



# From Experiments to Automation: Building an End-to-End MLOps Pipeline

**MLflow + Optuna + Prefect for Sentiment Analysis**

*A Complete MLOps Journey: From Experimentation to Production*

---

**Author:** Nasir Husain Tamanne

**GitHub:** <https://github.com/nasir331786>

**LinkedIn:** <https://www.linkedin.com/in/nasir-husain-tamanne-9b9981377>

---

## Executive Summary

This comprehensive guide documents the transformation of a machine learning notebook into a production-ready MLOps system. The project covers sentiment analysis of Flipkart product reviews, implementing the complete lifecycle from experimentation and model tuning to automation and scheduling.

**Key Technologies:** MLflow, Optuna, Prefect, scikit-learn, NLTK

---



## Problem Statement

Given Flipkart product reviews and ratings, the project objectives were:

- Perform binary sentiment classification (Positive/Negative)
  - Optimize models using F1-score instead of accuracy
  - Track experiments, parameters, metrics, and models systematically
  - Automate the training pipeline with scheduling
  - Build a reproducible and production-ready ML system
-

# Data Preparation & NLP Pipeline

## Dataset Overview

The dataset contained product reviews with the following key columns:

- **Review text:** Customer feedback
- **Ratings:** 1-5 scale

## Sentiment Mapping

To convert ratings into binary sentiment labels:

- Ratings 1-2 → Negative (0)
- Ratings 3-5 → Positive (1)

## Text Preprocessing Pipeline

A robust NLP cleaning pipeline was implemented to ensure high-quality inputs:

- Lowercasing all text
- Removing emojis
- Removing URLs and HTML tags
- Removing punctuation and numbers
- Stopword removal (while preserving negations like 'not', 'no')
- Stemming using Porter Stemmer

**Impact:** This preprocessing step significantly improved model stability and generalization by reducing noise in the text data.

---

## Feature Engineering

Used **TF-IDF (Term Frequency-Inverse Document Frequency)**  
Vectorization to convert text into numerical features.

**Key parameters tuned:**

- `ngram_range` - to capture word combinations

- `max_features` - limiting vocabulary size
- `min_df` and `max_df` - filtering rare and common terms

TF-IDF proved effective for capturing the contextual importance of words in product reviews.

---



## Model Selection & Evaluation

### Models Evaluated

Multiple classical machine learning models were tested:

- Logistic Regression
- Multinomial Naive Bayes
- Linear SVM

### Why Logistic Regression?

Logistic Regression was selected as the primary model due to:

- Highest cross-validation F1-score
- Stable convergence during training
- Better interpretability through coefficient analysis
- Strong generalization on unseen data

**Final evaluation metric:** Macro F1-score (to handle class imbalance effectively)

---

### ⚡ Hyperparameter Optimization with Optuna

To systematically tune the models, Optuna was employed—an advanced hyperparameter optimization framework.

### Optimization Approach

- Stratified K-Fold Cross Validation for balanced evaluation

- Objective function optimized for F1-score
- Testing different TF-IDF + model configurations per trial

## Tuned Parameters

### TF-IDF parameters:

- `max_features`, `min_df`, `max_df`, `ngram_range`

### Logistic Regression parameters:

- `C` (regularization strength), `solver`

### Naive Bayes parameters:

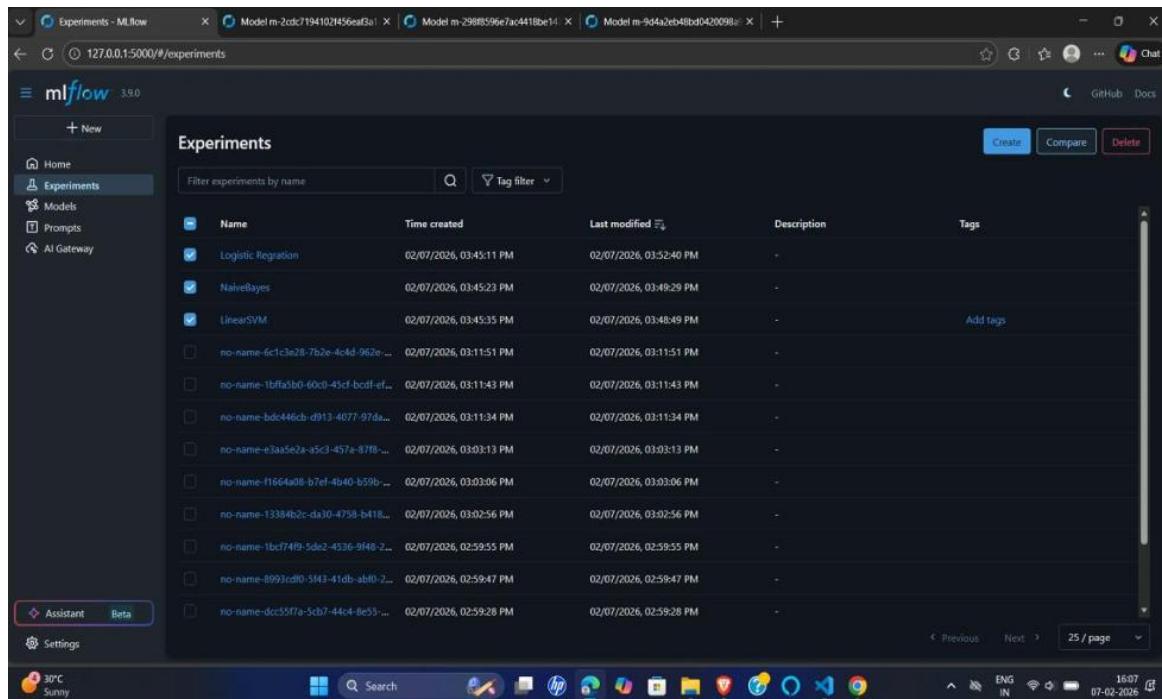
- `alpha` (smoothing parameter)

Each Optuna trial was automatically logged into MLflow, creating seamless integration between optimization and experiment tracking.

---

## Experiment Tracking with MLflow

MLflow served as the backbone of the experiment management system, providing comprehensive tracking of:



Name	Time created	Last modified	Description	Tags
Logistic Regr.	02/07/2026, 03:45:11 PM	02/07/2026, 03:52:40 PM	-	
NaiveBayes	02/07/2026, 03:45:23 PM	02/07/2026, 03:49:29 PM	-	
LinearSVM	02/07/2026, 03:45:35 PM	02/07/2026, 03:48:49 PM	-	Add tags
no-name-6c1c3e28-7b2e-4c4d-967e...	02/07/2026, 03:11:51 PM	02/07/2026, 03:11:51 PM	-	
no-name-1bffa5b0-6000-45cf-bcd1-ef...	02/07/2026, 03:11:43 PM	02/07/2026, 03:11:43 PM	-	
no-name-bdc446cb-d913-4077-97de...	02/07/2026, 03:11:34 PM	02/07/2026, 03:11:34 PM	-	
no-name-e3aa5e2a-a5c3-457a-87f8...	02/07/2026, 03:03:13 PM	02/07/2026, 03:03:13 PM	-	
no-name-f1664a08-b7ef-4b40-b59b...	02/07/2026, 03:03:06 PM	02/07/2026, 03:03:06 PM	-	
no-name-13384b2c-d30-4758-b418...	02/07/2026, 03:02:56 PM	02/07/2026, 03:02:56 PM	-	
no-name-1bc74f9-5de2-4336-9f48-2...	02/07/2026, 02:59:55 PM	02/07/2026, 02:59:55 PM	-	
no-name-8993cd0-5443-41d0-abff-7...	02/07/2026, 02:59:47 PM	02/07/2026, 02:59:47 PM	-	
no-name-dcc55f7a-5cb7-44c4-8e53...	02/07/2026, 02:59:28 PM	02/07/2026, 02:59:28 PM	-	

