# MLOPs Zoomcamp Week 2 (part 1): Experiment Tracking

Gabriel

# 2.1 Experiment Tracking Intro

# Experiment Tracking

- ML experiment: the process of building an ML model
- Experiment run: each trial in an ML exp
- Run artifact: any file that is associated with an ML run
- Experiment metadata

# Tracking relevant information

- Source code
- Environment
- Data
- Model
- Hyperparameters
- Metrics
- ...

For:

- Reproducibility
- Organization
- Optimization

# MLflow

- Open source tool for ML lifecycle (python package)
- Four main modules:
  - Tracking
  - Models (special type of model)
  - Model Registry
  - Projects

# Allows you to keep track of..

- Parameters
- Metrics
- Metadata: tags (name of the developer, algorithm)
- Artifacts: visualization, dataset (but hard to scale), dictionary vectorizer
- Models

Also auto log:

- Source code
- Version of the code (git commit)
- Start and end time
- Author

# Install and run (with sqlite backend)

$pip install mlflow

$mlflow ui –backend-store-uri sqlite:///mlflow.db

```
(exp-tracking-env) → experiment_tracking git:(main) × mlflow ui --backend-store-uri sqlite:/
//mlflow.db
[2022-05-09 16:53:53 +0200] [6995] [INFO] Starting gunicorn 20.1.0
[2022-05-09 16:53:53 +0200] [6995] [INFO] Listening at: http://127.0.0.1:5000 (6995)
[2022-05-09 16:53:53 +0200] [6995] [INFO] Using worker: sync
[2022-05-09 16:53:53 +0200] [6996] [INFO] Booting worker with pid: 6996
```

# 2.2 Getting Started with MLflow

# In notebook file

- Check duration-prediction.ipynb

# 2.3 Experiment Tracking with MLflow

# 2.3 Tags

# 2.3 compare Xgboost tags

# Auto logging

## Automatic Logging

Automatic logging allows you to log metrics, parameters, and models without the need for explicit log statements.

There are two ways to use autologging:

1. Call `mlflow.autolog()` before your training code. This will enable autologging for each supported library you have installed as soon as you import it.
2. Use library-specific autolog calls for each library you use in your code. See below for examples.
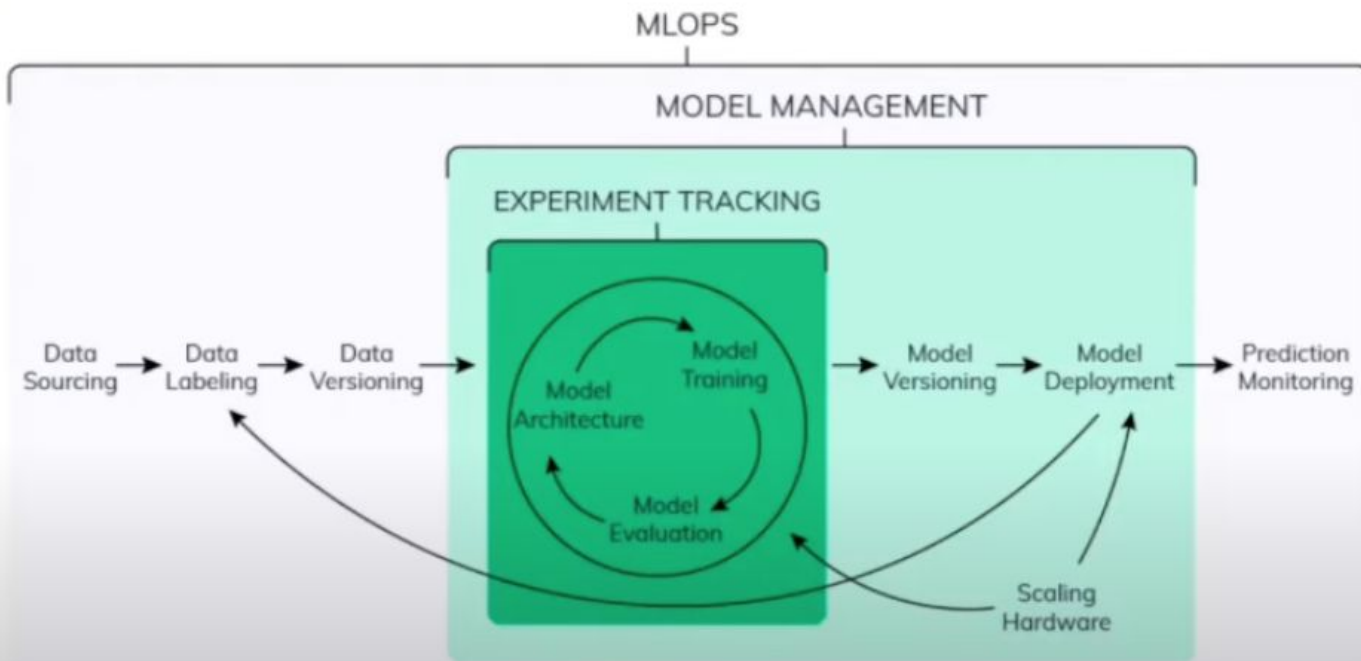
The following libraries support autologging:

