

Sales Forecasting Across Multiple Retail Stores

In this project, I developed a sales forecasting model for multiple retail stores. I began by cleaning the dataset, which contained over 1 million records, addressing missing values, removing duplicates, and handling outliers. This preparation enabled me to analyze sales trends over time.

For exploratory data analysis (EDA), I uncovered insights about customer behaviour & sales performance, & engineered features related to dates, promotions, and competition.

In the machine learning (ML) part, I built the forecasting model using Sklearn pipelines with the Random Forest Regressor.

In the deep learning (DL) part, I conducted the ADF test to assess the stationarity of the time series data. Then Built the LSTM model using TensorFlow Keras with MLflow autologging. Created a Python file for a Streamlit dashboard that allows users to upload CSV or Excel files, displaying the last available date, a data preview, & generating six-week sales predictions along with a prediction & feature importance graph.

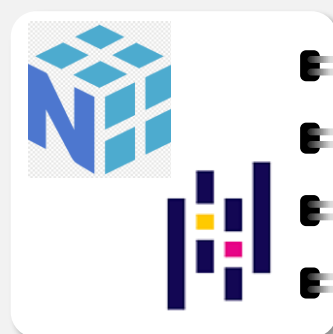
This project integrated data cleaning, EDA, ML & DL model building, and an interactive dashboard to provide a practical solution for sales forecasting.



by Debasis Baidya



Contributions of Key Libraries & Modules to the Project

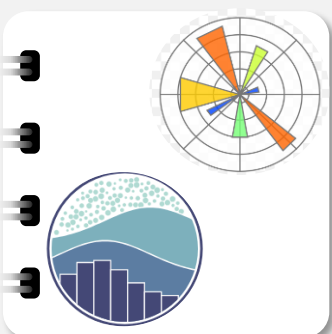


NumPy, Pandas:

I used NumPy for efficient numerical operations & Pandas for data manipulation, which made handling large datasets straightforward.

Matplotlib, Seaborn:

I leveraged Matplotlib and Seaborn to create insightful visualizations, helping me uncover patterns in sales data during my exploratory analysis.



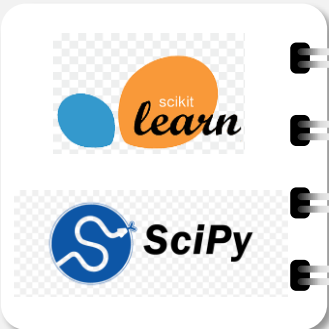
Logging, OS Module:

The Logging module helped me track the project's progress and debug issues, while the OS module managed file paths for data handling.



Sklearn, Scipy:

Sklearn provided tools for preprocessing and model evaluation, and Scipy offered essential statistical functions, which were vital for model development.



Statsmodels, Pickle:

I used Statsmodels for statistical analysis of sales data and Pickle to save my models, making it easy to load them later.

TensorFlow, Keras:

TensorFlow and Keras were crucial for building and training my LSTM model, allowing me to implement complex neural network architectures.



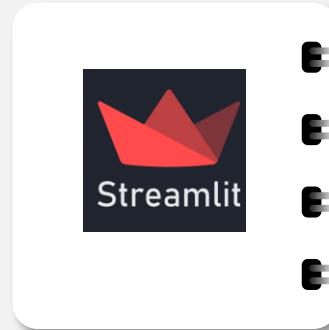
MLflow:

MLflow helped me track experiments and manage the machine learning lifecycle, ensuring reproducibility in my model evaluations.



Streamlit:

I used Streamlit to create an interactive web app for visualizing predictions, making it easy for stakeholders to explore the results.



Data Exploration, Cleaning, Processing Steps

1

Logging Configuration

I configured logging using Python's logging module to track the project's progress and errors, ensuring all important events were logged to a file for easy debugging.

2

Merging Train & Store Data

I merged the training dataset with store data during my exploratory data analysis (EDA) to better understand the relationships and patterns in the data.

3

Data Cleaning and Preprocessing

I focused on data cleaning and preprocessing to handle missing values, remove duplicates & standardize formats, ensuring the dataset was ready for analysis & modeling.

4

Feature Engineering

I performed feature engineering to create new variables from the existing data, which helped enhance the analysis during my exploratory data analysis (EDA).

5

Exploratory Data Analysis Visualizations

I plotted promo distribution, total sales by season, and sales distribution by promo etc. These visuals helped me identify key patterns during my EDA.

6

Dropping Feature Engineering Columns

I dropped the feature engineering columns from the DataFrame to revert Train Data to its default state, ensuring a clean dataset for further analysis and modeling.

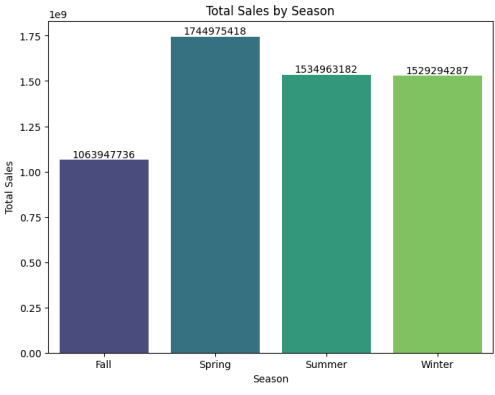
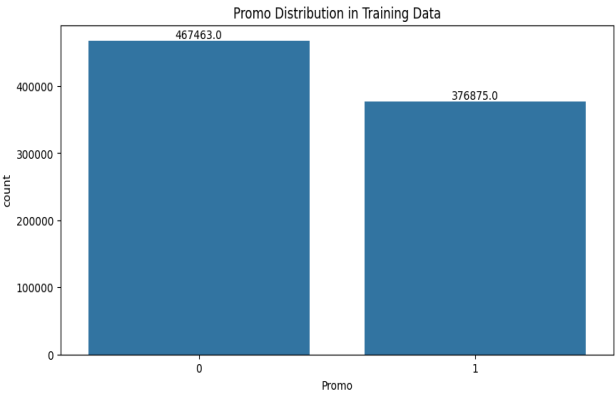
```
1 # Check for missing values and duplicates before handling them
2 print("Missing Values Count in Train Data:")
3 print(train_data.isnull().sum())
4
5 print("\nMissing Values Count in Test Data:")
6 print(test_data.isnull().sum())
7
8 # Check for duplicates
9 train_duplicates = train_data.duplicated().sum()
10 test_duplicates = test_data.duplicated().sum()
11
12 print(f"\nNumber of Duplicates in Train Data: {train_duplicates}")
13 print(f"Number of Duplicates in Test Data: {test_duplicates}")
```

Missing Values Count in Train Data:

Store	0
DayOfWeek	0
Date	0
Sales	0
Customers	0
Open	0
Promo	0
StateHoliday	0
SchoolHoliday	0
StoreType	0
Assortment	0
CompetitionDistance	2642
CompetitionOpenSinceMonth	323348
CompetitionOpenSinceYear	323348
Promo2	0
Promo2SinceWeek	508031
Promo2SinceYear	508031
PromoInterval	508031
dtype:	int64

Missing Values Count in Test Data:

Id	0
Store	0
DayOfWeek	0
Date	0
Open	11
Promo	0
StateHoliday	0
SchoolHoliday	0
dtype:	int64



Data Preparation & Visualization Summary

Renamed Customer Column

Renamed the "Customer" column in the training data to "Id" to match the test data column.