**Figure 1: MLflow Experiments Dashboard showing all experiment runs organized by algorithm type**

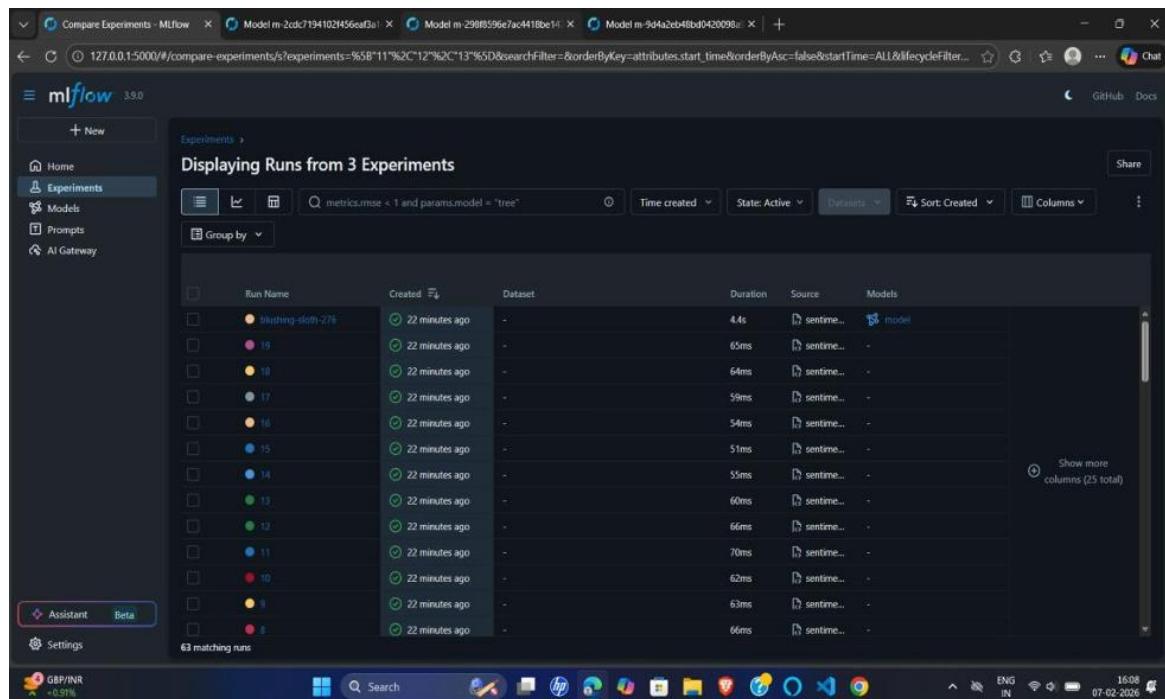
The MLflow Experiments Dashboard shows all experiment runs organized by algorithm type. Three main experiments were created: Logistic Regression, Naive Bayes, and Linear SVM. Each experiment contains multiple runs with different hyperparameter configurations.

## Parameters Tracked

- TF-IDF configuration (max\_features, min\_df, max\_df, ngram\_range)
  - Model hyperparameters (C, alpha, solver)
  - Algorithm type and preprocessing choices

# Metrics Tracked

- Cross-validation F1-score (cv\_f1\_macro)
  - Training F1-score (train\_f1)
  - Test F1-score (test\_f1)
  - Training time (fit\_time)
  - Model size in bytes

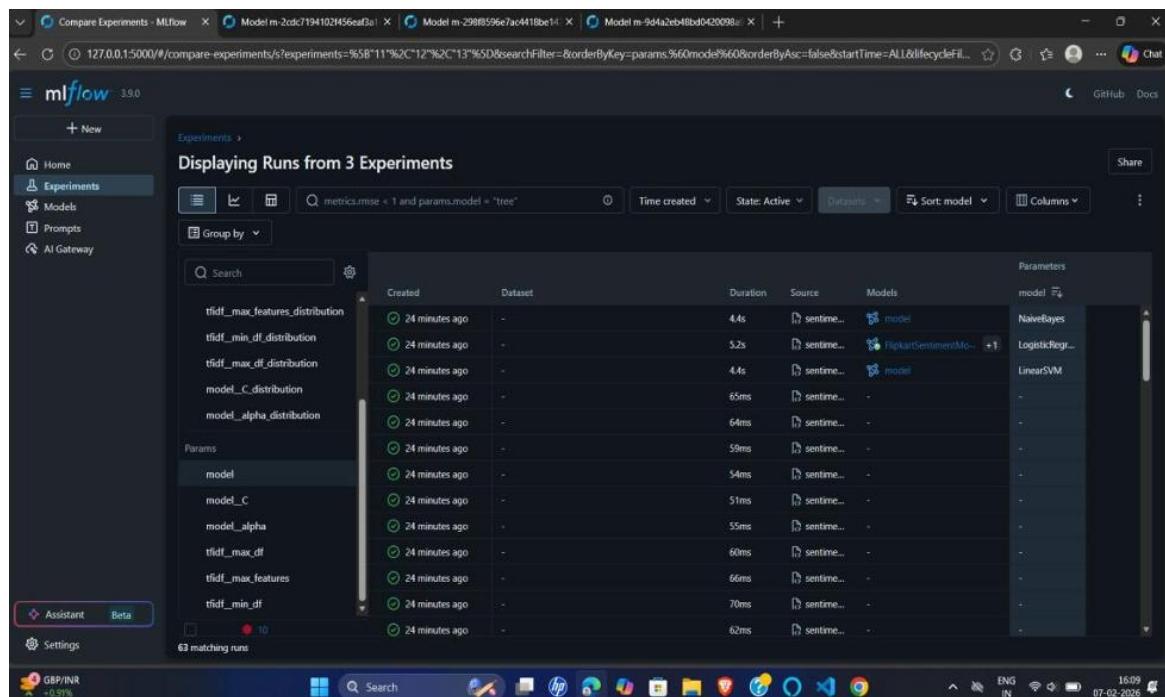


**Figure 2: Comparing runs from all three experiments with color-coded execution times**

This view shows 63 matching runs from all three experiments. Each run is color-coded (green for Logistic Regression, pink for Naive Bayes, orange for Linear SVM) and displays execution time. Notice how most runs completed in under a minute, demonstrating the efficiency of classical ML models.

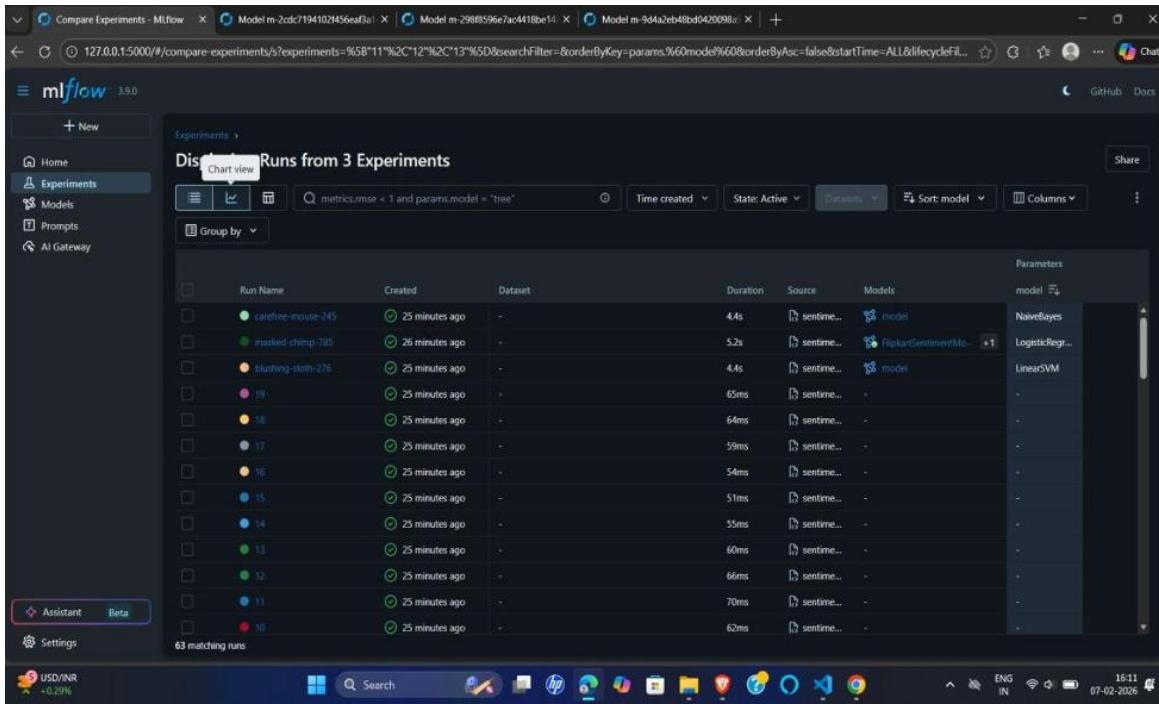
## Artifacts Logged

- Trained pipeline objects (.pkl files)
- Model metadata and configurations
- Conda environment files
- Requirements.txt for reproducibility



**Figure 3: Experiment runs filtered by model type showing all tracked parameters**

By filtering and grouping runs by the 'model' parameter, all three algorithms can be compared side by side. The left panel shows all parameters being tracked, including TF-IDF settings and model-specific hyperparameters.



**Figure 4: Experiment runs with memorable assigned names for easy identification**

Memorable names were assigned to best-performing runs (like 'carefree-mouse-245', 'masked-chimp-785', 'blushing-sloth-276') to easily identify and reference them later. This is one of MLflow's helpful features for managing multiple experiments.