- Scikit-learn
- TensorFlow and Keras
- Gluon
- XGBoost
- LightGBM
- Statsmodels
- Spark
- Fastai
- Pytorch

# 2.4 Model Management

# ML Lifecycle



Source: https://neptune.ai/blog/ml-experiment-tracking

# Model management

What's wrong with this?

- Error prone
- No versioning
- No model lineage

final_model     model_final_final     model_v1     model_v2_final_pr od

THIS_IS_THE_PRO D_MODEL     THIS_IS_THE_PRO D_MODEL_v2

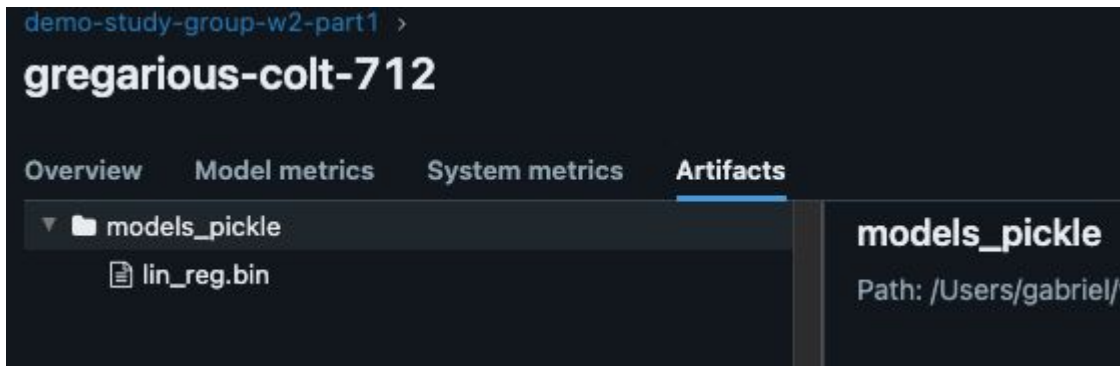# Log artifacts

```
# 2.4 log artifacts
 mlflow.log_artifact(local_path='./models/lin_reg.bin' ,artifact_path='models_pickle')
→ Not good, don't have instruction of how to use the model
```

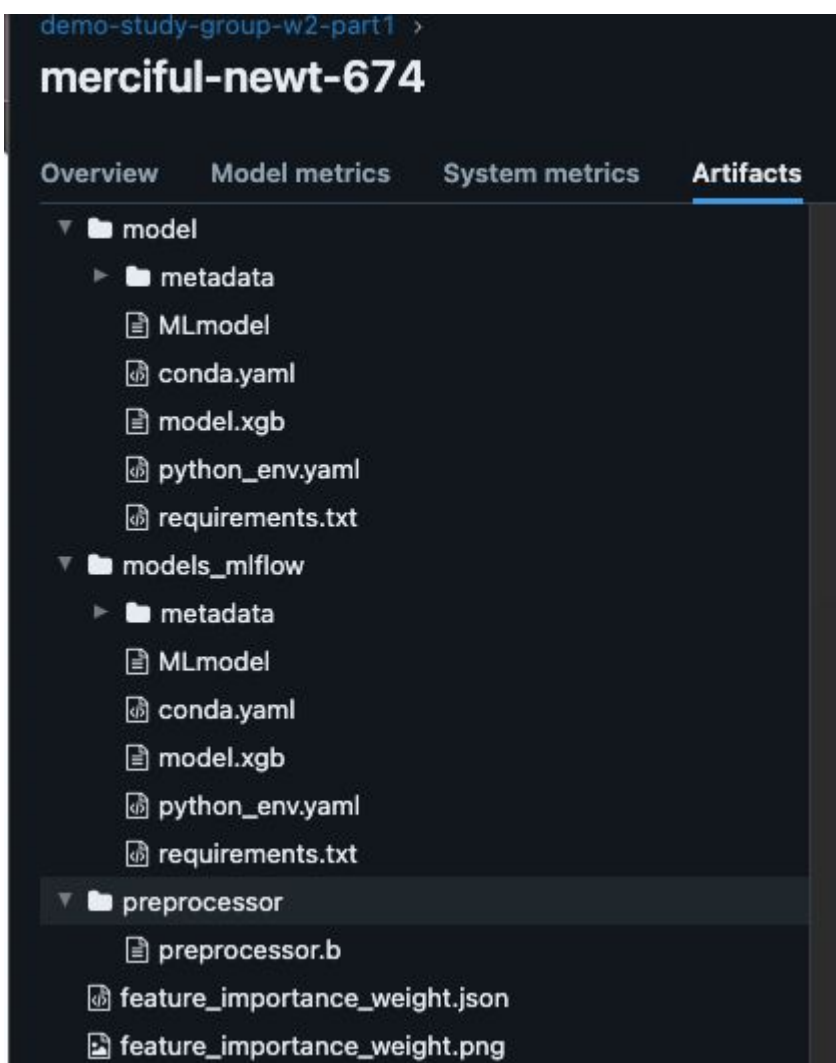# After logging the model, 2 ways to load

