



# Hands-on MLFlow

Managing the end-to-end machine learning lifecycle in practice.

[sit.org](https://sit.org)



# Toolkit

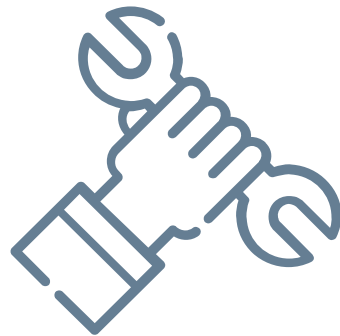
Python,

[Google Colab](#),

[Ngrok](#),

[GitHub](#),

ML libraries: scikit-learn, Pandas, Numpy, Matplotlib





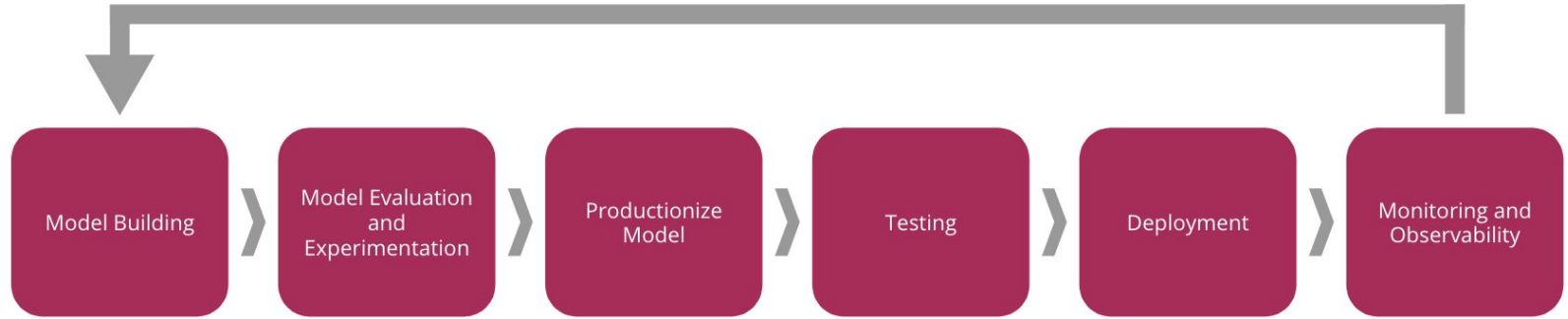
01



# Introduction to MLFlow

The Data Science Project Life  
Cycle, and how MLFlow fits in it

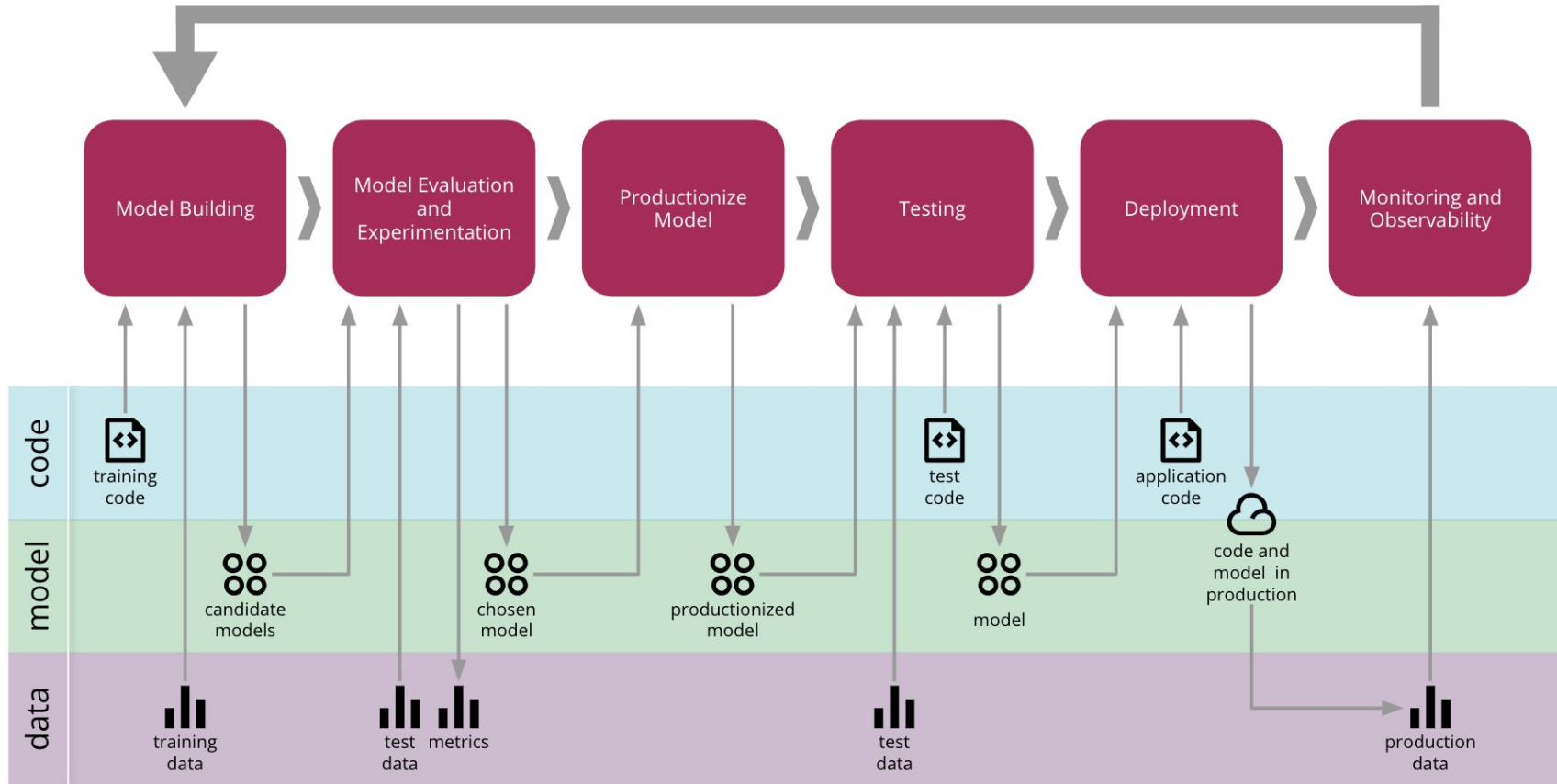
# The Data Science Process



[Source](#)



# The Data Science Process



# Challenges of Machine Learning Development

## Goal

= Optimizing metrics.  
Metrics can always be improved and can degrade if models are not updated.

## Quality

Depends on data, model, training parameters and code => we need to track all of these.



## Tools

Data Scientists use various tools/libraries for data preparation and modelling  
=> experiment **tracking** and **sharing** is challenging.

## Deployment

A same model might need to be deployed in various environments.