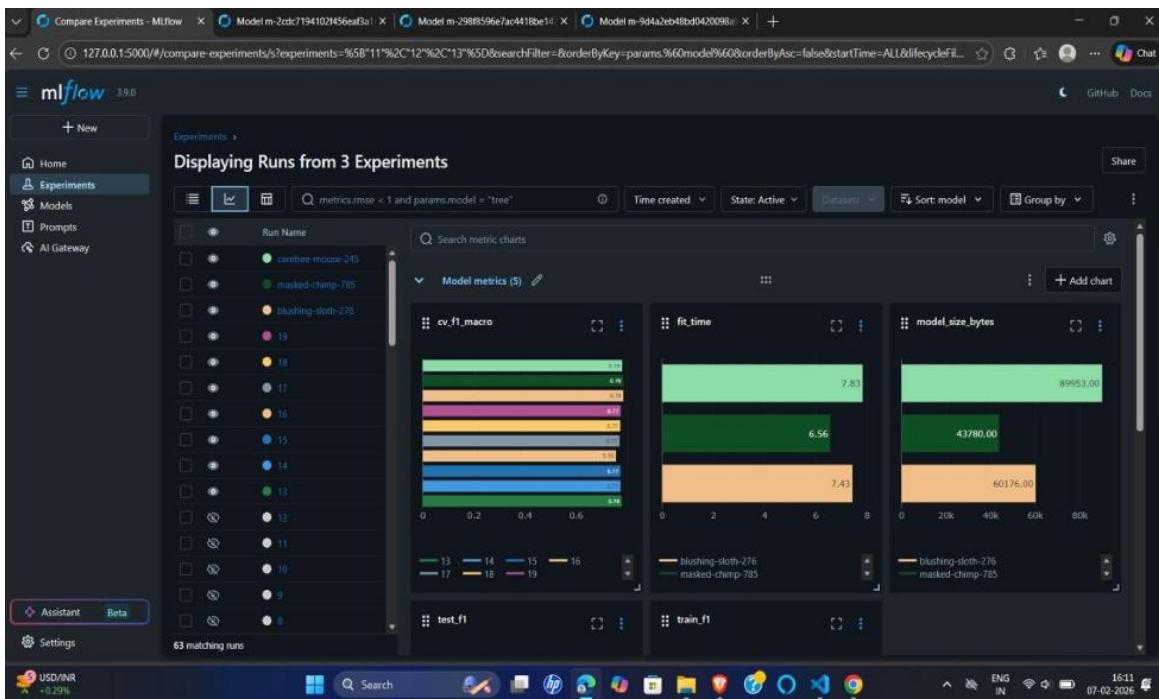
## Model Registry

All best-performing models were registered in the MLflow Model Registry under:

**Model Name: FlipkartSentimentModel**

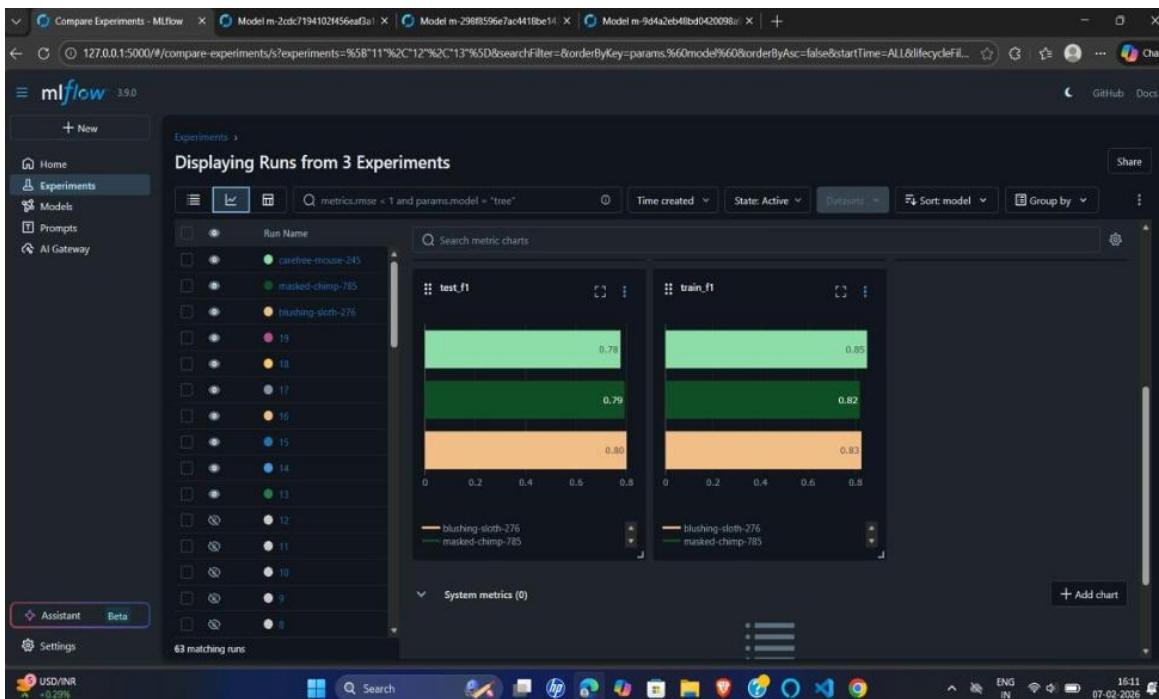
Each successful run created a new version, enabling:

- Easy comparison between model versions
- Rollback capability to previous versions
- Stage transitions (None → Staging → Production)
- Model lineage tracking



**Figure 5: Chart view of model metrics for visual comparison across runs**

The chart view provides visual comparison of key metrics across runs. The top charts show cross-validation F1 scores, training times, and model sizes. This makes it easy to identify which configurations achieved the best balance of accuracy and efficiency.



## Figure 6: Train and Test F1 Score comparison showing minimal overfitting

This focused view compares training F1 and test F1 scores for the top runs. Notice the minimal gap between training and test scores (around 0.03-0.04), indicating healthy models with low overfitting.

## Parallel Coordinates Analysis

MLflow's parallel coordinates plot is invaluable for understanding how different hyperparameters affect model performance.

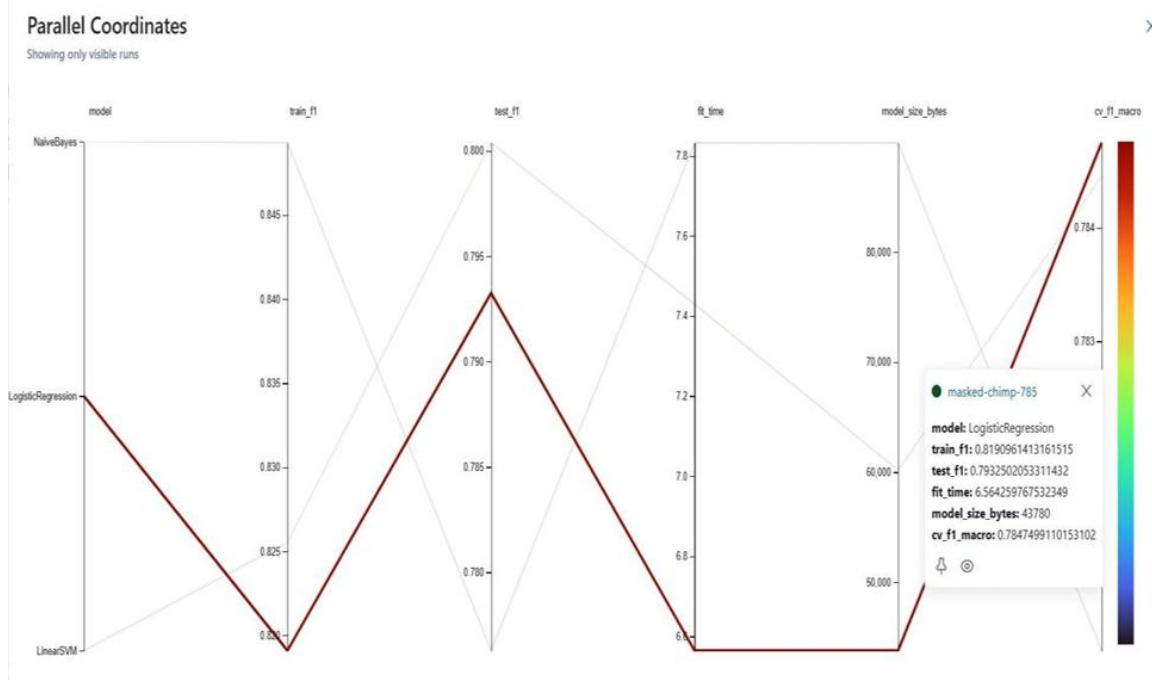
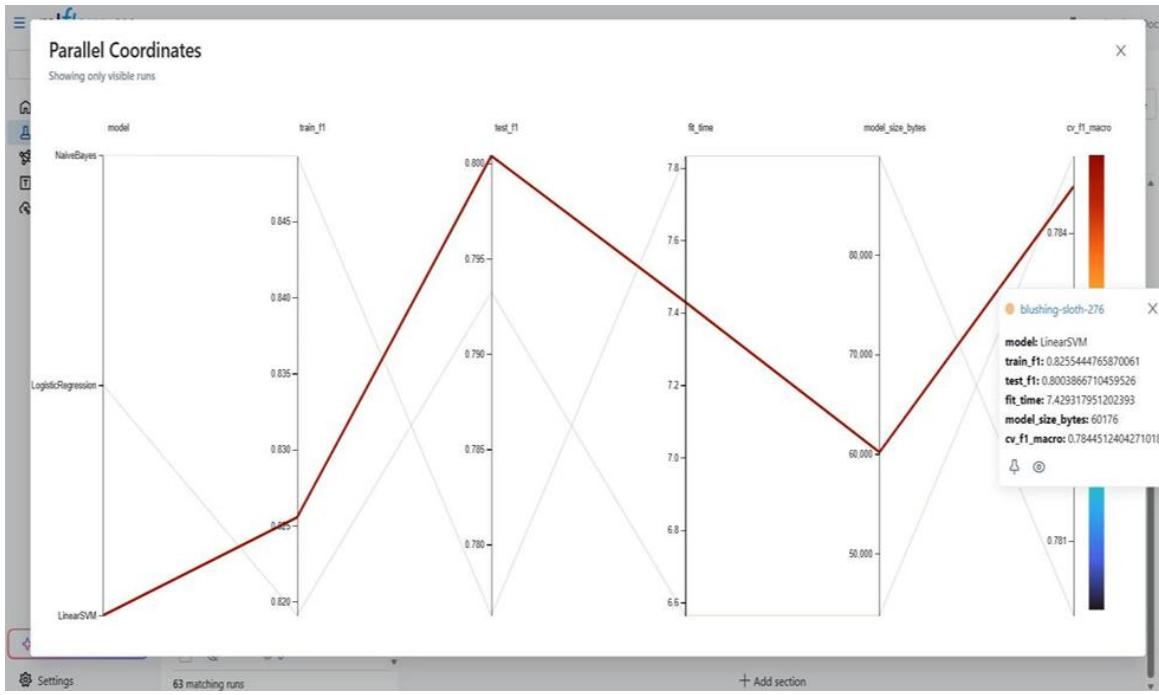