# Logged model, artifacts

Model: auto logging

Models_mlflow: manual logging

Preprocessor: artifacts for dictionary vectorizer

# Load model and make predictions

## Model schema

Input and output schema for your model. Learn more

| Name | Type |
| --- | --- |

No schema. See MLflow docs for how to include input and output schema with your model.

## Make Predictions

Predict on a Spark DataFrame:

```python
import mlflow
logged_model = 'runs:/b19e4d26363e4ce4a1c671b7ec
f83e9b/models_mlflow'

# Load model as a Spark UDF. Override result_typ
e if the model does not return double values.
loaded_model = mlflow.pyfunc.spark_udf(spark, mo
del_uri=logged_model, result_type='double')

# Predict on a Spark DataFrame.
columns = list(df.columns)
df.withColumn('predictions', loaded_model(*colum
ns)).collect()
```

Predict on a Pandas DataFrame:

```python
import mlflow
logged_model = 'runs:/b19e4d26363e4ce4a1c671b7ec
f83e9b/models_mlflow'

# Load model as a PyFuncModel.
loaded_model = mlflow.pyfunc.load_model(logged_m
odel)
```

# Recap of logging models in MLflow
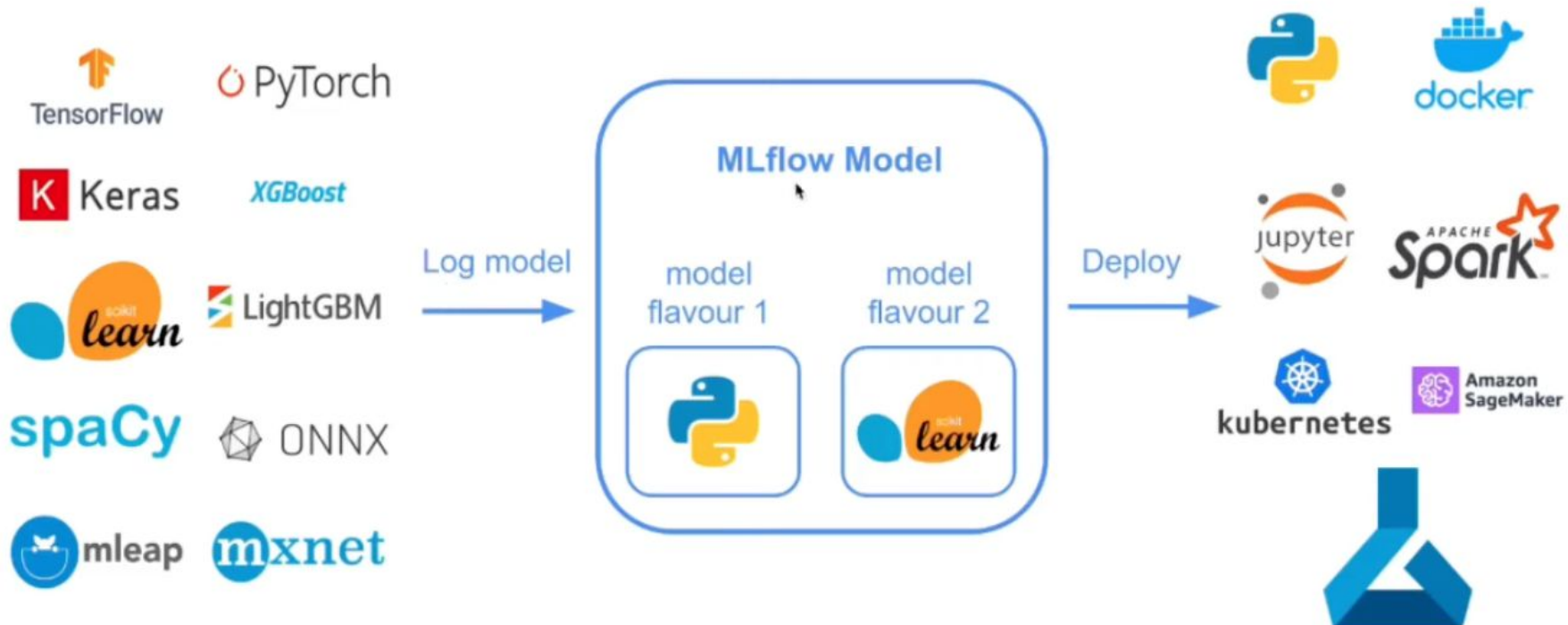
Two options:

- Log model as an artifact

```
mlflow.log_artifact("mymodel", artifact_path="models/")
```

- Log model using the method `log_model`

```
mlflow.<framework>.log_model(model, artifact_path="models/")
```
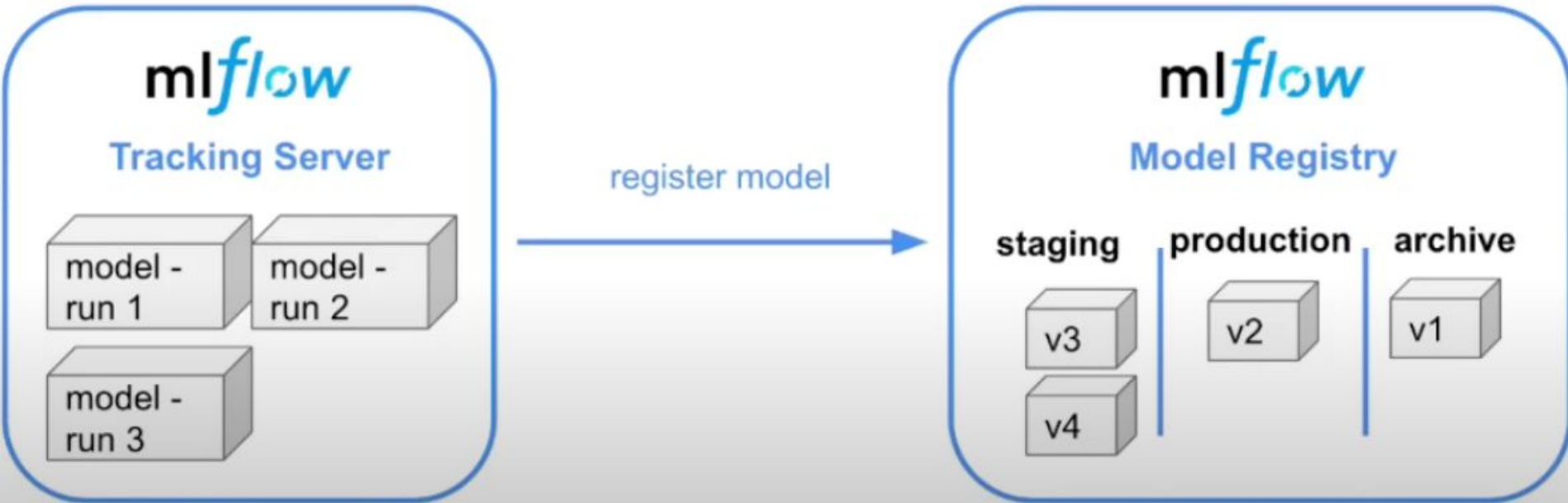
# MLflow Model Format

# 2.5 Model Registry

# Model Registry: version control

# **Registering the models**

- Register models (in artifact tab)
- Add to the group
- Check inside the registered models:
    - Lineage (which version links to which run)
    - Add tags (name and model)
- Check model-registry.ipynb

# MlflowClient Class

- A client of …
  - an MLflow Tracking Server that creates and manages experiments and runs.
  - an MLflow Registry Server that creates and manages registered models and model versions.
- To instantiate it we need to pass a tracking URI and/or a registry URI:

```python
from mlflow.tracking import MlflowClient

client = MlflowClient(tracking_uri="sqlite:///mlflow.db")
```