# MLFlow Solutions

---



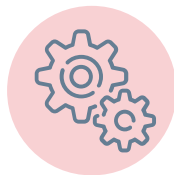
## MLFlow Tracking

Record and query experiments: code, data, config, and results



## MLFlow Projects

Package data science code in a format to reproduce runs on any platform



## MLFlow Models

Deploy machine learning models in diverse serving environments



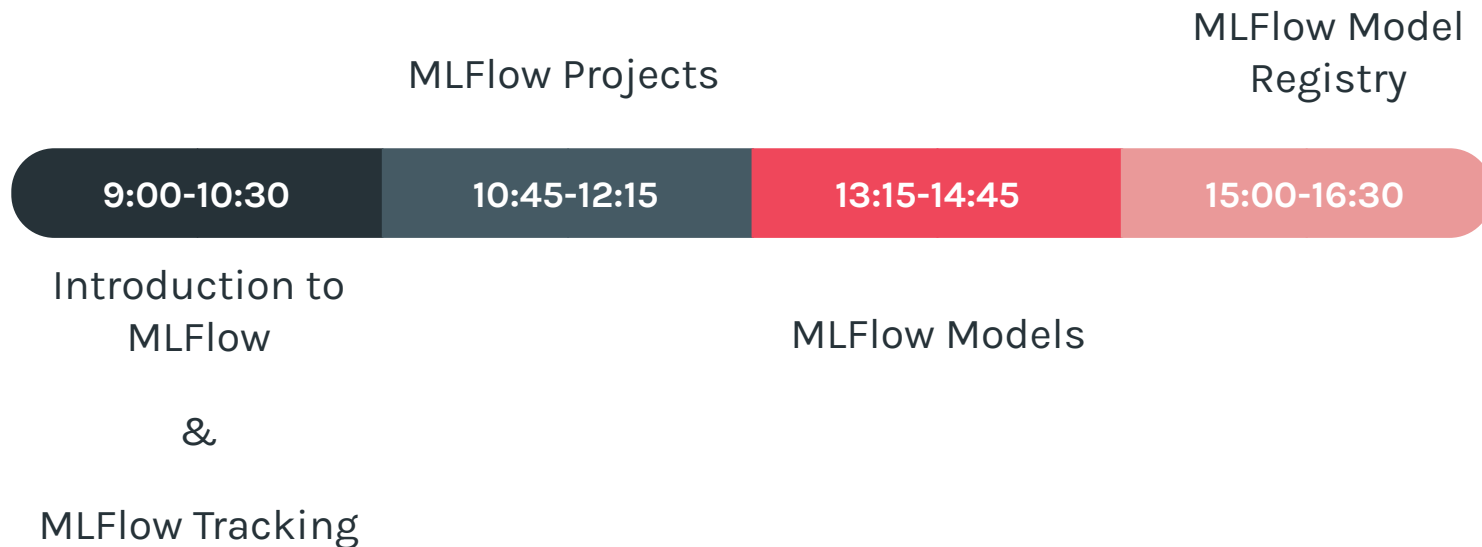
## MLFlow Registry

Store, annotate, discover, and manage models in a central repository

Source



# Outline of the day







02



# MLFlow Tracking

log parameters, code versions,  
metrics, and artifacts when running  
your machine learning code.

# MLFlow Tracking Concepts

---

In MLFlow tracking, a **run** is the execution of some piece of data science code. For each run, you can log:

## Code Version

Git commit hash of the run

## Start & End Time

Start and end Time of the run

