

Assignment: ML Lifecycle

Tracking, reproducibility and deployment using [MLflow](#)

March 13, 2023, Revision 2

1 INTRODUCTION

In this assignment you will extend the code for your wind power forecasting system you made previously to go through some of the steps in the data analytics lifecycle ([Figure 1.1](#)). This includes:

1. Organized experiment tracking using the [MLflow](#) library and tracking server
2. Deploy your model to the cloud such that others can use its predictions

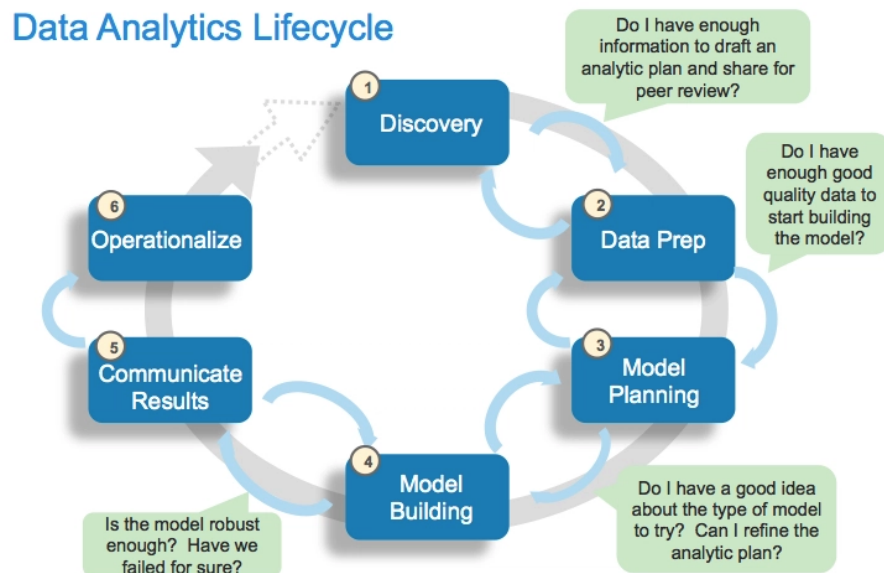


Figure 1.1: The data analytics life cycle. MLflow addresses tracking (2,3,4), reproducibility (4,5), and deployment/governance (6)

We increase the complexity of the system you will build, asking you to create several components that interoperate and communicate over the network. We're getting you further out of

the comfort-zone of your own local laptop environments. This will add some friction, but the purpose of this assignment is two-fold:

1. Familiarize you with a ML lifecycle tool, and
2. Expose you to the complexity/messiness of building scalable systems (motivating the myriad of tools for this).

This assignment also requires you to use tools or systems that you might not be familiar with from previous courses: Setting up a virtual machine (VM) with a public IP address, installing/-configuring software in a linux environment, managing long-running processes, automating recurrent tasks, doing basic networking configuration, version control using git, etc. Most of these sound more complicated than they are, but it can be a daunting task if you're new to it. A healthy dose of self-study is probably needed, and the TAs are your new best friends.

Again, we highly recommend the student-run MIT course [The Missing Semester of Your CS Education](#), which cover many of these aspects.

2 REQUIREMENTS AND HAND-IN

This section describes the (fairly loose) functional requirements of your system. The assignment is split into 4 parts: (1) Tracking, (2) Packaging and Reproducibility, (3) Deployment and (4) Model Governance.

We ask you to describe your design choices and their trade-offs (even if some of those decisions are defined by the assignment). In addition, each part contains one or more discussion points that we ask you to reflect on in your report.

The maximum length for the report handed in is five pages including figures.

2.1 Track your experiments and model governance

- Redo some of your experiments from Assignment 1, tracking the hyper-parameters and metrics using MLflow tracking.
- *Explain what you are logging and why.*

2.2 Package your setup and make it reproducible

- Your code should be able to automatically download data, train your model and repeat this every X days. Recall that you should always save the best model, if it is better than previously saved ones. Update your training pipeline so it:
 - Logs *metrics* every time a model is trained to a *remote* MLflow tracking server¹
 - Saves the final model using MLflow's API under the relevant MLflow experiment run on the remote tracking server.
- Package your code as an MLflow project and push it a repository you have created under [ITU's GitHub Enterprise server](#)² repo, such that your model can be [trained using a single command](#), e.g:

```
mlflow run git@github.itu.dk:LSDA-Spring-2023/mlflow-example.git -P alpha=0.5.
```
- *Discuss how each file in your repo affects the reproducibility of your project.*
- *Which steps does MLflow go through in order to run your project?*
- *Assume that you find your model to be producing low quality forecasts. What steps would you need to take to roll back the model to a previous version? How about if you did not use the model registry?*