Visualized Distributions of Features

Plotted distributions for various features to understand their characteristics.

Boxplot Comparison

Created boxplots to compare train and test data distributions, analyzing feature variability and outliers.

Analyzed Feature Variability

Assessed the spread of features across train and test datasets to identify differences.

Identified Outliers

Used boxplots to pinpoint outliers that may influence model performance.

Checked for Consistency

Ensured that the feature distributions were aligned between train and test datasets.

Renaming Train Data Column to Match Test Data Column

```
# Rename 'Customers' column in Train Data to 'Id'
train_data = train_data.rename(columns={'Customers': 'Id'}) # Rename Customers to Id

# Create a list of columns that exist in both DataFrames
common_columns = test_data.columns.intersection(train_data.columns).tolist()

# Rearranging Train Data columns to match Test Data columns
train_data = train_data[common_columns + [col for col in train_data.columns if col not in common_columns]]

# Print columns present in Train & Test Data after rearranging:
logging.info("Train columns: {}".format(train_data.columns.tolist()))
logging.info("Test columns: {}".format(test_data.columns.tolist()))

2025-03-10 06:39:34,711 - INFO - Train columns: ['Id', 'Store', 'DayOfWeek', 'Date', 'Open', 'Promo', 'StateHoliday', 'SchoolHoliday', 'Sales']
2025-03-10 06:39:34,712 - INFO - Test columns: ['Id', 'Store', 'DayOfWeek', 'Date', 'Open', 'Promo', 'StateHoliday', 'SchoolHoliday']
```

Visualize Distributions of Features

```
import matplotlib.pyplot as plt

# Plot histogram for Train Data excluding 'Date' column
train_data.hist(column=train_data.select_dtypes(include=['number']).columns.difference(['Date']),
                figsize=(12, 8), bins=30, edgecolor='black')
plt.suptitle("Train Data Distribution")
plt.show()

# Plot histogram for Test Data excluding 'Date' column
test_data.hist(column=test_data.select_dtypes(include=['number']).columns.difference(['Date']),
                figsize=(12, 8), bins=30, edgecolor='black')
plt.suptitle("Test Data Distribution")
plt.show()
```

Boxplot Comparison of Train and Test Data Distributions: Analyzing Feature Variability and Outliers

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 import matplotlib.gridspec as gridspec
4
5 # Create a figure with GridSpec
6 fig = plt.figure(figsize=(15, 6))
7 gs = gridspec.GridSpec(1, 2, width_ratios=[1, 1]) # 1 row, 2 columns
8
9 # Boxplot for Train Data
10 ax1 = fig.add_subplot(gs[0])
11 sns.boxplot(data=train_data.select_dtypes(include=['number']), ax=ax1)
12 ax1.set_title("Boxplot of Train Data (Numerical Features)")
13 ax1.tick_params(axis='x', rotation=45)
14
15 # Boxplot for Test Data
16 ax2 = fig.add_subplot(gs[1])
17 sns.boxplot(data=test_data.select_dtypes(include=['number']), ax=ax2)
18 ax2.set_title("Boxplot of Test Data (Numerical Features)")
19 ax2.tick_params(axis='x', rotation=45)
20
21 # Adjust layout
22 plt.tight_layout()
23 plt.show()
```

Winsorization of Sales Data

1

Outlier Treatment

Implemented Winsorization on the 'Sales' column to limit extreme values, enhancing data robustness.

2

Data Integrity Check

Verified the transformation by displaying the original & winsorized sales values side by side.

3

Comparative Visualization

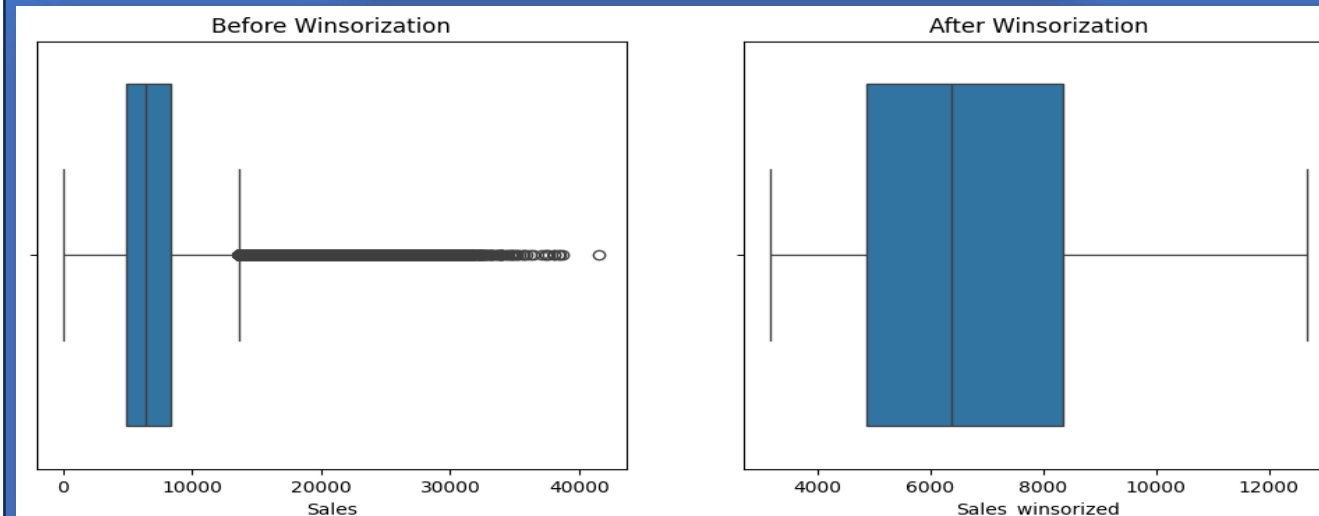
Created boxplots to illustrate the difference in sales distribution before and after Winsorization, highlighting reduced variability.

Capping Extreme Outliers with Winsorization: *Capping prevents extreme values from dominating the model*

```
[ ] 1 import numpy as np
    2 from scipy.stats.mstats import winsorize
    3
    4 # Apply Winsorization with higher limits to handle more outliers
    5 train_data['Sales_winsorized'] = winsorize(train_data['Sales'], limits=[0.05, 0.05]) # Cap 5% from both ends
    6
    7 # Convert winsorized values to integers if needed
    8 train_data['Sales_winsorized'] = train_data['Sales_winsorized'].astype(int)
    9
   10 # Display the first few entries to verify
   11 print(train_data[['Sales', 'Sales_winsorized']].head())
```

	Sales	Sales_winsorized
0	5263	5263
1	6064	6064
2	8314	8314
3	13995	12668
4	4822	4822