## Source

Name of the file to launch the run, or the project name and entry point for the run

## Parameters

Key-value input parameters of your choice. Both keys and values are strings.

## Metrics

Key-value metrics, where the value is numeric and can be updated throughout the course of the run

## Artifacts

Output files in any format (images, pickled models, data files...).

# MLFlow Tracking Python API

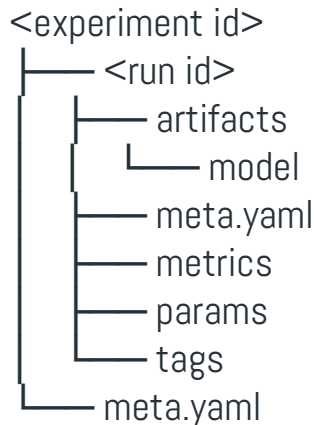
```
import numpy as np
from sklearn.linear_model import LinearRegression
import mlflow
```

```
X = np.array([[1, 1], [1, 2], [2, 2], [2, 3]])
y = np.dot(X, np.array([1, 2])) + 3
```

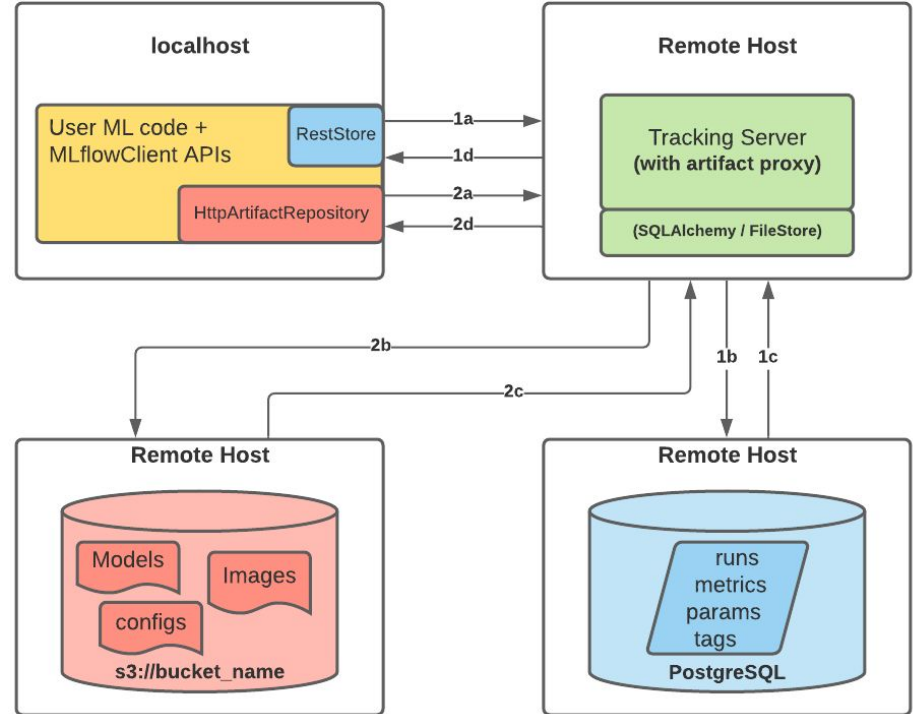
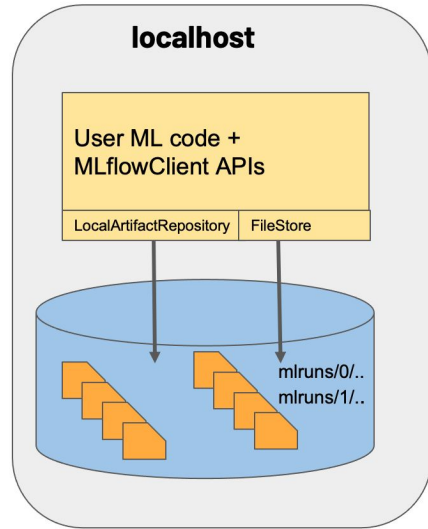
```
with mlflow.start_run():
    mlflow.sklearn.autolog()
    reg = LinearRegression().fit(X, y)
```

