



MLflow Pipelines

Accelerating MLOps from Development to Production



Xiangrui Meng Principal Software Engineer Databricks



Jin Zhang Staff Software Engineer Databricks

ORGANIZED BY 😂 databricks

Problems



ML maturity levels

Almost every company wants to be a Data+Al company, but the majority are new to machine learning (ML).





3 common problems with scaling ML

Getting started	Iteration speed	Productionization
Mundane steps with lots of boilerplate	Slow, redundant iterative development	Handoff to production is very manual and bottlenecked



Explore Solutions



TensorFlow Extended (TFX)

Google-production-scale ML pipelines





TensorFlow Extended (TFX)

What fits:

- ML pipeline abstraction
- predefined steps
- embedded best practices

What doesn't fit:

 Designed for ML engineers to solve complex and large-scale ML problems, hard for data scientists to get started.

trainer = tfx.components.Trainer(module_file=os.path.abspath(_taxi_trainer_module_file), examples=transform.outputs['transformed_examples'], transform_graph=transform.outputs['transform_graph'], schema=schema_gen.outputs['schema'], train_args=tfx.proto.TrainArgs(num_steps=10000), eval_args=tfx.proto.EvalArgs(num_steps=5000)) context.run(trainer, enable_cache=True)



Apache Maven

DATA+AI SUMMIT 2022

"Why is no one using make for Java?"



Apache Maven

What fits:

- Easy to get started with predefined build lifecycle and standard directory layout.
- Standardized on POM for collaboration and sharing.
- Build cache for incremental build.

What doesn't fit:

- One "build" lifecycle might not fit all ML problem types.
- Not designed to solve ML problems, where data is essential.



Notebook vs. IDE

Do we have to choose one?



MLflow Pipelines



Start quickly with pipeline templates





Start quickly with pipeline templates

Clustering the 80% of ML problems

Built-in pipeline templates:

- regression
- classification
- •••

Simple to get started:

- standard directory layout
- mostly config-driven
- code within DS' comfort zone

pipeline.yaml

```
template: regression/v1
```

```
data:
type: parquet
path: datasets/taxi.parquet
```

```
target_col: fare_amount
```

```
metrics:
 primary: rmse
```

```
steps:
 transform:
```

•••



Iterate fast with the pipeline engine





Accelerate with opinionated dev workflows

Embedded best practices

- deterministic splits
- data profiles
- transformed feature names
- feature importance
- automatic MLflow tracking

Notebook + IDE

- Notebook to trigger pipeline execution.
- Notebook to view rich format output.
- Notebook to analyze data.
- IDE to develop modularized code.



...

Productionize without refactoring

Notebook is only used for triggering execution, not on the critical path.
 Whenever one finds a good model, it is ready to ship.

• Default test suites to prevent unexpected changes.

• Standard command-line interface to integrate with CI/CD.



mlflow Pipelines

Pipeline templates to get started quickly

Pipeline engine for accelerated development Opinionated structure to automate handoffs Feature Spotlights







Getting started

- Download MLP template (scan)
 - https://github.com/mlflow/mlp-regression-template
- Discuss improvement topics
 - https://github.com/mlflow/mlflow/discussions
- Report issues
 - https://github.com/mlflow/mlflow/issues