```
1 # Checking by plotting a boxplot:
2 # Plot Before & After Winsorization
3 plt.figure(figsize=(12, 5))
4
5 # Before Winsorization
6 plt.subplot(1, 2, 1)
7 sns.boxplot(x=train_data['Sales'])
8 plt.title("Before Winsorization")
9
10 # After Winsorization
11 plt.subplot(1, 2, 2)
12 sns.boxplot(x=train_data['Sales_winsorized'])
13 plt.title("After Winsorization")
14
15 plt.show()
```



Model Training With Sklearn Pipeline

2.2 Building models with sklearn pipelines

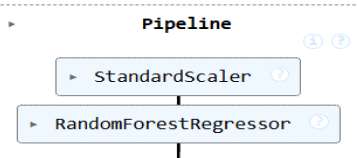
```
[ ] 1 # Import Required Libraries
    2 from sklearn.ensemble import RandomForestRegressor
    3 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
    4 import pickle
    5 from datetime import datetime
    6
    7 # Define features and target variable from Train Data
    8 X_train = train_data.drop(columns=['Sales', 'Sales_winsorized', 'Date']) # Drop unnecessary columns
    9 y_train = train_data['Sales_winsorized'] # Target variable
   10
   11 # Define Test Data
   12 X_test = test_data.drop(columns=['Sales', 'Date'], errors='ignore')
```

```
[ ] 1 # Check shapes
    2 print("X_train shape:", X_train.shape)
    3 print("y_train shape:", y_train.shape)
    4 print("X_test shape:", X_test.shape)
```

```
X_train shape: (844338, 7)
y_train shape: (844338,)
X_test shape: (41088, 7)
```

```
1 # Create Pipeline for Scaling & Model Training
2 from sklearn.pipeline import Pipeline
3 from sklearn.preprocessing import StandardScaler
4
5 ml_pipeline = Pipeline([
6     ('scaler', StandardScaler()),
7     ('model', RandomForestRegressor(n_estimators=100, random_state=42, n_jobs=-1))
8 ])
```

```
1 # Train Model Using Pipeline
2 ml_pipeline.fit(X_train, y_train)
```



2.3 Choose & Justify Loss Function

```
[ ] 1 # Compute Evaluation Metrics
    2 mae = mean_absolute_error(y_train, y_train_pred)
    3 mse = mean_squared_error(y_train, y_train_pred)
    4 rmse = np.sqrt(mse)
    5 r2 = r2_score(y_train, y_train_pred)
    6
    7 # Print Evaluation Results
    8 print("Model Evaluation Metrics (on Training Data):")
    9 print(f"Mean Absolute Error (MAE): {mae:.2f}")
   10 print(f"Mean Squared Error (MSE): {mse:.2f}")
   11 print(f"Root Mean Squared Error (RMSE): {rmse:.2f}")
   12 print(f"R² Score: {r2*100:.2f}%")
```

```
Model Evaluation Metrics (on Training Data):
Mean Absolute Error (MAE): 152.53
Mean Squared Error (MSE): 51475.51
Root Mean Squared Error (RMSE): 226.88
R² Score: 99.22%
```

The chosen loss function is MAE (Mean Absolute Error) because:

- MAE (152.53): The average difference between predicted and actual values is approximately 152.53 units.
- MSE (51475.51): The average squared difference between predicted and actual values is approximately 51475.51 units squared.
- RMSE (226.88): The square root of MSE, representing the average magnitude of errors, is approximately 226.88 units.
- R² Score (99.22%): The model explains about 99.22% of the variance in the training data, indicating an excellent fit.

Note: Although MSE seems suitable based on the high R² score and relatively low error values, **MAE** is chosen as the loss function.

01

Pipeline Construction

Built a Random Forest Regressor Pipeline to streamline the training process.

02

Feature Scaling

Implemented feature scaling to normalize the input data, enhancing model performance.

03

Training Process

Trained the pipeline on the processed data, ensuring that it learned from relevant features.

04

Performance Evaluation Metrics

Evaluated the model's performance using several metrics, including Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE).

05

Predictive Ability

Achieved an R² Score of 99.22, indicating a strong predictive ability and excellent fit to the training data.

Sales Prediction Steps for the Last Date

```
1 from datetime import timedelta
2
3 # Step 1: Identify the last date
4 last_date = train_data['Date'].max()
5 print("Last Date in Train Data:", last_date)
6
7 # Step 2: Filter data for the last date
8 last_date_data = train_data[train_data['Date'] == last_date]
9 print("Data for Last Date:")
10 print(last_date_data)
11
12 # Step 3: Prepare features (X) for prediction
13 X_last_date = last_date_data.drop(columns=['Sales', 'Sales_winsorized', 'Date'])
14 print("Features for Last Date:")
15 print(X_last_date)
16
17 # Step 4: Make predictions for the last date
18 y_last_date_pred = ml_pipeline.predict(X_last_date)
19 print("Predicted Sales for Last Date:", y_last_date_pred)
```

Last Date in Train Data: 2015-07-31 00:00:00

Prediction of Sales in various stores up to 6 weeks ahead of time using ML Pipeline

