

# COMPONENTS IN MLOPS

# Data Management and Processing

**Data Versioning:** Tools like DVC, Pachyderm, or Delta Lake track changes in datasets.

**Data Validation:** Ensures data consistency and integrity before training using tools like Great Expectations or TensorFlow Data Validation (TFDV).

**Feature Engineering & Feature Store:** Organizes reusable features for multiple models using platforms like Feast or Tecton.

# Modular Coding & Code Versioning

**Modular Code Structure:** Breaks down ML code into reusable modules for data preprocessing, training, and inference.

**Code Versioning:** Uses Git and DVC to track code and data versions for reproducibility.

**Automated Testing:** Implements unit tests (e.g., pytest) for data transformations, model predictions, and integrations.

```
main.py  [ ] [ ] [ ] Share Run
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.metrics import accuracy_score
5
6 # Load data
7 df = pd.read_csv("data.csv")
8 X = df.drop(columns=["target"])
9 y = df["target"]
10
11 # Split data
12 X_train, X_test, y_train, y_test = train_test_split(X, y,
13 test_size=0.2, random_state=42)
14
15 # Train model
16 model = RandomForestClassifier(n_estimators=100, random_state=42)
17 model.fit(X_train, y_train)
18
19 # Evaluate
20 y_pred = model.predict(X_test)
21 acc = accuracy_score(y_test, y_pred)
22 print(f"Accuracy: {acc}")
23
```

**Hard to reuse:** If we need another dataset, we must rewrite code.

**No separation of concerns:** Data loading, preprocessing, training, and evaluation are mixed together.

**Difficult debugging:** If something breaks, we have to go through the entire script.

```

from sklearn.ensemble import RandomForestClassifier

def train_model(X_train, y_train, n_estimators=100, random_state=42):
    model = RandomForestClassifier(n_estimators=n_estimators, random_state=
                                   =random_state)
    model.fit(X_train, y_train)
    return model

```

## model\_trainer.py

```

from sklearn.metrics import accuracy_score

def evaluate_model(model, X_test, y_test):
    y_pred = model.predict(X_test)
    return accuracy_score(y_test, y_pred)

```

## evaluate.py

```

1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3
4 def load_data(filepath):
5     df = pd.read_csv(filepath)
6     X = df.drop(columns=["target"])
7     y = df["target"]
8     return X, y
9
10 def split_data(X, y, test_size=0.2, random_state=42):
11     return train_test_split(X, y, test_size=test_size, random_state=
                               =random_state)
12

```

## data\_loader.py

```

from data_loader import load_data, split_data
from model_trainer import train_model
from evaluator import evaluate_model

# Load and preprocess data
X, y = load_data("data.csv")
X_train, X_test, y_train, y_test = split_data(X, y)

# Train model
model = train_model(X_train, y_train)

# Evaluate
accuracy = evaluate_model(model, X_test, y_test)
print(f"Model Accuracy: {accuracy}")

```

## main.py

# Experiment Tracking & Model Development

**Experiment Tracking:** Logs hyperparameters, datasets, and model metrics using MLflow, Weights & Biases, or Neptune.

**Automated Hyperparameter Tuning:** Uses Optuna, Ray Tune, or Hyperopt for optimizing model performance.

**Model Registry:** Stores trained models with metadata, ensuring traceability (MLflow Model Registry, Kubeflow Model Store).

# Model Deployment & Serving

**Model Packaging:** Converts models into deployable artifacts using Docker, TensorFlow Serving, or TorchServe.

## **Model Deployment Strategies:**

**A/B Testing:** Compares two models in production to determine the best one.

**Canary Deployment:** Gradually rolls out a new model to a small portion of traffic.

**Shadow Deployment:** Runs the new model alongside the old one without affecting real users.

# Model Monitoring & Maintenance

**Performance Monitoring:** Tracks metrics like latency, accuracy, and throughput (Prometheus, Grafana, or Aws Sagemaker etc).

**Drift Detection:** Detects changes in input data (data drift) or model performance (concept drift) using tools like WhyLabs or Fiddler.

**Model Retraining & Continuous Learning:** Automates model retraining when performance degrades (Kubeflow Pipelines, Airflow).



# CI/CD for Machine Learning

**Continuous Integration (CI):** Automates testing and validation of ML code and models (Jenkins, GitHub Actions, AWS Tools).

**Continuous Deployment (CD):** Deploys models automatically to production environments (ArgoCD, Azure DevOps).

**Infrastructure as Code (IaC):** Automates provisioning of ML infrastructure using Terraform, AWS CloudFormation.