- ❖ Creates a new **run**
- ❖ Logs all **parameters** for LinearRegression
- ❖ Logs training **Metrics**
- ❖ Logs **Model**

By default, all this information is stored in the folder **mlrun**:



# MLFlow Tracking Storage Configurations



# MLFlow Tracking UI

Once you have executed a **run** in MLFlow tracking, you can explore all the entities you logged using the MLFlow UI:

The screenshot displays the MLFlow Tracking UI interface. At the top, the 'mlflow 1.25.1' logo is visible alongside 'Experiments' and 'Models' tabs. The 'Experiments' tab is active, showing a search bar and a list of experiments with 'Default' selected. A notification banner states: 'Track machine learning training runs in experiments. Learn more'. Below this, the 'Experiment ID: 0' is shown. The 'Description' section includes buttons for 'Refresh', 'Compare', 'Delete', 'Download CSV', 'Start Time' (set to 'All time'), and a search bar with the query 'metrics.rmse < 1 and params.model = "tree"'. A table titled 'Showing 1 matching run' contains one entry:

	Start Time	Duration	Run Name	User	Source	Version	Models	Metrics
<input type="checkbox"/>	1 minute ago	2.2s	-	marie	ipykernel	-	sklearn	2.220

A 'Load more' button is located at the bottom of the table.

[Source](#)

# MLFlow Tracking: experiments

---

- ❖ Minimal Example
- ❖ Regression problem
- ❖ Classification Problem
- ❖ Exercise: Tracking with your own code

# Tracking Livecoding!

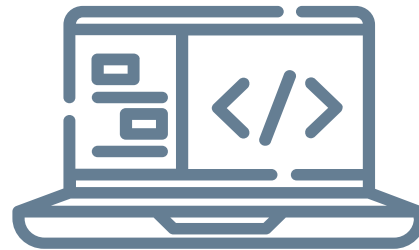
---

Go to Google Colab

The notebooks are available here:

[LC0\\_MLFlow\\_Tracking\\_MinExample](#)

[LC1\\_MLFlow\\_Tracking\\_Regression\\_SwissHousing](#)





03



# MLFlow Projects

package data science code in a  
reusable way



# Overview

---

- ❖ Concepts and motivation
- ❖ Documentation
- ❖ Build and share an MLflow Project with Git

# MLFlow Projects

---



**MLFlow  
Projects**

Package data science  
code in a format to  
reproduce runs on any  
platform



Reproducibility

...and how to not become  
another statistic in the  
**reproducibility crisis**

# MLFlow Projects - Virtual environments

---

## Reproducibility:

- Data (processing methods)
- Code (package versions, logging)