```
1 from datetime import datetime, timedelta
2
3 # Step 5: Generate future dates (next 6 weeks)
4 last_date = train_data['Date'].max()
5 future_dates = [last_date + timedelta(weeks=i) for i in range(1, 7)]
6 print("Future Dates:")
7 for i, date in enumerate(future_dates, start=1):
8     print(f"Next Week {i} from last date: {date}")
9
10 # Step 6: Prepare future data for predictions
11 future_data = []
12 feature_columns = X_train.columns # Get the feature column names from the training data
13 unique_stores = train_data['Store'].unique() # Get all unique store IDs
14
15 for date in future_dates:
16     for store in unique_stores:
17         # Create a dictionary for the future data
18         future_entry = {
19             'Date': date,
20             'Store': store
21         }
22         # Add the feature columns with placeholder values
23         for col in feature_columns:
24             future_entry[col] = 0
25         future_data.append(future_entry)
26
27 # Convert to DataFrame
28 future_df = pd.DataFrame(future_data)
29
30 # Ensure the feature columns are correctly aligned
31 future_df['Date'] = pd.to_datetime(future_df['Date'])
32
33 # Predict future sales
34 y_future_pred = ml_pipeline.predict(future_df[feature_columns])
35 future_df['Predicted_Sales'] = y_future_pred
36
37 # Print the future predictions
38 print("\nFuture Predictions:")
39 for i, (date, store, pred) in enumerate(zip(future_dates, future_df['Store'], y_future_pred), start=1):
40     print(f"Next Week {i} from last date: {date} - Predicted Sales: {int(pred)}")
```

Predictions From The Last Date To The Next 6 Weeks

1

Last Date Analysis

The last training date is July 31, 2015, with significant sales variability across stores, indicating diverse consumer preferences.

2

Sales Variability

Sales ranged from 4,822 to 27,508, highlighting the impact of location and store strategies.

3

Future Predictions

Store 0's future sales are consistently predicted at 3,174 for the next 6 weeks, suggesting a stable baseline.

Post-Prediction Analysis & Confidence Intervals

Feature Importance Analysis

Identified key predictors, with Store, Promo, and DayOfWeek having the highest impact on sales.

Analysis & Confidence Intervals

Estimated confidence intervals (5th & 95th percentile) for predictions, providing uncertainty estimates.

Model Serialization

Saved the trained model as a timestamped .pkl file for future deployment. “rf_model_2025-03-10-06-41-47.pkl”.

Dropping Additional Columns

Dropped Unwanted Columns: Predicted_Sales (irrelevant for raw test data), Lower & Upper_Bound (confidence interval bounds not needed for predictions).

Renaming "Sales_winsorized"

The column 'Sales_winsorized' in train_data has been renamed to 'Sales'.

Generated Files

train_file.csv (Processed Training Data)
test_file.csv (Processed Test Data)

2.4 Post Prediction Analysis - Feature Importance

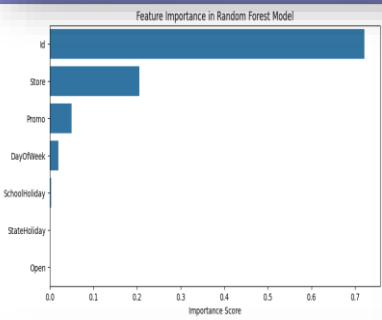
```
1 # Extract feature importance from the trained Random Forest model
2 feature_importances = ml_pipeline.named_steps['model'].feature_importances_
3
4 # Create DataFrame for visualization
5 feature_importance_df = pd.DataFrame({'Feature': X_train.columns, 'Importance': feature_importances})
6 feature_importance_df = feature_importance_df.sort_values(by='Importance', ascending=False)
7
8 # Plot Feature Importance
9 plt.figure(figsize=(10, 5))
10 sns.barplot(x=feature_importance_df['Importance'], y=feature_importance_df['Feature'])
11 plt.title("Feature Importance in Random Forest Model")
12 plt.xlabel("Importance Score")
13 plt.ylabel("Feature Name")
14 plt.show()
```

2.4 Post Prediction Analysis - Confidence Interval Estimation

```
1 # Generate predictions using multiple trees in the Random Forest
2 all_predictions = np.array([tree.predict(X_test) for tree in ml_pipeline.named_steps['model'].estimators_])
3
4 # Compute confidence intervals (5th and 95th percentile)
5 lower_bound = np.percentile(all_predictions, 5, axis=0)
6 upper_bound = np.percentile(all_predictions, 95, axis=0)
7
8 # Store confidence intervals in test data
9 test_data['Predicted_Sales'] = ml_pipeline.predict(X_test)
10 test_data['Lower_Bound'] = lower_bound
11 test_data['Upper_Bound'] = upper_bound
12
13 # Save predictions with confidence intervals
14 test_data[['Predicted_Sales', 'Lower_Bound', 'Upper_Bound']].to_csv('predicted_sales_with_ci.csv', index=False)
15
16 # ----- Print Output to Verify -----
17 print("Sample Predictions with Confidence Intervals:")
18 print(test_data[['Predicted_Sales', 'Lower_Bound', 'Upper_Bound']].head())
```

Sample Predictions with Confidence Intervals:

	Predicted_Sales	Lower_Bound	Upper_Bound
0	3174.0	9120.4	12668.0
1	3174.0	9408.0	12668.0
2	3174.0	12668.0	12668.0
3	3174.0	12668.0	12668.0
4	3174.0	12668.0	12668.0



2.5 Serialize the Model: Save Trained Model as .pkl File with Timestamp

```
1 import pickle
2 from datetime import datetime
3
4 # Save Trained Model as '.pkl' File with Timestamp
5 model_filename = f"rf_model_{datetime.now().strftime('%Y-%m-%d-%H-%M-%S')}.pkl"
6 pickle.dump(ml_pipeline, open(model_filename, 'wb'))
7
8 print(f"Model saved as {model_filename}")
```

Model saved as rf_model_2025-03-10-06-41-47.pkl

Dropping Additional Columns Created and Retaining Original Columns to Save Train and Test Data for Future Reference

```
1 # Drop 'Sales_log' from train_data before saving
2 train_data = train_data.drop(columns=['Sales', 'Sales_log'], errors='ignore')
3
4 # Drop 'Predicted_Sales', 'Lower_Bound', 'Upper_Bound' from test_data before saving
5 test_data = test_data.drop(columns=['Predicted_Sales', 'Lower_Bound', 'Upper_Bound'], errors='ignore')
6
7 # Print columns of Train & Test Datasets
8 print("Train File\n", train_data.columns.to_list())
9 print("Test File\n", test_data.columns.to_list())
```

Train File
['Id', 'Store', 'DayOfWeek', 'Date', 'Open', 'Promo', 'StateHoliday', 'SchoolHoliday', 'Sales_winsorized']
Test File
['Id', 'Store', 'DayOfWeek', 'Date', 'Open', 'Promo', 'StateHoliday', 'SchoolHoliday']

Time Series Analysis



Sales Trend

Clear seasonal patterns with periodic peaks and drops.



Stationary Data

ADF test confirms sales data is stationary (p-value: 0.0001).



ACF/PACF

Strong weekly seasonality (lags of 7, 14, 21, 28).



Residual Analysis

Examined residuals for heteroscedasticity and autocorrelation.



Decomposition

Identified trend, seasonal, and residual components.



30-Day Sliding Window Forecast

Different sequences overlap, reflecting diverse sales patterns across various periods.

Aggregate sales by date for time series analysis

```
1  ### Step 3: Aggregate sales by date for time series analysis
2
3  import matplotlib.pyplot as plt
4
5  # Aggregate sales by date for time series analysis
6  sales_trend = train_df.groupby("Date")["Sales"].sum()
7
8  plt.figure(figsize=(12,5))
9  plt.plot(sales_trend, label='Sales Trend')
10 plt.xlabel("Date")
11 plt.ylabel("Sales")
12 plt.title("Sales Trend Over Time")
13 plt.legend()
14 plt.show()
```

Autocorrelation (ACF) & Partial Autocorrelation (PACF) Analysis

```
1  import matplotlib.pyplot as plt
2  from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
3
4  # Plot ACF
5  plt.figure(figsize=(12, 6))
6  plot_acf(sales_trend, lags=30, ax=plt.gca())
7  plt.title('Autocorrelation Function (ACF)')
8  plt.show()
9
10 # Plot PACF
11 plt.figure(figsize=(12, 6))
12 plot_pacf(sales_trend, lags=30, ax=plt.gca())
13 plt.title('Partial Autocorrelation Function (PACF)')
14 plt.show()
```