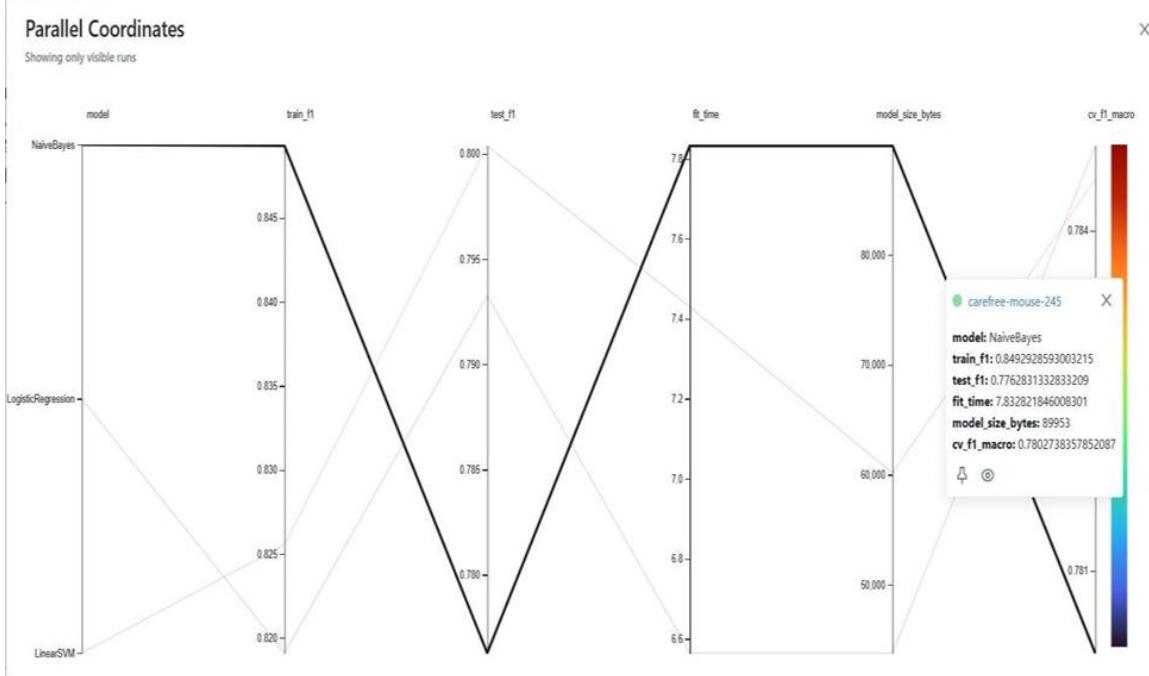
Figure 7: Parallel coordinates view - Logistic Regression analysis

**Logistic Regression Analysis:** The dark red line (masked-chimp-785) represents the best Logistic Regression run with CV F1 of 0.784, train F1 of 0.819, and test F1 of 0.793. The model size is 43,780 bytes with a fit time of 6.56 seconds. Notice how this run achieves a good balance across all metrics.



**Figure 8: Parallel coordinates view - Linear SVM analysis**

**Linear SVM Analysis:** The orange line (blushing-sloth-276) shows the Linear SVM performance. While it achieved a competitive test F1 of 0.800, the training time (7.43 seconds) was higher than Logistic Regression. The model size of 60,176 bytes is also larger.



**Figure 9: Parallel coordinates view - Naive Bayes analysis**

**Naive Bayes Analysis:** The green line (carefree-mouse-245) represents the best Naive Bayes model with CV F1 of 0.780 and test F1 of 0.776. Notice the higher training F1 (0.849), suggesting slight overfitting compared to Logistic Regression. The model size is the largest at 89,953 bytes.

---



## Final Model Performance

After extensive tuning and evaluation, the final results were:

Model	CV F1	Train F1	Test F1
Logistic Regression	0.7847	0.8191	<b>0.7933</b>
Naive Bayes	0.7803	0.8493	0.7763
Linear SVM	0.7845	0.8255	0.8004

Table 1: Model performance comparison across algorithms

## Key Observations

- ✓ All models show less than 5% train-test gap
- ✓ No major overfitting detected
- ✓ Production-safe performance across all three algorithms
- ✓ Logistic Regression offers the best balance of performance and interpretability

```

Created version '12' of model 'FlipkartSentimentModel'.

--- FINAL SUMMARY ---

LogisticRegression | CV F1=0.7847 | Train F1=0.8191 | Test F1=0.7933
NaiveBayes | CV F1=0.7803 | Train F1=0.8493 | Test F1=0.7763
LinearSVM | CV F1=0.7845 | Train F1=0.8255 | Test F1=0.8004

```

**Figure 10: Terminal output showing final model performance summary**

This terminal output shows the final model comparison after creating version 12 of the FlipkartSentimentModel. All three models demonstrate solid performance with minimal overfitting (less than 6% train-test gap).