## Other benefits:

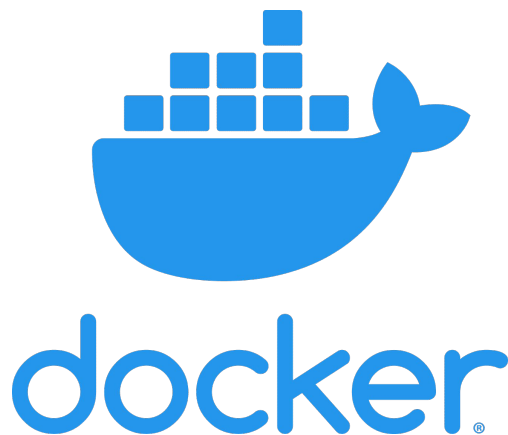
- Platform-agnostic
- Automation
- Troubleshooting

# MLFlow Projects - Virtual environments

---

Reproducibility:

- Code (package versions, logging)

The Conda logo is displayed in a bright green color. It features the word "CONDA" in a bold, sans-serif font. The letter "C" is stylized with a circular pattern of white lines on a green background, resembling a DNA helix or a network structure.

# MLFlow Projects - *MLProject* YAML file

---

## *MLProject.yaml*

- **Entry points**
  - Commands run within project
  - ***Need at least one!***
- **Parameters**
  - Variables to be passed to entry point
    - *File paths, options, etc.*
- **Environment**
  - Conda, Docker, venv, etc.

# MLFlow Projects - *MLProject* YAML file

*MLProject.yaml*

Project name

```
name: My Project
```

```
conda_env: my_env.yaml
```

Environment type  
& config file

```
# Can have a docker_env instead of a conda_env, e.g.
```

```
# docker_env:
```

```
#   image: mlflow-docker-example
```

```
entry_points:
```

```
  main:
```

```
    parameters:
```

```
      data_file: path
```

```
      regularization: {type: float, default: 0.1}
```

Final command to be  
run in terminal

```
    command: "python train.py -r {regularization} {data_file}"
```

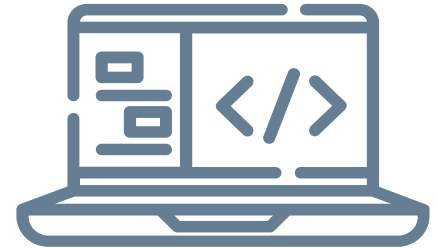
# Projects Livecoding!

---

Go to Google Colab

The notebook is available here:

[LC2\\_MLFlow\\_Projects\\_Regression\\_SwissHousing](#)





