters` to ensure that the parameter will be properly recorded:

```
execution_date: ${execution_date}
```

## 3.5 Opening the UI

### 3.5.1 The mlflow user interface

Mlflow offers a user interface (UI) that enable to browse the run history.

### 3.5.2 The kedro-mlflow helper

When you use a local storage for kedro mlflow, you can call a `mlflow cli command` to launch the UI if you do not have a `mlflow tracking server configured`.

To ensure this UI is linked to the tracking uri specified configuration, `kedro-mlflow` offers the following command:

```
kedro mlflow ui
```

which is a wrapper for `kedro ui` command with the tracking uri of the `mlflow.yml` file.

Opens `http://localhost:5000` in your browser to see the UI after calling previous command.

## INDICES AND TABLES

- `genindex`
- `modindex`
- `search`