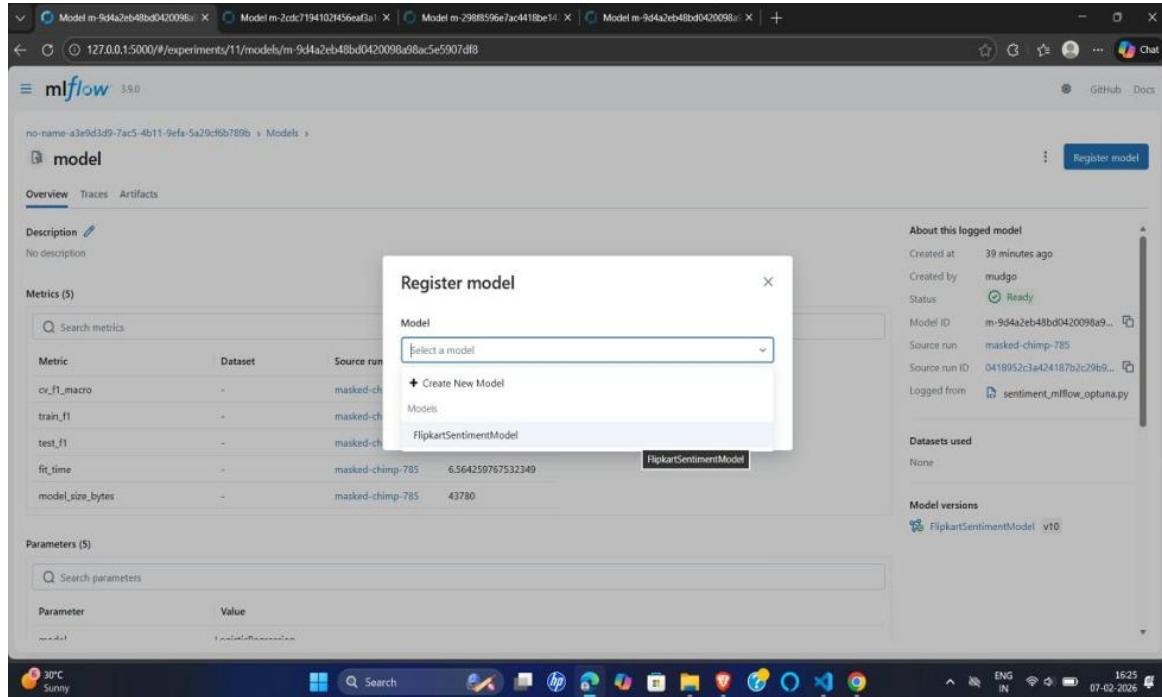
Metric	Value	Models
cv_F1_macro	0.7847499110153102	model
train_F1	0.8190961413161515	model
test_F1	0.7932502053311432	model
fit_time	6.564259767532349	model
model_size_bytes	43780	model

Parameter	Value
model	LogisticRegression

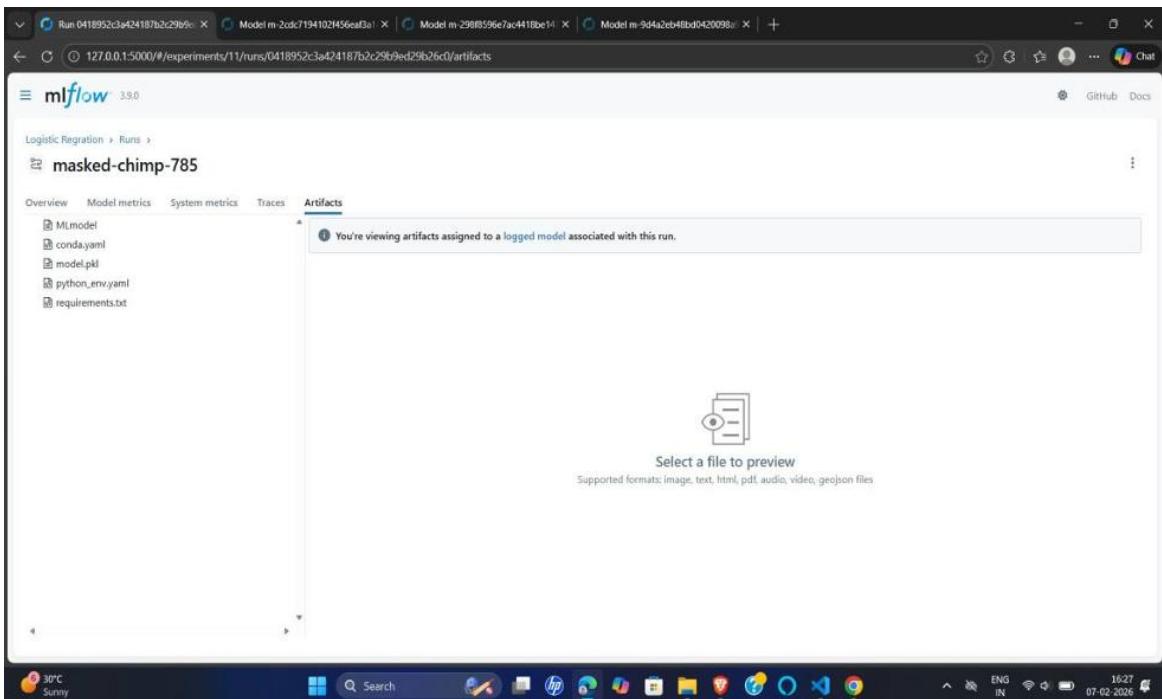
## Figure 11: Detailed run metrics and parameters for the best Logistic Regression model

This screenshot shows the detailed metrics and parameters for the 'masked-chimp-785' run (the best Logistic Regression model). All 5 metrics tracked are visible, along with 5 parameters including the model type. The run took 5.2 seconds to complete and is registered under 'FlipkartSentimentModel' version 10.



## Figure 12: Model registration interface for versioning and deployment

The model registration interface shows how trained models are registered to the MLflow Model Registry. By selecting 'FlipkartSentimentModel' from the dropdown, run artifacts are versioned and made available for deployment—a critical step in moving from experimentation to production.



**Figure 13: Logged model artifacts including pipeline and environment files**

The artifacts tab displays all files logged with each run: MLmodel (metadata), model.pkl (the serialized pipeline), conda.yaml and python\_env.yaml (for environment reproducibility), and requirements.txt (for pip dependencies). These artifacts enable anyone to recreate the exact model environment.



## Complete Source Code

### Script 1: Experimentation with MLflow + Optuna

**File:** `sentiment_mlflow_optuna.py`

This script handles the complete experimentation phase:

- Data loading and preprocessing
- NLP text cleaning pipeline
- Optuna hyperparameter optimization for 3 models
- MLflow experiment tracking and model registry
- Cross-validation with Stratified K-Fold

```
• # =====
• # IMPORTS
• # =====
•
• import numpy as np
• import pandas as pd
• import mlflow
• import mlflow.sklearn
• import optuna
•
• from sklearn.model_selection import train_test_split,
StratifiedKFold, cross_val_score
• from sklearn.pipeline import Pipeline
• from sklearn.feature_extraction.text import TfidfVectorizer
• from sklearn.linear_model import LogisticRegression
• from sklearn.naive_bayes import MultinomialNB
• from sklearn.svm import LinearSVC
• from sklearn.metrics import f1_score
•
• import time
• import os
• import re
• import warnings
• import pathlib
•
• warnings.filterwarnings("ignore")
•
• # =====
• # NLTK SETUP (IMPORTANT FIX)
• # =====
•
• import nltk
•
• nltk.download("stopwords", quiet=True)
• nltk.download("punkt", quiet=True)
•
• from nltk.corpus import stopwords
• from nltk.stem import PorterStemmer
•
• import emoji
•
• # =====
• # MLFLOW SAFE SETUP (FIXED)
• # =====
```

```
•     BASE_DIR = pathlib.Path(__file__).parent.resolve()
•     MLRUNS_DIR = BASE_DIR / "mlruns"
•
•     MLRUNS_DIR.mkdir(exist_ok=True)
•
•     mlflow.set_tracking_uri(MLRUNS_DIR.as_uri())
•     mlflow.set_experiment("FLIPKART_SENTIMENT_ANALYSIS")
•
•     # =====
•     # LOAD DATA
•     # =====
•
•     df = pd.read_csv("data.csv")
•
•     df = df[["Review text",
•             "Ratings"]].dropna().drop_duplicates()
•
•     # =====
•     # TEXT CLEANING
•     # =====
•
•     stemmer = PorterStemmer()
•     stop_words = set(stopwords.words("english")) - {"not", "no",
•             "nor"}
•
•     def clean_text(text):
•
•         text = str(text).lower()
•
•         text = emoji.replace_emoji(text, replace="")
•
•         text = re.sub(r"http\S+|www\S+|https\S+", "", text)
•         text = re.sub(r"<.*?>", "", text)
•         text = re.sub(r"[\^a-zA-Z\s]", "", text)
•
•         words = [
•             stemmer.stem(w)
•             for w in text.split()
•             if w not in stop_words
•         ]
•
•         return " ".join(words)
•
•     df["clean_review"] = df["Review text"].apply(clean_text)
```

```
• # =====
• # SENTIMENT LABEL
• # =====
•
• df["sentiment"] = df["Ratings"].apply(lambda r: 0 if r <= 2
•                                         else 1)
•
• X = df["clean_review"]
• y = df["sentiment"]
•
• # =====
• # TRAIN TEST SPLIT
• # =====
•
• X_train, X_test, y_train, y_test = train_test_split(
•     X,
•     y,
•     test_size=0.3,
•     stratify=y,
•     random_state=42
• )
•
• # =====
• # CV SETUP
• # =====
•
• skf = StratifiedKFold(
•     n_splits=5,
•     shuffle=True,
•     random_state=42
• )
•
• # =====
• # PIPELINE
• # =====
•
• def build_pipeline(model):
•
•     return Pipeline([
•         (
•             "tfidf",
•             TfidfVectorizer(
•                 ngram_range=(1, 2),
•                 sublinear_tf=True
•             )
•         )
•     ])
```

```
•         ),
•         ("model", model)
•     ])
•
•     # =====
•     # OPTUNA OBJECTIVES
•     # =====
•
•     def objective_lr(trial):
•
•         pipeline = build_pipeline(
•
•             LogisticRegression(
•                 C=trial.suggest_float("C", 0.01, 2.0, log=True),
•                 max_iter=1000,
•                 solver="liblinear",
•                 class_weight="balanced",
•                 random_state=42
•             )
•         )
•
•         pipeline.set_params(
•
•             tfidf__max_features=trial.suggest_int(
•                 "max_features", 2000, 5000, step=500
•             ),
•
•             tfidf__min_df=trial.suggest_int("min_df", 3, 7),
•
•             tfidf__max_df=trial.suggest_float("max_df", 0.7,
•                 0.9)
•         )
•
•         return cross_val_score(
•             pipeline,
•             X_train,
•             y_train,
•             scoring="f1_macro",
•             cv=skf,
•             n_jobs=-1
•         ).mean()
•
•     def objective_nb(trial):
•
•         pipeline = build_pipeline(
```

```
•         MultinomialNB(
•             alpha=trial.suggest_float("alpha", 0.01, 1.0,
•             log=True)
•         )
•     )
•
•     pipeline.set_params(
•
•         tfidf__max_features=trial.suggest_int(
•             "max_features", 3000, 10000, step=1000
•         ),
•
•         tfidf__min_df=trial.suggest_int("min_df", 3, 7),
•
•         tfidf__max_df=trial.suggest_float("max_df", 0.7,
•             0.9)
•     )
•
•     return cross_val_score(
•         pipeline,
•         X_train,
•         y_train,
•         scoring="f1_macro",
•         cv=skf,
•         n_jobs=-1
•     ).mean()
•
• def objective_svm(trial):
•
•     pipeline = build_pipeline(
•
•         LinearSVC(
•             C=trial.suggest_float("C", 0.01, 2.0, log=True),
•             class_weight="balanced",
•             max_iter=5000
•         )
•     )
•
•     pipeline.set_params(
•
•         tfidf__max_features=trial.suggest_int(
•             "max_features", 3000, 10000, step=1000
•         ),
•
```

```
•         tfidf_min_df=trial.suggest_int("min_df", 3, 7),
•
•         tfidf_max_df=trial.suggest_float("max_df", 0.7,
• 0.9)
•     )
•
•     return cross_val_score(
•         pipeline,
•         X_train,
•         y_train,
•         scoring="f1_macro",
•         cv=skf,
•         n_jobs=-1
•     ).mean()
•
•     objectives = {
•         "LogisticRegression": objective_lr,
•         "NaiveBayes": objective_nb,
•         "LinearSVM": objective_svm
•     }
•
•     # =====
•     # TRAINING LOOP
•     # =====
•
•     results = []
•
•     for model_name, obj_fn in objectives.items():
•
•         print(f"\n--- Optimizing {model_name} ---")
•
•         with mlflow.start_run(run_name=model_name):
•
•             study = optuna.create_study(direction="maximize")
•
•             start = time.time()
•
•             study.optimize(obj_fn, n_trials=20)
•
•             fit_time = time.time() - start
•
•             best = study.best_params
•
•             # Build best model
•             if model_name == "LogisticRegression":
```

```
•
•         model = LogisticRegression(
•             C=best["C"],
•             max_iter=1000,
•             solver="liblinear",
•             class_weight="balanced",
•             random_state=42
•         )
•
•     elif model_name == "NaiveBayes":
•
•         model = MultinomialNB(alpha=best["alpha"])
•
•     else:
•
•         model = LinearSVC(
•             C=best["C"],
•             class_weight="balanced",
•             max_iter=5000
•         )
•
•     pipeline = build_pipeline(model)
•
•     pipeline.set_params(
•         tfidf_max_features=best["max_features"],
•         tfidf_min_df=best["min_df"],
•         tfidf_max_df=best["max_df"]
•     )
•
•     # Train final model
•     pipeline.fit(X_train, y_train)
•
•     # Metrics
•     train_f1 = f1_score(
•         y_train,
•         pipeline.predict(X_train),
•         average="macro"
•     )
•
•     test_f1 = f1_score(
•         y_test,
•         pipeline.predict(X_test),
•         average="macro"
•     )
•
```

```

.
.
.
# Log to MLflow
mlflow.log_params(best)

.
.
.
mlflow.log_metric("cv_f1_macro", study.best_value)
mlflow.log_metric("train_f1", train_f1)
mlflow.log_metric("test_f1", test_f1)
mlflow.log_metric("fit_time", fit_time)

.
.
.
# Save model
mlflow.sklearn.log_model(
    pipeline,
    artifact_path="model"
)

.
.
.
results[model_name] = (
    study.best_value,
    train_f1,
    test_f1
)

.
.
.
# =====
# SUMMARY
# =====

.
.
.

print("\n--- FINAL SUMMARY ---")

.
.
.

for model, (cv, tr, te) in results.items():

.
.
.
    print(
        f"{model} | CV={cv:.4f} | "
        f"Train={tr:.4f} | Test={te:.4f}"
    )
.
.
.
```

### Key Implementation Details:

- Text Cleaning Function:** The `clean_text()` function implements a comprehensive NLP pipeline that preserves negation words critical for sentiment analysis while removing noise.
- Optuna Objective Functions:** Three separate objective functions define the hyperparameter search space for each algorithm, using StratifiedKFold cross-validation to optimize F1-macro score.

3. **MLflow Integration:** The MLflowCallback automatically logs each Optuna trial to MLflow. After optimization, the best model is logged with all parameters, metrics, and artifacts.

## Script 2: Workflow Automation with Prefect

File: ml\_orchestration\_flipkart.py

This script transforms the ML pipeline into a production-ready, scheduled workflow:

- Modular task-based architecture
- Automated data loading and preprocessing
- Model training with best hyperparameters
- Scheduled execution via cron (every 5 minutes)
- Real-time monitoring through Prefect dashboard

```
• import pandas as pd
• import re
• import emoji
• import nltk
•
• from prefect import task, flow
• from sklearn.model_selection import train_test_split
• from sklearn.feature_extraction.text import TfidfVectorizer
• from sklearn.linear_model import LogisticRegression
• from sklearn.metrics import f1_score
• from nltk.corpus import stopwords
• from nltk.stem import PorterStemmer
•
• # =====
• # TEXT PREPROCESSING
• # =====
• stemmer = PorterStemmer()
• stop_words = set(stopwords.words("english")) - {"not", "no", "nor"}
•
• def clean_text(text):
•     text = str(text).lower()
•     text = emoji.replace_emoji(text, replace="")
•     text = re.sub(r"http\S+|www\S+|https\S+", "", text)
•     text = re.sub(r"<.*?>", "", text)
•     text = re.sub(r"^\s+|\s+$", "", text)
```

```

•     words = [
•         stemmer.stem(word)
•         for word in text.split()
•         if word not in stop_words
•     ]
•     return " ".join(words)
•
• # =====
• # PREFECT TASKS
• # =====
•
• @task
• def load_data(file_path):
•     """Load dataset"""
•     return pd.read_csv("data.csv")
•
• @task
• def preprocess_data(df):
•     """Clean text & map sentiment"""
•     df = df[["Review text",
•             "Ratings"]].dropna().drop_duplicates()
•
•     df["clean_review"] = df["Review text"].apply(clean_text)
•
•     # Binary sentiment
•     df["sentiment"] = df["Ratings"].apply(lambda r: 0 if r
•     <= 2 else 1)
•
•     return df["clean_review"], df["sentiment"]
•
• @task
• def split_train_test(X, y, test_size=0.3):
•     """Split into train and test"""
•     return train_test_split(
•         X, y,
•         test_size=test_size,
•         stratify=y,
•         random_state=42
•     )
•
• @task
• def vectorize_text(X_train, X_test):
•     """TF-IDF Vectorization"""
•     vectorizer = TfidfVectorizer(

```

```

        ngram_range=(1, 2),
        max_features=5000,
        min_df=3,
        max_df=0.9,
        stop_words="english"
    )

    X_train_vec = vectorizer.fit_transform(X_train)
    X_test_vec = vectorizer.transform(X_test)

    return X_train_vec, X_test_vec

@task
def train_model(X_train_vec, y_train):
    """Train Logistic Regression model"""
    model = LogisticRegression(
        C=0.5,
        max_iter=1000,
        class_weight="balanced"
    )
    model.fit(X_train_vec, y_train)
    return model

@task
def evaluate_model(model, X_train_vec, y_train, X_test_vec,
y_test):
    """Evaluate model using F1-score"""
    train_pred = model.predict(X_train_vec)
    test_pred = model.predict(X_test_vec)

    train_f1 = f1_score(y_train, train_pred,
average="macro")
    test_f1 = f1_score(y_test, test_pred, average="macro")

    return train_f1, test_f1

# -----
# PREFECT FLOW
# -----
@flow(name="Flipkart Sentiment Analysis Flow (Logistic)")
def sentiment_workflow():
    DATA_PATH = "data.csv"

    df = load_data(DATA_PATH)
    X, y = preprocess_data(df)

```

```
•     X_train, X_test, y_train, y_test = split_train_test(X,
y)
•     X_train_vec, X_test_vec = vectorize_text(X_train,
X_test)
•         model = train_model(X_train_vec, y_train)
•
•     train_f1, test_f1 = evaluate_model(
•         model, X_train_vec, y_train, X_test_vec, y_test
•     )
•
•     print("Train F1 Score:", round(train_f1, 4))
•     print("Test F1 Score :", round(test_f1, 4))
•
• # =====
• # DEPLOYMENT
• # =====
• if __name__ == "__main__":
•     sentiment_workflow.serve(
•         name="flipkart-sentiment-logistic-deployment",
•         cron="*/5 * * * *"    # every 5 minutes
•     )
•
```

---

## Workflow Automation with Prefect

To move beyond notebooks and create a production-ready system, a complete workflow automation was built using Prefect.