¹We recommend starting with local tracking, then try our public one <http://training.itu.dk:5000/>. Optionally use the [Azure Machine Learning Studio](#) or setup your [own MLflow tracking server on your VM](#)

²Ensure that you are connected to the ITU Wi-Fi or using the [Forti client VPN](#). Login using your ITU credentials, and [add your public SSH key such that you can push/pull without passwords/https to your repo](#)

2.3 Deploy your model in the cloud

- Procure a virtual machine with a public IP address³.
- Setup miniconda and mlflow
- Setup a recurring job (e.g. using [cron](#)) that:
 - Runs your retraining project from the git repo⁴
 - Serves the saved model⁵
- *How does the continuous logging of metrics of your retraining script help to identify data drift?*
- *Assume that your system had an extreme increase in usage. How would you scale your system to handle the large number of requests?*

3 MLFLOW AND YOUR ASSIGNMENT

MLflow covers three main areas: [Experiment tracking](#), [reproducibility](#), and [model deployment](#).

Prior to working on this assignment, we highly recommend you go through the [MLflow Tutorial](#) playing around with saving different metrics and models.

3.1 Experiment tracking

3.1.1 Method A: Start with local tracking

When developing ML models, you go through a lot of experiments, trying many different combinations of models, hyper-parameters and feature extraction. This means you need to keep track of a lot of metrics, and how you created each metric. This is what MLflow Tracking is trying to solve.

In MLflow each *experiment* can consist of many *runs*, and in each run you can log *parameters*, *metrics*, *tags* and *artifacts*. To get a better explanation of this, visit the [MLflow Tracking introduction](#)

A simple experiment could look like this (you should use variables instead of constants like 0.5 though):

```
1 import mlflow
2 with mlflow.start_run():
3     mlflow.log_param("Param one", 1)
4     mlflow.log_metric("Accuracy", 0.5)
```

By default, MLflow will log your runs to the default experiment, and will save all your runs to your local file system in the folder `mlruns`. After running the example above, the directory structure could look like this:

```
1 mlruns/
2   0/
3     1f880fec49d64bffa1fdd4e7600f7c5b/
4       artifacts/
5         meta.yaml
6         metrics/
7           Accuracy
8         params/
9           Param one
10        tags/
```

³We recommend using Azure, see [subsubsection 3.3.3](#). If you run out of credits on Azure, contact us and we will figure something out.

⁴To start long-running processes that don't stop when you terminate your ssh connection, have a look at the [nohup command](#). To install new services (like MinIO), you might want to use [docker](#). Think about how these steps could be automated.

⁵You might need to kill the previous job. Have a look at [killall](#) (e.g. `killall -g mlflow`)

```

11         mlflow.source.git.commit
12         mlflow.source.name
13         mlflow.source.type
14         mlflow.user
15     meta.yaml

```

The experiment id is 0, the run id is 1f880fec49d64bffa1fdd4e7600f7c5b and everything is just saved as text files. Try opening some of these files in your text editor. When running MLflow from within a git repository, the current commit is also saved as a tag, making it easier to recreate experiments.

MLflow also have (more-or-less experimental) [autologging](#) features, which you are welcome to try out.

3.1.2 *Method B (Recommended): Use MLflow with ITU's tracking server instead of local folders*

In addition to tracking your experiments on your local computer, you can also use a public tracking server. We have set up such a server at <http://training.itu.dk:5000/>. This tracking server stores runs in a [PostgreSQL](#) database instead of the local file system. For the tracking of artifacts using the ITU server to work, it needs the AWS credentials for training:

```

1 # Setting the MLflow tracking server
2 mlflow.set_tracking_uri('http://training.itu.dk:5000/')
3 # Setting the required environment variables
4 os.environ['MLFLOW_S3_ENDPOINT_URL'] = 'http://130.226.140.28:5000'
5 os.environ['AWS_ACCESS_KEY_ID'] = 'training-bucket-access-key'
6 os.environ['AWS_SECRET_ACCESS_KEY'] = 'tqvdsSEdNBWTDuGkZYVsRKnTeu'

```

MLflow is not suitable as a competition framework, and is not really intended for many users. There is no authentication and results are not validated. This means that other users can delete your runs/experiments (possibly by mistake) and it is trivial to log fake metrics. Do not use the public tracking server for results that you cannot afford to lose. Take it with a grain of salt.

3.1.3 *OPTIONAL Method C: Use MLflow with Azure*

Microsoft Azure (See [Azure](#) section for how to get access) provides a similar interface to the MLFlow UI which can [be used as a tracking server](#). To use that, you can:

1. [Create Azure Machine Learning workspace](#)
2. Go in the ML workspace you just created, download config.json file and store it locally together with your code.
3. Launch studio
4. Run your experiments locally (or on Azure VM, see below sections to set up the VM to be able to run mlflow there) and you should be able to see the results in the Experiments tab on Microsoft Azure Machine Learning Studio.