```
1  ### Step 4: Perform ADF Test for Stationarity
2  from statsmodels.tsa.stattools import adfuller
3
4  # Perform the ADF test on the sales trend data
5  adf_result = adfuller(sales_trend)
6
7  # Output the ADF statistic and p-value
8  print(f"ADF Statistic: {adf_result[0]}")
9  print(f"p-value: {adf_result[1]}")
10
11 # Interpret the result
12 if adf_result[1] < 0.05:
13     print("Time series is stationary.")
14 else:
15     print("Time series is not stationary, differencing may be required.")
```

ADF Statistic: -4.6062197626111505
p-value: 0.00012580255594144246
Time series is stationary.

Create 30-Day Sliding Window Forecast

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  # Function to create sliding window data for visualization
5  def create_sliding_window_forecast(df, seq_length=30, sample_size=5):
6      forecast_data = []
7      indices = np.linspace(0, len(df) - seq_length - 1, sample_size, dtype=int) # Pick evenly spaced sequences
8
9      for i in indices:
10         forecast_data.append(df[['Sales']].iloc[i:i + seq_length].values)
11
12     return np.array(forecast_data)
13
14 # Get sample sliding window sequences
15 X_sliding_sample = create_sliding_window_forecast(train_df, sample_size=5) # Only visualize 5 sequences
16
17 # Plot the sequences
18 plt.figure(figsize=(12, 6))
19 for i, sequence in enumerate(X_sliding_sample):
20     plt.plot(sequence, label=f'Sequence {i+1}')
21
22 plt.title("Visualization of 30-Day Sliding Window Forecast")
23 plt.xlabel("Days")
24 plt.ylabel("Sales")
25 plt.legend()
26 plt.show()
```

Building the LSTM Model for Time Series Forecasting

Data Preparation

Converted 'Date' to numerical format
(days since the first date).
Selected key features for LSTM
training.

Selected Features

Features: ['Id', 'Store', 'DayOfWeek',
'Date', 'Open', 'Promo', 'StateHoliday',
'SchoolHoliday']

Target: ['Sales']

Key Takeaways

Date conversion is crucial for LSTM.

Diverse features enhance model
performance.

Conclusion

Features: ['Id', 'Store', 'DayOfWeek',
'Date', 'Open', 'Promo', 'StateHoliday',
'SchoolHoliday']

Target: ['Sales']

Building the LSTM Model for Time Series Forecasting

```
1 # Prepare Data for LSTM Training
2 from sklearn.preprocessing import MinMaxScaler
3 import numpy as np
4 import pandas as pd
5
6 # Convert Date to datetime format
7 train_df['Date'] = pd.to_datetime(train_df['Date'])
8
9 # Define feature and target columns
10 feature_columns = ['Id', 'Store', 'DayOfWeek', 'Date', 'Open', 'Promo', 'StateHoliday', 'SchoolHoliday']
11 target_column = 'Sales'
12
13 # Extract features and target variable
14 X = train_df[feature_columns].copy()
15 y = train_df[target_column].values
16
17 # Display Max Columns
18 pd.set_option('display.max_columns', None)
19
20 # Print the first 5 rows of X
21 print("First 5 rows of Feature Matrix X:")
22 print(X[:5])
23
24 # Print the first 5 values of y
25 print("\nFirst 5 values of Target Array y:")
26 print(y[:5])
```

First 5 rows of Feature Matrix X:

	Id	Store	DayOfWeek	Date	Open	Promo	StateHoliday	SchoolHoliday
0	555	1	5	2015-07-31	1	1	0	1
1	625	2	5	2015-07-31	1	1	0	1
2	821	3	5	2015-07-31	1	1	0	1
3	1498	4	5	2015-07-31	1	1	0	1
4	559	5	5	2015-07-31	1	1	0	1

First 5 values of Target Array y:
[5263 6064 8314 12668 4822]

Data Scaling for Better LSTM Performance

Purpose of Scaling

Scaled numerical features and target variable to a range of -1 to 1 using MinMaxScaler.

Scaling Process

Created two MinMaxScaler objects: scaler_X: For numerical features. | scaler_y: For target variable (y).

Applied Scaling to Numerical Columns

Id, Store, DayOfWeek, Open, Promo, StateHoliday, SchoolHoliday.

Excluded the Date column from scaling.

Reshaped Scaled Input Data

Reshaped scaled input data (X_scaled) to the format required by LSTM: (samples, timesteps, features)

Outcome

Numerical features and target variable successfully scaled to -1 to 1.

Input data reshaped for LSTM training.

Conclusion

Scaling improves stability and efficiency of LSTM training.

Ensures all features and target variables are on the same scale, enhancing the model's learning process.

Data Scaling for better LSTM Performance

```
1 # Scale data
2 from sklearn.preprocessing import MinMaxScaler
3
4 # Scale numerical features (excluding Date) and target variable to range (-1, 1) for better training stability
5 scaler_X = MinMaxScaler(feature_range=(-1, 1))
6 scaler_y = MinMaxScaler(feature_range=(-1, 1))
7
8 # Scale all columns except 'Date'
9 X_scaled = X.copy()
10 num_cols = ['Id', 'Store', 'DayOfWeek', 'Open', 'Promo', 'StateHoliday', 'SchoolHoliday']
11 X_scaled[num_cols] = scaler_X.fit_transform(X[num_cols])
12 y_scaled = scaler_y.fit_transform(y.reshape(-1, 1))
13
14 # Reshape Data for LSTM (samples, timesteps, features)
15 X_scaled = np.array(X_scaled[num_cols]).reshape((X_scaled.shape[0], 1, len(num_cols)))
16
17 print("\nScaled Target Array (y_scaled) with shape:", y_scaled.shape)
18 print(y_scaled)
```

Scaled Target Array (y_scaled) with shape: (844338, 1)

```
[[-0.55993259]
 [-0.39119444]
 [ 0.08278913]
 ...
 [-0.60796292]
 [-0.72256162]
 [-0.41289235]]
```


Build Optimized LSTM Model

STEP 01

Model Framework

Utilized **Sequential** model from **TensorFlow Keras**.

STEP 02

Layer Implementation

First LSTM Layer:

- 64 units
- Dropout (30%) to mitigate overfitting.

Second LSTM Layer:

- 64 units
- Dropout (30%)

STEP 03

Output Layer

Dense layer with **linear activation** for regression tasks.

STEP 04

Model Compilation

- Optimizer:** Adam for efficient training.
- Loss Function:** Mean Squared Error (MSE) for regression.