### Pipeline Architecture

Each step of the ML workflow was modularized as a Prefect task:

1. Load data from source
2. Clean and preprocess text using NLP pipeline
3. Perform train-test split with stratification
4. Apply TF-IDF vectorization
5. Train model with optimal hyperparameters
6. Evaluate using F1-score metrics

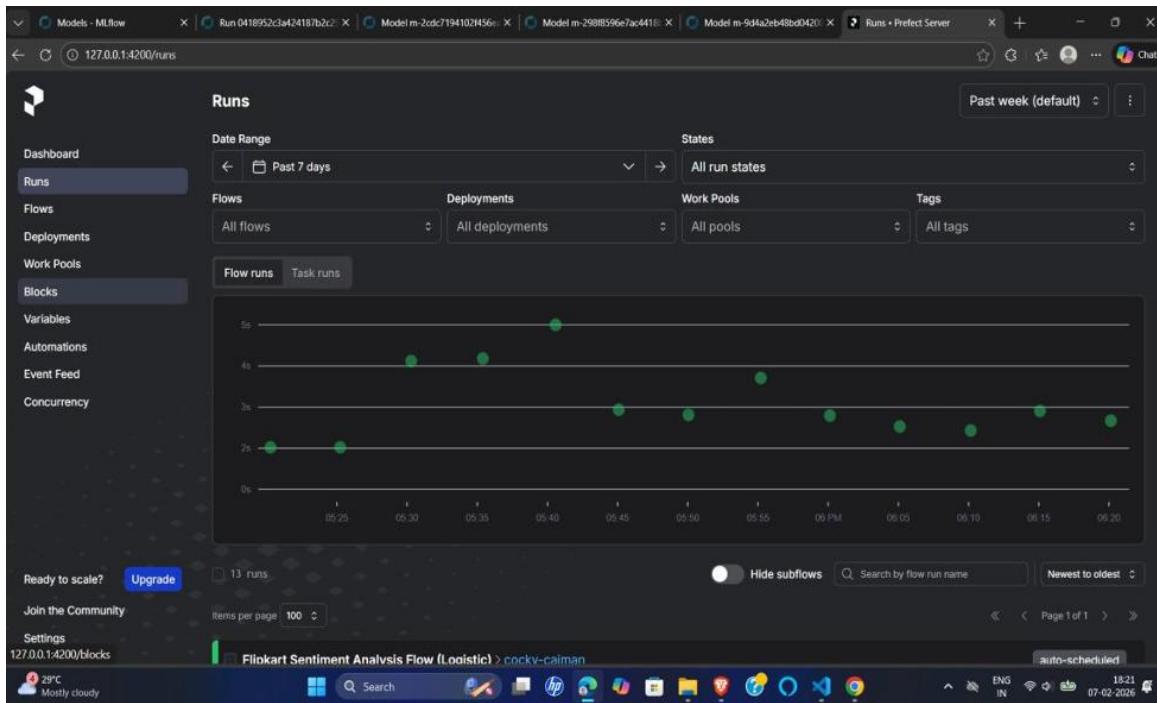
## 7. Log results to MLflow

# Prefect Architecture Breakdown

**Task Decorators (@task):** Each step of the pipeline is decorated with `@task`, making it a unit of work that Prefect can monitor, retry, and log independently.

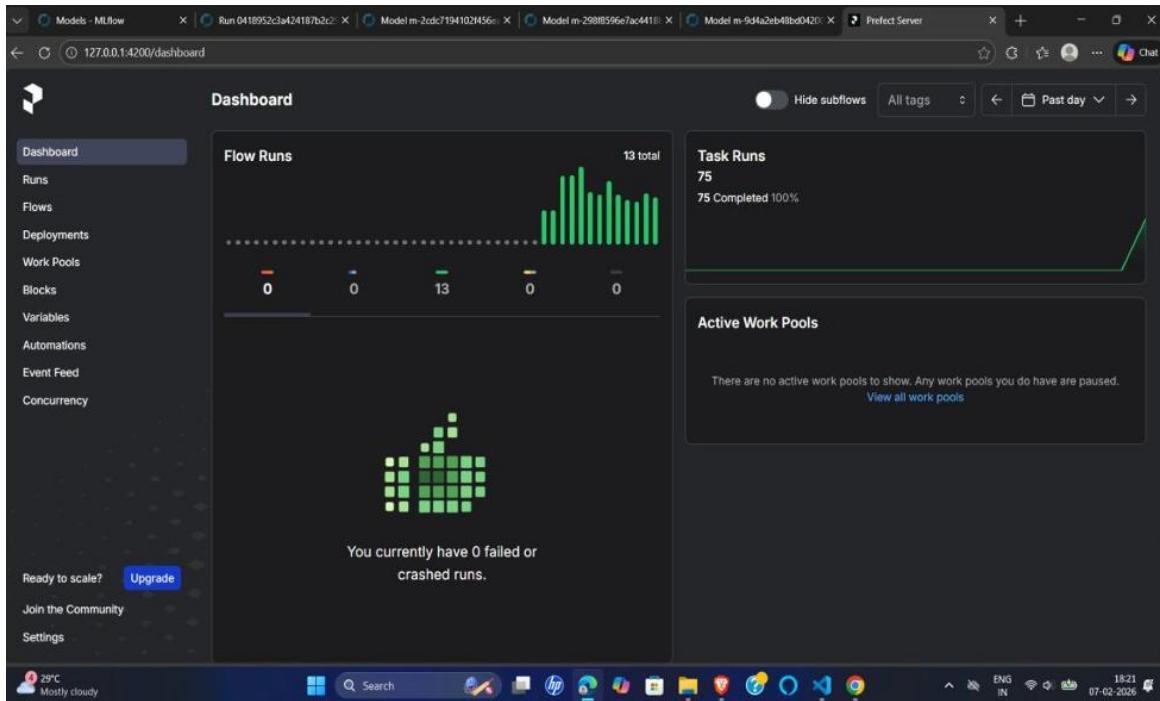
**Flow Decorator (@flow):** The `sentiment_workflow()` function orchestrates all tasks in the correct order, managing data dependencies automatically.

**Deployment with Scheduling:** The `.serve()` method deploys the flow with a cron schedule (`*/5 * * * *` = every 5 minutes), enabling continuous model retraining without manual intervention.



**Figure 14: Prefect Flow Runs Dashboard showing all successful executions**

The Prefect Runs dashboard shows 13 total flow runs over the past week, all successfully completed (shown in green). The timeline at the top visualizes execution patterns. Notice the deployment is tagged as 'auto-scheduled', indicating automated execution.



