



Optimizing MLOps for Enhanced Accident Prediction Models

Presented by: Abhishikth Peri, Ameya Rahurkar, Kexin Zhang, Ryan Hsieh, Venkata Sai Krishna Abbaraju

Overview



Current Situation

Understanding existing
Model Development
processes



Our Solution

Improve model
development using best
MLOps practices



Impact Created

Discussing benefits in
development and
management

Understanding the landscape



01

Improve model development using best MLOps practices

02

Manage ML models efficiently by monitoring them using open source tools

03

Make a calculated approach to retrain the model and save man-hours.

Introduction

- Streamlining ML Development
 - Adopting MLOps best practices.
 - Ensuring scalability, efficiency, and reproducibility.
- Enhanced Model Management
 - Utilizing advanced open-source tools.
 - Continuous monitoring for improved performance and reliability.
- Strategic Model Retraining
 - Data-driven methodology for retraining decisions.
 - Employing drift detection for timely insights.

Streamlined Workflow



History Data Store
Amazon S3



Data Streaming
Apache Kafka



Data Processing
Apache Spark



Data Storage
Hive warehouse

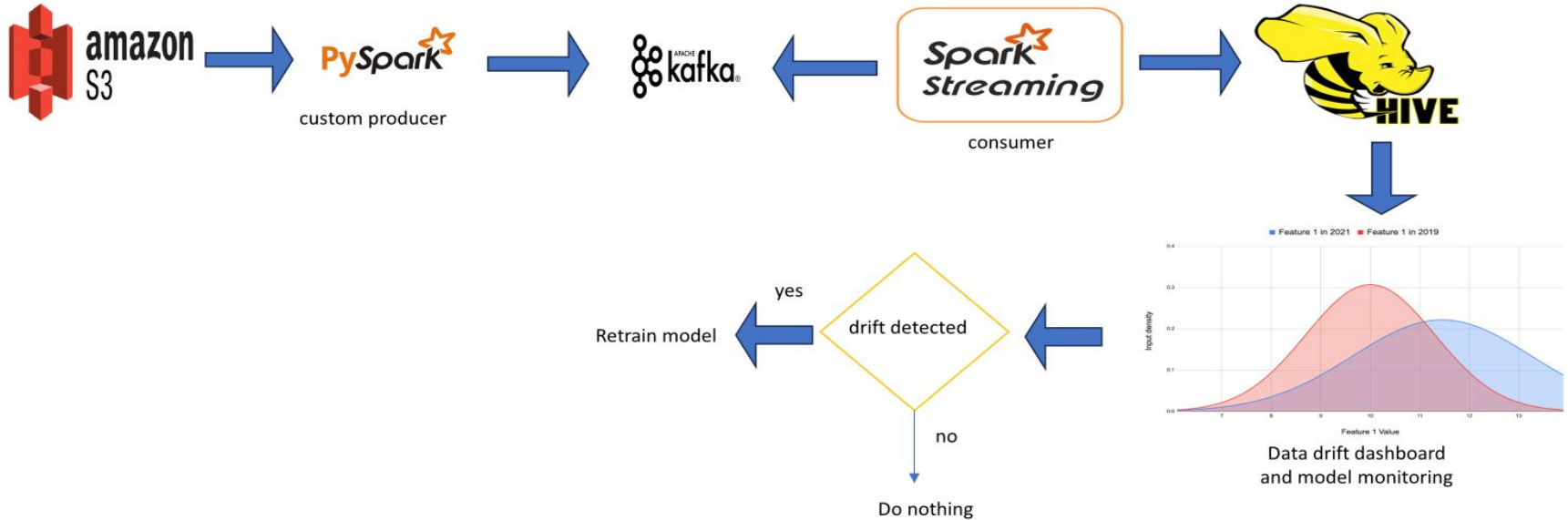


Severity Prediction



**Data Drift
Detection**

Architecture



Feature Engineering

- 1) *Converted data type to better fit the model and apply the pipeline*
 - 2) *Identified and addressed data quality issues – dropped rows with nulls in columns like Location information and weather information such as Wind Chill, Temperature*
 - 3) *Detected and excluded unrealistic outlier values in Precipitation, and Distance, e.g., 80% of Distance data erroneously recorded as 0.*
 - 4) *Developed new features like 'Accident Duration Time' to deepen analysis. However, 'Comfort Index' deemed less influential in predicting accidents*
 - 5) *Random Forest identified Humidity, Pressure, Temperature, and Wind Speed as significant numerical predictors.*
 - 6) *Analysis revealed a higher incidence of severe accidents near crossings, highlighting the importance of location-based binary features.*
-

MLOps PIPELINE

Design



Business
Requirements

Development



Data Engineering
Model Building

Deployment



Model Deployment

Operations



System Monitoring



MLOps: Streamlining the Machine Learning Lifecycle

Databricks for ML Experiments:

- *"Leveraged Databricks as our primary environment for running scalable PySpark ML experiments."*

Hyperparameter Tuning with Cross-Validation:

- *"Utilized cross-validation techniques to fine-tune hyperparameters, ensuring optimal model performance."*

Experiment Tracking with MLflow:

- *"Employed MLflow for comprehensive tracking of experiments, including hyperparameters, metrics, and model comparisons."*

Model Comparison and Visualization:

- *"Analyzed and visualized model performance and hyperparameter impacts using MLflow's comparison charts."*

Model Registration and Deployment:

- *"Successfully registered the final model in MLflow and explored deployment options for real-time predictions."*



Proactive Data Drift Detection for Model Reliability"

Implementing Drift Detection with Evidently AI:

- *"Integrated Evidently AI to systematically monitor and detect data drift in incoming data streams."*

Comparative Data Analysis:

- *"Conducted weekly data quality checks by comparing current data against a 2-week historical sample."*

Automated Reporting in MLflow:

- *"Automatically generated and logged detailed drift reports in MLflow for traceability and review."*

Alerting Mechanism:

- *"Set up alerts to promptly notify stakeholders about detected drifts and potential data quality issues."*

Decision Making on Model Retraining:

- *"Deterministic process for model retraining based on drift detection outcomes."*

What we did?



Enhanced crisis management

- Real-time predictions of traffic accident severity
- Optimize resource allocation and response times



Effortless model maintenance

- Drift detection mechanism to streamline model re-training
- Alerting mechanism to notify variations in batch vs stream

How it's better?



Efficient Model Development

- Long term sustainability of prediction accuracy
- Save man-hours to retrain and recompute



Enhanced Risk Assessment

- Simplify insurance claims and improve risk modeling
- Strengthen real-time risk assessment capabilities