STEP 05

Model Summary

LSTM Layer 1: (None, 1, 64) | Params: 18K | Dropout 1: (None, 1, 64)
LSTM Layer 2: (None, 64) | Params: 33K | Dropout 2: (None, 64)
Dense Layer: (None, 1) | Params: 65 | Total Parameters: 52K

STEP 06

Conclusion

Optimized LSTM regression model designed to capture sequential patterns while minimizing overfitting. 🚀

Build Optimized LSTM Model

```
1 from tensorflow.keras.models import Sequential
2 from tensorflow.keras.layers import LSTM, Dense, Dropout
3
4 # Define LSTM model architecture
5 model = Sequential([
6     LSTM(64, activation='tanh', return_sequences=True, input_shape=(1, len(num_cols))),
7     Dropout(0.3),
8     LSTM(64, activation='tanh'),
9     Dropout(0.3),
10    Dense(1, activation='linear') # Linear activation for regression
11 ])
12
13 # Compile the model using Adam optimizer and Mean Squared Error loss function
14 model.compile(optimizer='adam', loss='mse', metrics=['mae'])
15
16 # Print confirmation message
17 print("Optimized LSTM regression model with two layers has been built and compiled successfully.\n")
18
19 # Print Model Summary
20 model.summary()
```

Optimized LSTM regression model with two layers has been built and compiled successfully.

Model: "sequential"

Layer (type)	Output Shape	Param #
lstm (LSTM)	(None, 1, 64)	18,432
dropout (Dropout)	(None, 1, 64)	0
lstm_1 (LSTM)	(None, 64)	33,024
dropout_1 (Dropout)	(None, 64)	0
dense (Dense)	(None, 1)	65

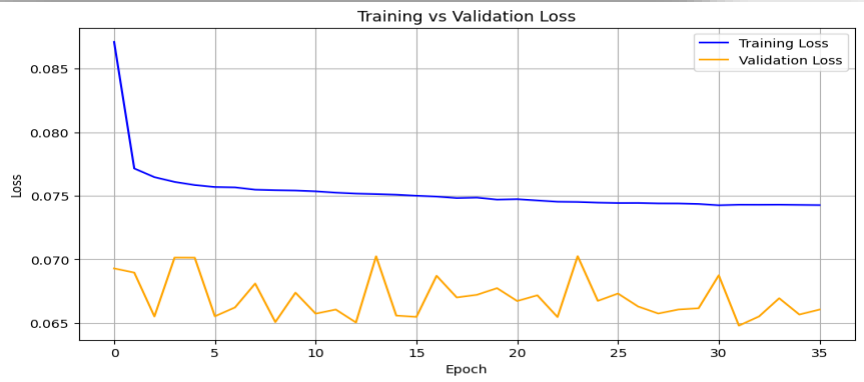
Total params: 51,521 (201.25 KB)
Trainable params: 51,521 (201.25 KB)
Non-trainable params: 0 (0.00 B)

LSTM Model Training & MLflow Logging

Training the LSTM Model and Logging with MLflow

```
1 # Uninstall & Re-Install the mlflow Library
2 !pip uninstall mlflow
3 !pip install mlflow

1 # Train Model and Log with MLflow
2 import mlflow.tensorflow
3 import matplotlib.pyplot as plt
4 from tensorflow.keras.callbacks import EarlyStopping
5 import os
6
7 # Prevent GPU errors in certain environments
8 os.environ["CUDA_VISIBLE_DEVICES"] = "-1"
9
10 # Enable automatic logging with MLflow
11 mlflow.tensorflow.autolog()
12
13 # Early stopping to prevent overfitting
14 early_stop = EarlyStopping(monitor='loss', patience=5, restore_best_weights=True)
15
16 # Train the model and save history
17 with mlflow.start_run():
18     history = model.fit(X_scaled, y_scaled, epochs=100, batch_size=128, validation_split=0.2, callbacks=[early_stop])
19     mlflow.tensorflow.log_model(model, "lstm_model")
20
21 # Plot Training vs Validation Loss
22 plt.figure(figsize=(10, 5))
23 plt.plot(history.history['loss'], label='Training Loss', color='blue')
24 plt.plot(history.history['val_loss'], label='Validation Loss', color='orange')
25 plt.title('Training vs Validation Loss')
26 plt.ylabel('Loss')
27 plt.xlabel('Epoch')
28 plt.legend()
29 plt.grid(True)
30 plt.show()
```



01

Model

LSTM for future sales prediction using time series data.

02

Tracking

Utilized MLflow autologging for performance metrics.

03

Key Techniques

EarlyStopping (patience = 5) to prevent overfitting.
Disabled GPU: `os.environ["CUDA_VISIBLE_DEVICES"] = "-1"`.
Trained for 100 epochs, batch size of 128, 20% validation split.

04

Outcomes

Training Loss: Decreased steadily; final: 0.0743 (MAE: 0.2110).
Validation Loss: Stabilized after initial fluctuations; final: 0.0660 (MAE: 0.1970).

05

Analysis

Training Curve: Smooth decline, indicating effective learning.
Validation Curve: Initial spikes but overall stabilization, suggesting good generalization.

06

Key Takeaways

- ✓ Effective learning with minimal overfitting.
- ✓ Minor validation fluctuations are expected but manageable.

07

Conclusion

LSTM model captured sequential patterns well; training logged with MLflow for reproducibility. 🚀

Determining Last Available Date & Saving the Model & Scaler



Overview

Developed an LSTM model to predict future sales based on historical data.

Model & Scaler Saving

- ❑ Saved LSTM model as **lstm.keras**.
- ❑ MinMaxScaler objects saved as **scaler_X.pkl** and **scaler_y.pkl**

Last Training Data Date

Last available date in training data is: 2015-07-31

Determine the last available date in training data

```
1 import pandas as pd
2
3 # Check the first few rows of the Date column
4 print(train_df['Date'].head())
5
6 # Check the data type of the Date column
7 print(train_df['Date'].dtype)
8
9 # Step 2: Identify the last date
10 last_date = train_df['Date'].max()
11 print("\nLast Date in Train Data:", last_date)
```

```
0 2015-07-31
1 2015-07-31
2 2015-07-31
3 2015-07-31
4 2015-07-31
Name: Date, dtype: datetime64[ns]
datetime64[ns]

Last Date in Train Data: 2015-07-31 00:00:00
```

Predict Next 6 Weeks Sales

```
1 from datetime import timedelta
2
3 # Determine the last available date in training data
4 last_date = train_df['Date'].max()
5
6 # Generate future dates for prediction (next 6 weeks = 42 days)
7 pred_dates = [last_date + timedelta(days=i) for i in range(1, 43)]
8
9 # Make predictions
10 y_pred_scaled = model.predict(X_scaled)
11 y_pred = scaler_y.inverse_transform(y_pred_scaled)
12
13 # Print lengths for debugging
14 print(f"Length of pred_dates: {len(pred_dates)}")
15 print(f"Length of y_pred: {len(y_pred)}")
16
17 # Print future dates alongside predicted values
18 print("\nFuture Prediction Dates and Predicted Sales:")
19 for i in range(min(len(pred_dates), len(y_pred))): # Adjusting to the minimum length
20     print(f"{pred_dates[i]} - Predicted Sales: {y_pred[i]}")
```

