# Model Monitoring Pipeline and Drift Tracking

## Introduction

Machine learning models are trained to learn patterns from historical data so they can generalize to new, unseen data. However, their performance can degrade over time—a phenomenon known as model drift. Model drift occurs in two main forms: data drift, where the distribution of input features changes from the training data, and concept drift, where the relationship between inputs and outputs evolves. This can lead to poor decision-making and unreliable predictions if not addressed promptly.

To maintain trust, fairness, and accuracy in real-world applications, models must be continuously monitored post-deployment. A robust model monitoring pipeline enables detection of drift, triggers alerts and prompts corrective actions.

## Model Monitoring Pipeline

### Pre-deployment

During model development, validate the model with datasets and check for bias and overfitting. Generate evaluation reports using baseline metrics (e.g., accuracy, F1-score, confusion matrix), and save these for future comparisons post-deployment.

### Deployment phase

The deployment pipeline reuses those baseline reports for automated drift tracking through the following steps:

### 1. Ingestion and Logging

All real-time inference requests (inputs and model outputs) are logged. Logs are stored in a time-stamped format in a centralized data lake or monitoring database. This forms the basis for batch comparisons.

## 2. Input Batch Drift Detection

On a scheduled basis (e.g., nightly or weekly, during 'down times'), run statistical tests to detect changes in input feature distributions using:

- **Kolmogorov-Smirnov test** – to determine whether datasets originate from the same data distribution
- **Chi-squared test** – for categorical features.
- **Population Stability Index (PSI)** – to detect feature distribution shifts.

Automated scripts calculate these metrics. If drift exceeds predefined thresholds, alerts are triggered.

## 3. Prediction Monitoring

Track output drift by examining prediction probabilities. Use measures such as:

- **Entropy** – to measure prediction uncertainty.
- **Kullback-Leibler (KL) Divergence** – to compare predicted distributions over time.
- **Prediction Confidence/Spread** – to detect signs of overfitting or uncertainty.

Threshold violations can trigger notifications, prompting manual inspection or retraining.

## 4. Compare Accuracy with Training (Label Drift)

Once ground truth becomes available (e.g., from user feedback or delayed labels), compare predictions to actual outcomes. Schedule accuracy checks using:

- **Accuracy, Precision, Recall** – for classification tasks.
- **RMSE, MAE** – for regression tasks.
- **BLEU, ROUGE, or Word Error Rate (WER)** – for NLP or transcription tasks.

Calculate these metrics on a rolling basis and compare them with training-time benchmarks. If performance drops by a significant margin, set up alerts for retraining or data rebalancing.

**Conclusion**

Model monitoring and drift detection are crucial to ensuring reliable AI systems. A pipeline with automated statistical checks, prediction monitoring, and ground truth evaluation enables teams to catch and respond to drift early. With proper thresholds and alerting in place, organizations can reduce risk, improve fairness, and ensure sustained model quality in production.