04



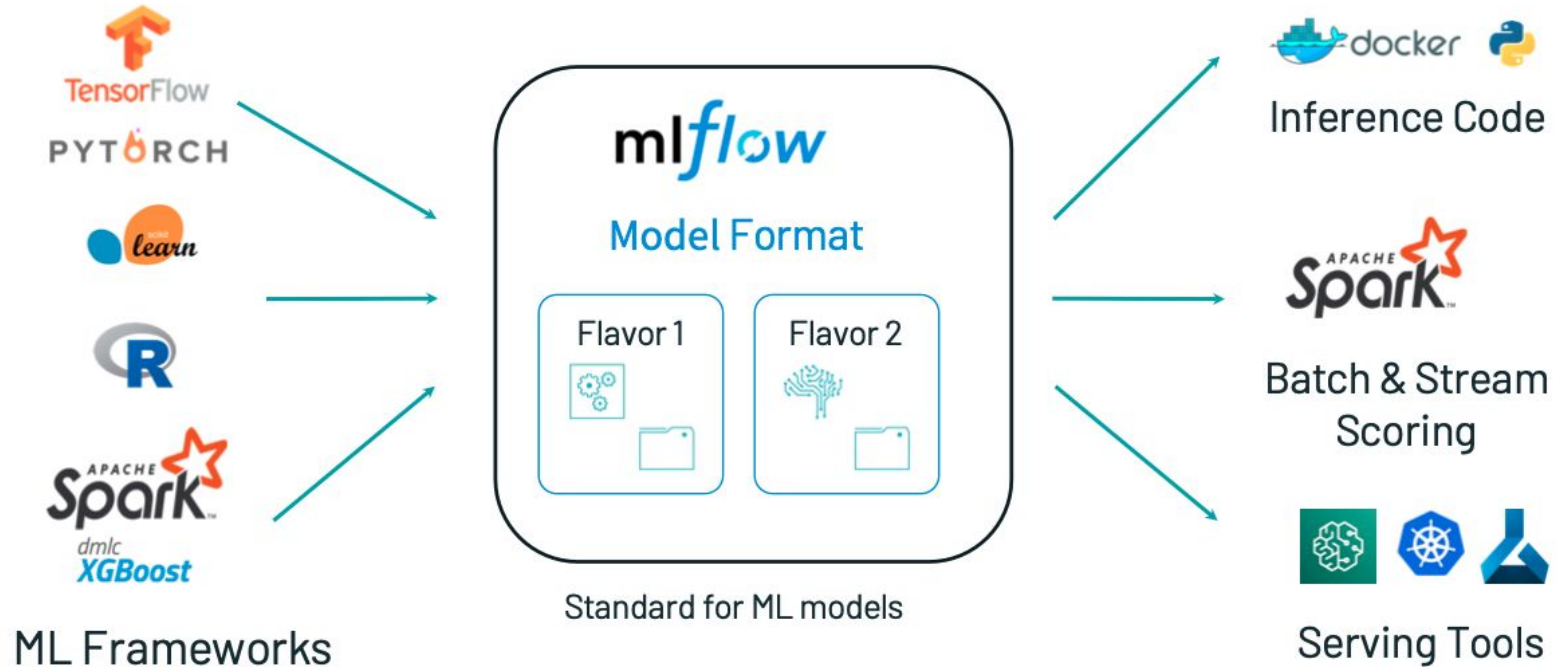
# MLFlow Models

Package Machine Learning Models  
and deploy them



# MLFlow Models Motivation:

Deploy machine learning models in diverse serving environments



Source

# MLFlow Models Documentation:

---

Deploy machine learning models in diverse serving environments

- ❖ Storage Format, Model Flavors
- ❖ Model API, Model Evaluation
- ❖ Deployment Tools

[Source](#)

# Example MLFlow Model (option 1):

By default, all this information is stored in the folder model:

model/  
├── **MLmodel**  
├── model.pkl  
├── conda.yaml  
├── python\_env.yaml  
└── requirements.txt



artifact\_path: model

flavors:

**python\_function:**

env: conda.yaml

loader\_module: mlflow.sklearn

model\_path: model.pkl

python\_version: 3.7.13

**sklearn:**

code: null

pickled\_model: model.pkl

serialization\_format: cloudpickle

sklearn\_version: 1.0.2

mlflow\_version: 1.26.1

model\_uuid: 07d32968a75345a9b28a8cd8f0b3f093

run\_id: f5d7ac67f6454523a880adc52f72242b

utc\_time\_created: '2022-06-09 11:12:27.594199'

Usable by:

Any tool that can  
run Python (Docker,  
Spark, etc)

Tools that under-  
stand Scikit-Learn  
model format

## Built-In Model Flavors

# Model Scikit-Learn Flavor Example



→ `mlflow.sklearn.log_model(pipeline_rf, "model")`

Train a model



model  
format



Flavor 1:  
python\_function



```
predict =  
mlflow.pyfunc.load_model(...)  
predict(pandas.input_dataframe)
```



Flavor 2:  
sklearn



```
model =  
mlflow.sklearn.load_model(...)  
model.predict(data['X_test'])
```

# Model Scikit-Learn API and Evaluation

## Model API

Save and load MLflow Models in multiple ways:

- 1) `mlflow.sklearn` contains `save_model`, `log_model`, and `load_model` functions for scikit-learn models.
- 2) `mlflow.models.Model` class to create and write models.

## Model Evaluation

`mlflow.evaluate()`\* API to evaluate MLflow Model performance

\*currently supports the `python_function (pyfunc)` model flavor for classification and regression tasks, evaluation results are logged to `MLflow Tracking`

## Examples:

```
import mlflow
import mlflow.sklearn

iris = load_iris()
sk_model = tree.DecisionTreeClassifier()
sk_model = sk_model.fit(iris.data, iris.target)
# set the artifact_path to location where experiment
artifacts will be saved

...

# log model
mlflow.sklearn.log_model(sk_model, "sk_models")

with mlflow.start_run() as run:
    model_info = mlflow.sklearn.log_model(model, "model")
    result = mlflow.evaluate(
        model_info.model_uri,
        eval_data,
        targets="label",
        model_type="classifier",
        dataset_name="adult",
        evaluators=["default"],
    )
```

# MLFlow Models Deployment:

## Built-In Deployment Tools

- on a local machine
- to several production environments

### Deploy MLflow models

Deploy a python\_function model on Microsoft Azure ML

Deploy a python\_function model on Amazon SageMaker

Export a python\_function model as an Apache Spark UDF

Not all deployment methods are available for all model flavors.