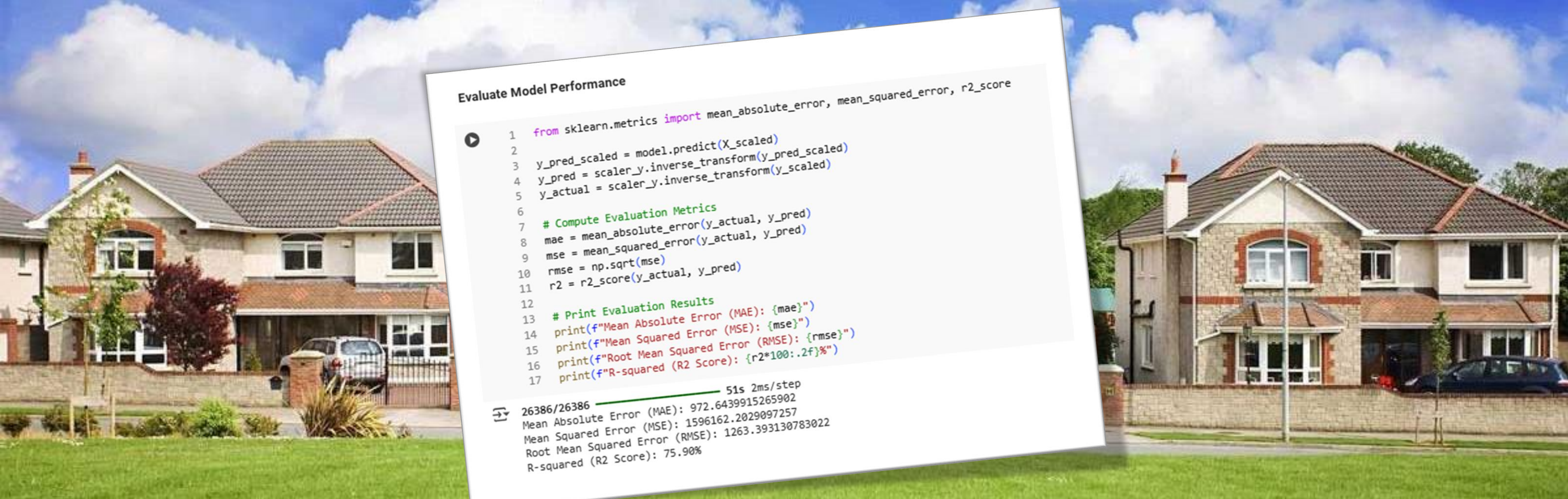
26386/26386 ————— 60s 2ms/step
Length of pred_dates: 42
Length of y_pred: 844338

Saving the Trained LSTM Model and Scaler

```
1 # Save trained model and scalers
2 import pickle
3
4 # Define filenames
5 model_filename = "lstm.keras"
6 scaler_x_filename = "scaler_X.pkl"
7 scaler_y_filename = "scaler_y.pkl"
8
9 # Save the model and scalers
10 model.save(model_filename)
11 with open(scaler_x_filename, "wb") as f:
12     pickle.dump(scaler_X, f)
13 with open(scaler_y_filename, "wb") as f:
14     pickle.dump(scaler_y, f)
15
16 # Confirm saving with formatted strings using placeholders
17 print(f"Model '{model_filename}' saved successfully.")
18 print(f"Scaler '{scaler_x_filename}' saved successfully.")
19 print(f"Scaler '{scaler_y_filename}' saved successfully.")
```

Model 'lstm.keras' saved successfully.
Scaler 'scaler_X.pkl' saved successfully.
Scaler 'scaler_y.pkl' saved successfully.

Future Prediction Dates and Predicted Sales:	
2015-08-01 00:00:00	- Predicted Sales: [6088.965]
2015-08-02 00:00:00	- Predicted Sales: [6720.2705]
2015-08-03 00:00:00	- Predicted Sales: [8231.525]
2015-08-04 00:00:00	- Predicted Sales: [11670.013]
2015-08-05 00:00:00	- Predicted Sales: [6125.136]
2015-08-06 00:00:00	- Predicted Sales: [6401.5405]
2015-08-07 00:00:00	- Predicted Sales: [11404.154]
2015-08-08 00:00:00	- Predicted Sales: [8313.937]
2015-08-09 00:00:00	- Predicted Sales: [7229.507]
2015-08-10 00:00:00	- Predicted Sales: [7181.1855]
2015-08-11 00:00:00	- Predicted Sales: [10688.756]
2015-08-12 00:00:00	- Predicted Sales: [9173.822]
2015-08-13 00:00:00	- Predicted Sales: [6410.3784]
2015-08-14 00:00:00	- Predicted Sales: [7408.187]
2015-08-15 00:00:00	- Predicted Sales: [7830.8154]
2015-08-16 00:00:00	- Predicted Sales: [9279.649]
2015-08-17 00:00:00	- Predicted Sales: [9069.689]
2015-08-18 00:00:00	- Predicted Sales: [9004.695]
2015-08-19 00:00:00	- Predicted Sales: [7468.012]
2015-08-20 00:00:00	- Predicted Sales: [9311.368]
2015-08-21 00:00:00	- Predicted Sales: [7184.624]
2015-08-22 00:00:00	- Predicted Sales: [6954.079]
2015-08-23 00:00:00	- Predicted Sales: [6127.023]
2015-08-24 00:00:00	- Predicted Sales: [9891.59]
2015-08-25 00:00:00	- Predicted Sales: [11897.265]
2015-08-26 00:00:00	- Predicted Sales: [6588.5454]
2015-08-27 00:00:00	- Predicted Sales: [10805.882]
2015-08-28 00:00:00	- Predicted Sales: [7027.8447]
2015-08-29 00:00:00	- Predicted Sales: [7608.9785]
2015-08-30 00:00:00	- Predicted Sales: [6708.262]
2015-08-31 00:00:00	- Predicted Sales: [6877.4985]
2015-09-01 00:00:00	- Predicted Sales: [6965.9214]
2015-09-02 00:00:00	- Predicted Sales: [8861.8955]
2015-09-03 00:00:00	- Predicted Sales: [10320.075]
2015-09-04 00:00:00	- Predicted Sales: [7903.761]
2015-09-05 00:00:00	- Predicted Sales: [10765.211]
2015-09-06 00:00:00	- Predicted Sales: [8680.574]
2015-09-07 00:00:00	- Predicted Sales: [6557.0664]
2015-09-08 00:00:00	- Predicted Sales: [8951.579]
2015-09-09 00:00:00	- Predicted Sales: [7144.589]
2015-09-10 00:00:00	- Predicted Sales: [5707.478]
2015-09-11 00:00:00	- Predicted Sales: [10674.02]