3.2 Making your setup reproducible using MLflow "projects"

[MLflow Projects](#) is a standardized way to encapsulate an experiment in a reproducible manner. This is done by specifying all the dependencies as either a conda or docker environment. Here we will only talk about the conda approach.

An MLflow project consists of:

- The code you want to package, including the data for your experiments
- An environment file specifying the dependencies for the code. In this case a YAML file with the conda environment.
- An MLproject file specifying which environment file to use, and the entry points to the code.

We've made a [small project](#) that shows an example of polynomial regression, which we will use to show how MLflow projects work. See appendix for more detailed explanation.

3.3 Deploying your model

3.3.1 Laptop deployment

Once you have a model you are satisfied with, you can log and deploy it using MLflow Models model format. There's detailed descriptions of the deployment options in [the documentation](#).

As an example, have a look at [this repo](#). This is an MLflow Project that trains a CAISO model for electricity load forecasting and saves it to the local filesystem. The saved model can then be distributed and deployed using MLflow Models.

To use the saved artifacts, we need to clone the repo to your filesystem before running (otherwise the model will just be saved in a temporary folder).

```
1 git clone git@github.com:NielsOerbaek/PolyRegExample.git
2 cd caiso-mlflow
3 mlflow run .
4 mlflow models serve -m <name of model folder>
```

The model will now be served on your local machine. You can then query your model by opening another terminal and running:

```
1 curl http://127.0.0.1:5000/invocations -H 'Content-Type: application/json' -d '{
2   "columns": ["Time"],
3   "data": [["2021-04-15T20:00:00"]]
4 }'
```

You will then get an answer like `[21.492166666666666]`, meaning that your model thinks that the electricity demand in Orkney at 8 PM on Saturday the 1st of May 2021 will be 21.49 MW.

3.3.2 To the cloud!

Deploying to your own machine might seem a bit roundabout. The interesting thing comes when deploying to a server and can be queried by many users.

You can look into deploying your model to a cloud-based virtual machine with a public IP address. One such option is to use Microsoft Azure. MLflow has methods for deploying directly to Azure, but we've had mixed experiences with it. Again, have a look at [the documentation](#).

If you already have a Virtual Machine on Azure (if not, follow section [Azure](#) steps 1-2) you can ssh into it and serve your model:

1. Push your MLflow project to Github if you develop your project locally. In case you did your MLflow stuff (running experiments/training model) on Azure VM, you can skip steps 3-6.
2. ssh into your VM on Azure `ssh <username>@<Public IP address>`. You can find VM Public IP address in the Overview.
3. [Install miniconda on the VM](#).
4. Exit VM and ssh again so that conda environment is activated.
5. Get MLflow from the community "conda-forge" channel:

```
1 conda install -c conda-forge mlflow
```

6. Download your model to the VM using the GitHub repository
7. [Open port 5000](#)⁶. Skip Create a network security group, go to your Azure Resource group and select existing Network security group.
8. Serve the model:

```
1 mlflow run <name of folder with MLproject file in>
2 mlflow models serve -m <name of model folder> -h 0.0.0.0 -p 5000
```

9. Query served model from your local machine:

⁶Port 5000 is default for MLflow

```

1 curl http://<Public VM IP address>:5000/invocations -H 'Content-Type: application/
  json' -d '{
2   "columns": ["Time"],
3   "data": [["2020-11-14T20:00:00"]]
4 }'
```

You will then get an answer like [16.07111111111111], meaning that your model thinks that the electricity demand in Orkney at 8 PM on Saturday the 14th of November 2020 will be 16.07 MW.

3.3.3 Azure

One option for getting a VM with a public IP is to use Microsoft Azure. Azure is enterprise-oriented and can be quite complicated, so you need to go through some steps that might feel like a bit over the top for this task. But on the plus-side, you should all be eligible for [free student credits](#), giving you some credits to experiment with.

You can follow these guides to get you started:

1. [Setup your student account](#)
2. [Setup a VM with a public IP](#):
 - a) You could pick a VM residing in a region powered by much sustainable energy. Like Norway or Sweden.
 - b) In the Image choose Ubuntu Server 20.04.
 - c) You should pick a cheaper image. Like "B2s".
 - d) We recommend you re-use an existing public key located on your laptop if you have one. If not, generate one. There's plenty on guides online for this, again we recommend the ["CS Missing Semester" on SSH keys and remote machines](#). Note that Azure might require an RSA key, so use `-t rsa` instead of `-t ed25519` if using the missing semester guide to create SSH keys.
 - e) If you generate a new ssh key pair, you need to download it to your laptop, such that your chosen ssh client can find use it.
 - f) [Create a network security group to open the ports you need, port 5000 is default for the MLflow serve command](#),. Should be possible during the creation of the VM, if done later, click the previous link.