## Example:

- Serve a model saved with MLflow.
- Launch a webserver on the specified host and port. Supports models with the python\_function or crate (R Function) flavor.
- Make requests to POST /invocations in pandas split- or record-oriented formats.

```
$ mlflow models serve -m  
runs:/my-run-id/model-path &
```

```
$ curl http://127.0.0.1:5000/invocations -H  
'Content-Type: application/json' -d '{  
  "columns": ["a", "b", "c"],  
  "data": [[1, 2, 3], [4, 5, 6]]  
}'
```

# MLFlow Models Summary:

---

- ❖ Packaging format for ML Models
  - Any directory with MLmodel file
- ❖ Defines dependencies for reproducibility
  - Conda environment can be specified in MLmodel configuration
- ❖ Model creation and loading utilities
  - `mlflow.<model_flavor>.save_model(...)` or `log_model(...)`
  - `mlflow.<model_flavor>.load_model(...)`
- ❖ Deployment APIs
  - CLI/Python/R/Java
  - `mlflow models [OPTIONS] COMMAND [ARGS]...`
    - `mlflow models serve [OPTIONS [ARGS] ....`
    - `mlflow models predict [OPTIONS [ARGS] ...`

# Models Livecoding!

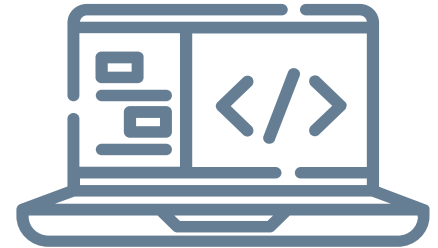
---

Go to Google Colab

The notebooks are available here:

LC3\_MLFlow\_Models\_Regression\_SwissHousing

LC3\_MLFlow\_Models\_Exe







05



# MLFlow Model Registry

get familiar with the centralized  
model store

# MLFlow Model Registry – Data Stores



## Artifact store

Model files, plots,  
config files, etc.

(./mlruns, S3  
bucket, etc.)

## Backend store

Metrics, IDs,  
parameters, etc.

(Local file, SQL DB,  
etc.)



Register different versions  
(*v.1, v.2, v.n*)



Register different stages  
(*staging, production, etc.*)

# MLFlow Model Registry – Model Store

---

- Model management API
  - From the beginning to the end of a project
- Full model lineage
  - Versioning, stage, MLFlow experiment metadata (run, date, metrics, etc.)
- Higher level annotations
- Easy accessibility for deployment



Default > Run ed089cf2d5d24b8880446761987b7ac2 ▾

Date: 2019-10-16 22:49:25

Source: train\_predict.py

Git Commit:

fcfefecd6f830c43dc8cdc048b83ac92fd8d5d7e

User: sid.murching

Duration: 6.2s

▾ Notes

None

▸ Parameters

▸ Metrics

▾ Tags

Name	Value	Actions
No tags found.		

Add Tag

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

▾ Artifacts

<div><div>▾ model</div><div><div> MLmodel</div><div> conda.yaml</div><div> tfmodel</div></div></div>	<div>Full Path: /tmp/swag/0/ed089cf2d5d24b8880446761987b7ac2/artifacts/model</div> <div>Size: 0B</div>	<div>Register Model</div>
--	--	---------------------------





**mlflow**GitHubDocs

**Default > Run ed089cf2d5d24b8880446761987b7ac2** ▾

Date: 2019-10-16 22:49:25

Source: train\_predict.py

Git Commit:  
fcfefcd6f830c43dc8cdc048b83ac92fd8d5d7e

User: sid.murching

Duration: 6.2s

▼ Notes

None

▶ Parameters

▶ Metrics

▼ Tags

Name

Value

Ne

Add Tag

Name

Value

Add

▼ Artifacts

▼ model

MLmodel

conda.yaml

▶ tfmodel

Full Path: /tmp/swag/0/ed089cf2d5d24b8880446761987b7ac2/artifacts/model  
Size: 0B

Register Model

Register Model

×

Once registered, the model will be available in the model registry and become public.