Evaluate Model Performance

```
1 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
2
3 y_pred_scaled = model.predict(X_scaled)
4 y_pred = scaler_y.inverse_transform(y_pred_scaled)
5 y_actual = scaler_y.inverse_transform(y_scaled)
6
7 # Compute Evaluation Metrics
8 mae = mean_absolute_error(y_actual, y_pred)
9 mse = mean_squared_error(y_actual, y_pred)
10 rmse = np.sqrt(mse)
11 r2 = r2_score(y_actual, y_pred)
12
13 # Print Evaluation Results
14 print(f"Mean Absolute Error (MAE): {mae}")
15 print(f"Mean Squared Error (MSE): {mse}")
16 print(f"Root Mean Squared Error (RMSE): {rmse}")
17 print(f"R-squared (R2 Score): {r2*100:.2f}%")
```

```
26386/26386 51s 2ms/step
Mean Absolute Error (MAE): 972.6439915265902
Mean Squared Error (MSE): 1596162.2029097257
Root Mean Squared Error (RMSE): 1263.393130783022
R-squared (R2 Score): 75.90%
```



FUTURE PREDICTIONS

- Generated dates for the next 6 weeks (42 days): from 2015-08-01 to 2015-09-11.
- Predictions scaled back to original values.



PERFORMANCE METRICS

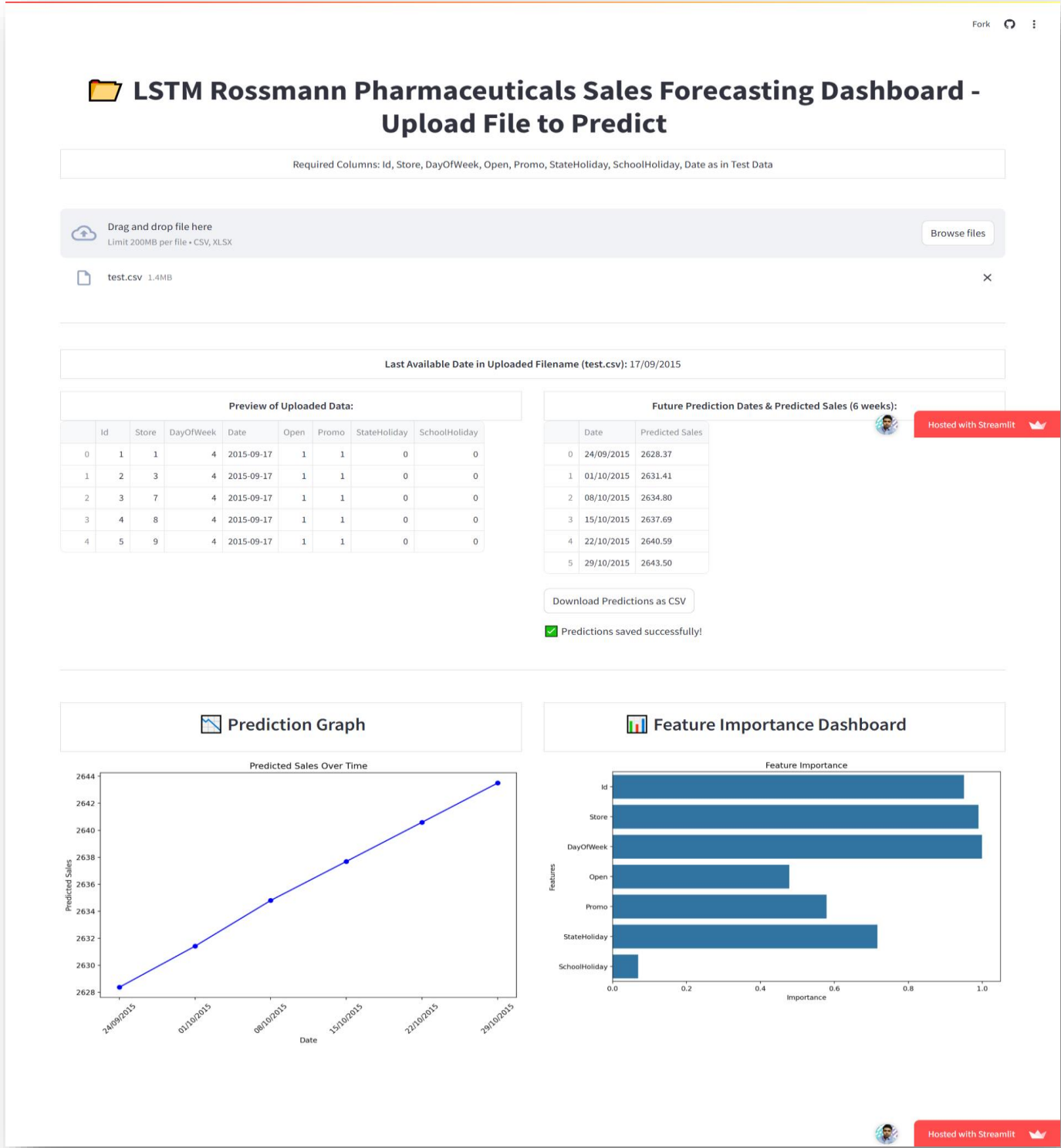
- **MAE:** 972.64
- **MSE:** 1,596,162.20
- **RMSE:** 1,263.39
- **R² Score:** 75.90%



KEY TAKEAWAYS

- Model explains 75.90% of variance; reasonable error margins.
- Ready for future sales forecasting. 🚀

Interactive Sales Forecasting with Streamlit and LSTM



1

Dashboard Overview

Developed & Deployed an interactive dashboard for forecasting sales using LSTM neural networks & Streamlit.

2

Technologies Used

- ❑ **Streamlit:** Web app framework for user interface.
- ❑ **Pandas:** Data manipulation and analysis.
- ❑ **NumPy:** Numerical operations.
- ❑ **TensorFlow/Keras:** LSTM model handling.
- ❑ **Scikit-learn:** Data preprocessing (MinMaxScaler).
- ❑ **Matplotlib & Seaborn:** Data visualization.
- ❑ **MLflow:** Model tracking and management.

3

Key Features

- ❑ **Last Most Recent Date:** Displays initial rows for user verification.
- ❑ **Uploaded Data Preview:** Shows the latest date in the dataset.
- ❑ **Predictions Download as csv:** Option to download predictions as CSV.
- ❑ **Prediction Graph:** Visualizes predicted sales trends.
- ❑ **Feature Importance Dashboard:** Displays feature importance.

4

User Interface

- ❑ Clear instructions for users regarding data requirements.
- ❑ Intuitive layout with separate sections for data preview, predictions, and visualizations.

Thank You

I appreciate your time and attention. 😊

I hope this presentation was informative and helpful.

Please don't hesitate to contact me with any questions or feedback you may have.

I am more than eager to hear your thoughts and suggestions. 😊