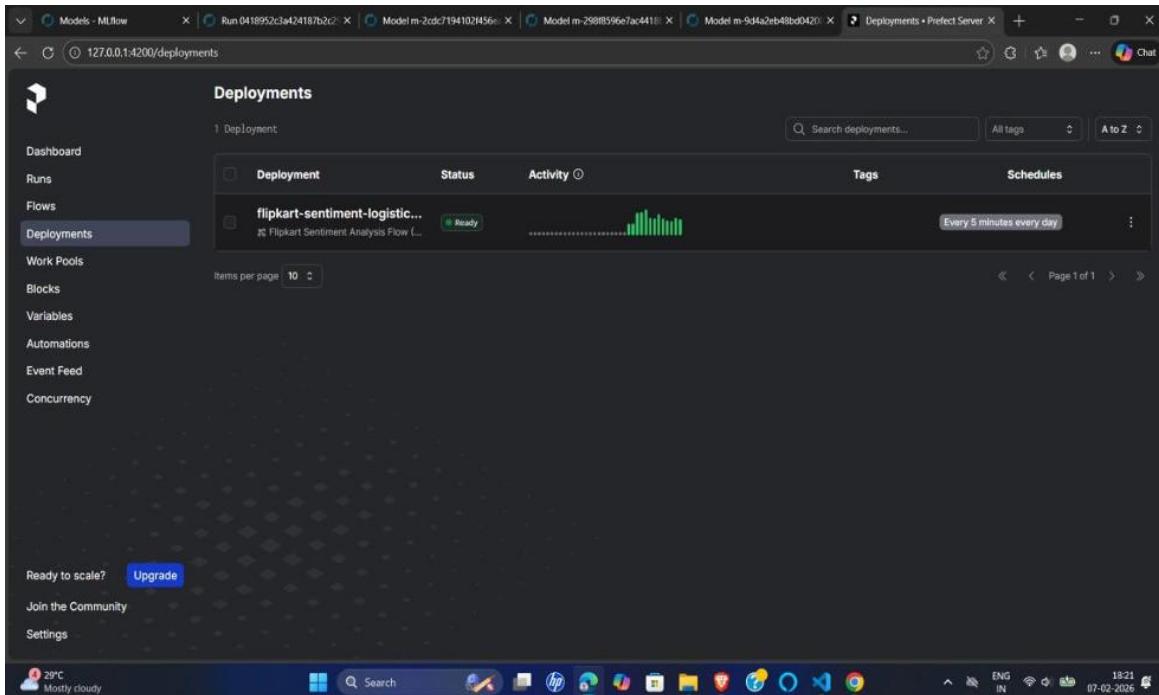
**Figure 15: Prefect Dashboard overview with Flow Runs and Task Runs metrics**

The Prefect Dashboard provides an overview of all workflow activity. The Flow Runs section shows 13 total runs with 0 failures—demonstrating pipeline reliability. The Task Runs section shows 75 completed tasks (100% success rate). The activity heatmap shows consistent daily execution.

## Scheduling & Deployment

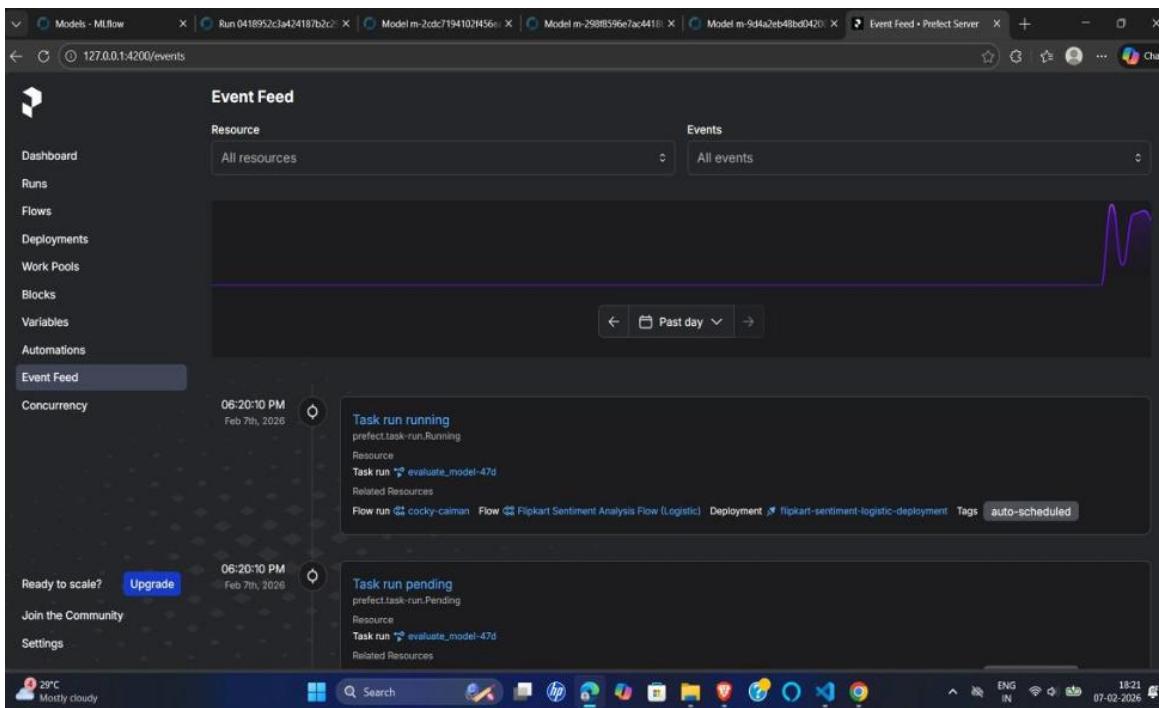
The workflow was deployed with automated scheduling:

- Cron-based scheduling for periodic re-training
- No manual intervention required
- Real-time monitoring via Prefect dashboard
- Comprehensive logging of all executions



**Figure 16: Prefect Deployment configuration showing scheduled execution every 5 minutes**

The Deployments page shows the 'flipkart-sentiment-logistic-deployment' configured to run every 5 minutes. This deployment is in 'Ready' status and linked to the 'Flipkart Sentiment Analysis Flow (Logistic)'. The activity chart shows consistent execution patterns.



## **Figure 17: Prefect Event Feed providing real-time workflow execution updates**

The Event Feed provides real-time updates on workflow execution. Task runs transition through states: 'running', 'pending', and completed. Each event shows the associated flow run, deployment, and tags. This granular logging helps with debugging and monitoring.

---

## **Why This Matters: MLOps Perspective**

This project demonstrates a real-world MLOps pipeline that bridges the gap between model development and production deployment:

**MLflow** → Comprehensive experiment tracking, model registry, and versioning

**Optuna** → Intelligent hyperparameter optimization with systematic search

**NLP Pipeline** → Reproducible text processing with careful feature engineering

**Prefect** → Workflow orchestration, scheduling, and automation

## **System Benefits**

This stack creates a system where:

- ✓ Every experiment is tracked and reproducible
  - ✓ Model selection is data-driven, not intuition-based
  - ✓ Deployment is automated and monitored
  - ✓ The entire pipeline can be versioned and shared
- 

## **Key Learnings**

### **1. Accuracy is Misleading for Imbalanced NLP Tasks**

F1-score provides a much more honest evaluation of model performance, especially when dealing with sentiment classification where class distribution may not be perfectly balanced.

## **2. Experiment Tracking Saves Enormous Debugging Time**

Being able to compare 63 different runs side-by-side, filter by parameters, and visualize metrics made it trivial to identify what works. Without MLflow, this would require manual spreadsheet tracking.

## **3. Workflow Orchestration Turns Scripts into Systems**

Prefect transformed a Jupyter notebook into a production pipeline. The ability to schedule, monitor, and version entire workflows is what separates exploratory work from production-grade ML.

## **4. Classical ML Models Still Perform Extremely Well for Text**

With proper preprocessing and feature engineering, Logistic Regression and Naive Bayes achieved 0.79+ F1 scores. These models train in seconds, are interpretable, and require minimal infrastructure compared to deep learning approaches.

## **5. Clean Preprocessing Matters More Than Complex Models**

The NLP pipeline (stopword removal, stemming, careful handling of negations) had a bigger impact on performance than trying different algorithms. Good features beat fancy models.

---



### **What's Next**

The immediate next steps to make this a complete production system:

## **1. Load MLflow Production Models Directly**

Instead of retraining, automatically pull the latest Production-stage model from the registry for inference.

## 2. Integrate with Flask Web Application

Create a REST API endpoint where users can submit reviews and receive real-time sentiment predictions.

## 3. Add Data Drift Monitoring

Track distribution changes in input features and model performance over time to detect when retraining is needed.

## 4. Containerize with Docker

Package the entire pipeline (Prefect + MLflow + model) into Docker containers for consistent deployment across environments.

## 5. Implement A/B Testing Framework

Compare different model versions in production to continuously validate improvements.

---

## ✨ Final Thoughts

**Building models is only half the job. Building systems around models is what makes them usable in the real world.**

This project demonstrates how modern MLOps tools work together:

- Optuna finds optimal configurations
- MLflow tracks and versions everything
- Prefect automates and monitors execution

The result is a **reproducible, maintainable, and production-ready** machine learning system.

Today's work was a solid step toward understanding what it truly means to build ML systems, not just ML models.

---



## Call to Action

If you're learning MLflow, Optuna, or Prefect—build projects like this.

*Nothing teaches better than wiring everything end-to-end.*

---



# Connect

**GitHub:** <https://github.com/nasir331786>

**LinkedIn:** <https://www.linkedin.com/in/nasir-husain-tamanne-9b9981377>

---