After these steps you should be able to ssh into your VM and serve your model or MLflow UI⁷. If you are inexperienced using SSH, we recommend you read through ["CS Missing Semester" on SSH keys and remote machines](#).

If you for some reason do not have student credits, contact us and we will figure something out.

3.3.4 Trying out your deployed model

We've made a little webpage where you can try to use your deployed model in something that resembles a real use case a little more.

On <https://orkneycloud.itu.dk/mlflow/> you can insert the url for your invocation end-point and it will fetch the newest weather forecasts, use your model to make predictions, and draw a little graph.

⁷MLflow also has methods for deploying directly to Azure, but I've had mixed experiences with it in the past. Again, have a look at [the documentation](#).

4 DATA

For this assignment we have frozen the dataset to make comparisons easier. This means that the data will not be fetched from the influx database, but simply loaded from the attached json-file. The json file can be loaded into a pandas dataframe by using:

```
1 df = pd.read_json("path/to/file.json", orient="split")
```

The data format is the same as for Assignment 1; the only difference being that the generation and wind dataframe have been joined by an outer join. Thus you only need to load a single dataframe as your dataset.

The attached dataset covers 180 days of data. Below is the output of `df.info()`

```
1 <class 'pandas.core.frame.DataFrame'>
2 DateTimeIndex: 254967 entries, 2020-05-15 12:55:00 to 2020-11-11 12:54:00
3 Data columns (total 7 columns):
4 #   Column          Non-Null Count  Dtype
5 ---  -
6 0   ANM              254967 non-null float64
7 1   Non-ANM          254967 non-null float64
8 2   Total            254967 non-null float64
9 3   Direction        1379 non-null  object
10 4   Lead_hours       1379 non-null  float64
11 5   Source_time      1379 non-null  datetime64[ns]
12 6   Speed            1379 non-null  float64
13 dtypes: datetime64[ns](1), float64(5), object(1)
14 memory usage: 15.6+ MB
```

5 APPENDIXES

5.1 Using MLflow Projects: PolyRegExample repo

We've made a [small project](#) that shows an example of polynomial regression, which we will use to show how MLflow projects work. The main experiment, in which we simulate some data⁸ and model it using different degrees of polynomial regression, is defined in `experiment.py`:

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.preprocessing import PolynomialFeatures
3 from sklearn.pipeline import Pipeline
4 from matplotlib import pyplot as plt
5 import numpy as np
6
7 import sys
8 num_samples = int(sys.argv[1]) if len(sys.argv) > 1 else 100
9
10 def get_ys(xs):
11     signal = -0.1*xs**3 + xs**2 - 5*xs - 5
12     noise = np.random.normal(0,100,(len(xs),1))
13     return signal + noise
14
15 X = np.random.uniform(-20,20,num_samples).reshape((num_samples,1))
16 y = get_ys(X)
17
18 plt.scatter(X,y,label="data")
19
20 for degree in range(1,4):
21     model = Pipeline([
22         ("Poly", PolynomialFeatures(degree=degree)),
23         ("LenReg", LinearRegression())
24     ])
```

⁸The data is simulated using the polynomial function $f(x) = -0.1x^3 + x^2 - 5x + 5 + \epsilon$, where ϵ is Gaussian noise.

```

25     model.fit(X,y)
26     plotting_x = np.linspace(-20,20,num=50).reshape((50,1))
27     preds = model.predict(plotting_x)
28     plt.plot(plotting_x, preds, label=f"degree={degree}")
29
30 plt.legend()
31 plt.show()

```

The Conda environment for this experiment is specified in PolyReg.yaml:

```

1 name: PolyReg
2 channels:
3   - defaults
4 dependencies:
5   - scikit-learn>0.23
6   - numpy>1.19
7   - pip>20
8   - python=3.8
9   - pip:
10     - mlflow
11     - matplotlib>3

```

Finally, the MLproject-file specifies how to run the project:

```

1 name: PolyReg
2
3 conda_env: PolyReg.yaml
4
5 entry_points:
6   main:
7     parameters:
8       num_samples: {type: int, default: 100}
9     command: "python experiment.py {num_samples}"

```

With these three things in place, you can run the experiment using `mlflow run <path to project>`. So if the project is located on your computer, you can navigate to the directory and do:

```

1 mlflow run .

```

Because this project is hosted as a git repository, you can simply do:

```

1 mlflow run git@github.com:NielsOerbaek/PolyRegExample.git

```

This will fetch the project, resolve the environment, and run the main entry point with the default parameters.

If you want to run the experiment with 500 samples, instead of the default 100, you can do:

```

1 mlflow run git@github.com:NielsOerbaek/PolyRegExample.git -P num_samples=500

```