\* Model

+ Create New Model ▾

\* Model Name

Model A|

Cancel

Register

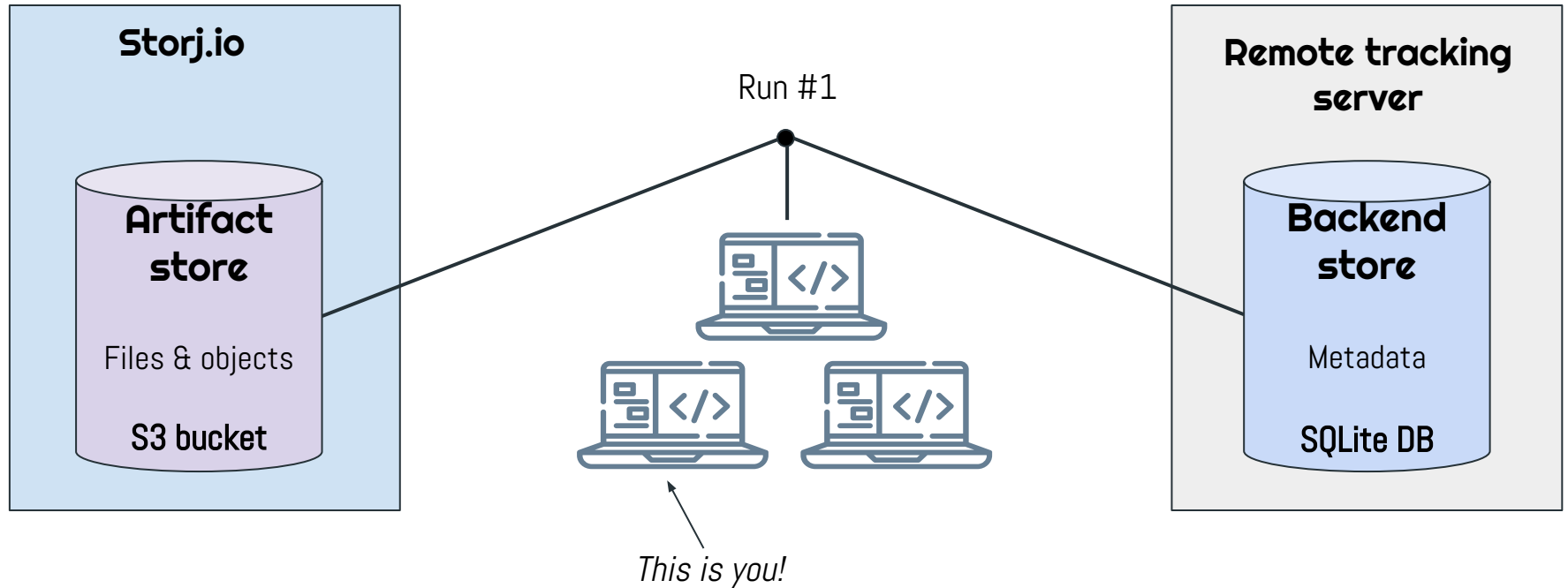


## Registered Models



Name ▾	Latest Version	Staging	Production	Last Modified ▾
<a href="#">Model A</a>	<a href="#">Version 1</a>	<a href="#">Version 1</a>	—	2019-10-16 22:51:19
<a href="#">Model B</a>	<a href="#">Version 1</a>	—	—	2019-10-16 22:51:52

# Registry Livecoding!





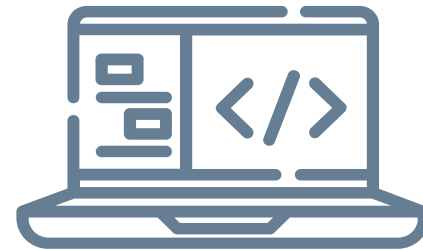
# Registry Livecoding!

---

Go to Google Colab

The notebooks are available here:

LC4 MLFlow Registry Regression SwissHousing





- [Continuous Delivery for Machine Learning](#)
- Databrick's 3-part tutorial: [Managing the Machine Learning Lifecycle using MLFlow](#)
- [MLFlow's official Documentation](#)



# SDS 2022 WORKSHOPS



SWISS CONFERENCE  
ON DATA SCIENCE

Workshop Feedback

**d+i**  
data innovation alliance