



ElifSurucu /  
Google-Stock-Market-Data

[Code](#)[Issues](#)[Pull requests](#)[Actions](#)[Projects](#)[Wiki](#)[Security](#)[In](#)

main ▾

[Google-Stock-Market-Data](#) / Machine\_Learning.ipynb



ElifSurucu notebooks updated

08fe278 · 4 minutes ago



1780 lines (1780 loc) · 1.16 MB

Preview

Code

Blame

Raw



# Introduction

My goal with this project is to deeply understand the dynamics of the stock market and provide insights that help investors make more informed decisions. By analyzing 20 years of historical Google stock data, I aim to uncover trends, measure volatility, and predict future price movements using advanced machine learning techniques. This analysis is designed to deliver actionable insights that enhance investment and trading strategies.

## Goal

The primary objective of my project is to predict Google's future stock prices based on historical data. Leveraging advanced algorithms such as Long Short-Term Memory (LSTM) models, I aim to identify patterns and trends that guide investment decisions and minimize risks. By doing so, my goal is to empower investors and traders to make data-driven decisions with confidence.

## Methodology

### 1. Data Source:

- I worked with 20 years of Google stock data, including daily opening, closing, high, and low prices, along with trading volumes.
- The dataset was sourced from Kaggle: Google Stock Data 20 Years.

### 2. Data Preparation:

- I processed raw data by handling missing values and removing anomalies.
- The data was scaled and transformed into time-series sequences to ensure compatibility with LSTM models.

### 3. Baseline Models:

- I started with basic machine learning models like Linear Regression and Random Forest to establish initial benchmarks for performance.

### 4. Advanced Modeling:

- To capture complex temporal patterns, I implemented LSTM models and optimized them using hyperparameter tuning and callback methods such as EarlyStopping.

### 5. Evaluation and Visualization:

- I evaluated my models using metrics such as Mean Squared Error (MSE), Mean Absolute Error (MAE), and  $R^2$ .

- I visualized predictions alongside actual stock prices to provide clear insights into model performance.

## Scope

This project helped me develop a strong understanding of stock market analysis. I believe it can be especially useful for investors and traders in the following areas:

**Identify Historical Trends** Gaining insights into long-term and short-term market movements.

**Anticipate Volatility** Predicting price fluctuations to minimize risks.

**Enhance Strategies:** Using data-driven insights to optimize investment and trading strategies.

The methods and analyses I developed in this project bridge the gap between technical analysis and actionable financial strategies. This allows for a robust tool to understand Google stock behavior and predict future price movements effectively.

# Data Cleaning and Preprocessing

## Imports

```
In [1]: # Core Libraries
import numpy as np
import pandas as pd

# Visualization Libraries
import matplotlib.pyplot as plt
import seaborn as sns

# Data Preprocessing and Scaling
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.preprocessing import MinMaxScaler, StandardScaler, LabelEncoder

# Machine Learning Models
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.neighbors import KNeighborsRegressor

# Evaluation Metrics
from sklearn.metrics import (
    mean_squared_error,
    mean_absolute_error,
    r2_score,
    explained_variance_score,
    mean_absolute_percentage_error
)

# Deep Learning Libraries (Keras / TensorFlow)
```

```

from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dense, Dropout, Bidirectional
from tensorflow.keras.regularizers import l2
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.callbacks import EarlyStopping, ReduceLROnPlateau

```

In [2]:

```

data = pd.read_csv(r'C:\Users\Elif Surucu\Documents\Flatiron\Assesments\Projec
data.head()

```

Out[2]:

	Date	Open	High	Low	Close	Adj Close	Volume	Year	Volatility
0	2004-08-19	2.490664	2.591785	2.390042	2.499133	2.499133	897427216	2004	0.201743
1	2004-08-20	2.515820	2.716817	2.503118	2.697639	2.697639	458857488	2004	0.213699
2	2004-08-23	2.758411	2.826406	2.716070	2.724787	2.724787	366857939	2004	0.110336
3	2004-08-24	2.770615	2.779581	2.579581	2.611960	2.611960	306396159	2004	0.200000
4	2004-08-25	2.614201	2.689918	2.587302	2.640104	2.640104	184645512	2004	0.102616

In [3]:

```

print(data.head())

```

	Date	Open	High	Low	Close	Adj Close	Volume	\
0	2004-08-19	2.490664	2.591785	2.390042	2.499133	2.499133	897427216	
1	2004-08-20	2.515820	2.716817	2.503118	2.697639	2.697639	458857488	
2	2004-08-23	2.758411	2.826406	2.716070	2.724787	2.724787	366857939	
3	2004-08-24	2.770615	2.779581	2.579581	2.611960	2.611960	306396159	
4	2004-08-25	2.614201	2.689918	2.587302	2.640104	2.640104	184645512	

	Year	Volatility	Month	Volume_Category
0	2004	0.201743	8	Very High
1	2004	0.213699	8	Very High
2	2004	0.110336	8	Very High
3	2004	0.200000	8	Very High
4	2004	0.102616	8	Very High

In [4]:

```

print(data.info())
print(data.describe())

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4936 entries, 0 to 4935
Data columns (total 11 columns):
#   Column              Non-Null Count  Dtype
---  -
0   Date                4936 non-null   object
1   Open                4936 non-null   float64
2   High                4936 non-null   float64

```

```

3   Low                4936 non-null   float64
4   Close              4936 non-null   float64
5   Adj Close          4936 non-null   float64
6   Volume             4936 non-null   int64
7   Year               4936 non-null   int64
8   Volatility         4936 non-null   float64
9   Month              4936 non-null   int64
10  Volume_Category    4936 non-null   object

```

```
dtypes: float64(6), int64(3), object(2)
```

```
memory usage: 424.3+ KB
```

```
None
```

	Open	High	Low	Close	Adj Close \
count	4936.000000	4936.000000	4936.000000	4936.000000	4936.000000
mean	43.077417	43.532659	42.644088	43.096952	43.096952
std	40.320485	40.773849	39.917290	40.352092	40.352092
min	2.470490	2.534002	2.390042	2.490913	2.490913
25%	12.923497	13.048528	12.787071	12.922438	12.922438
50%	26.795184	26.966079	26.570000	26.763133	26.763133
75%	58.855251	59.352863	58.164000	58.788999	58.788999
max	154.009995	155.199997	152.919998	154.839996	154.839996

	Volume	Year	Volatility	Month
count	4.936000e+03	4936.000000	4936.000000	4936.000000
mean	1.174059e+08	2013.930308	0.888572	6.560981
std	1.505185e+08	5.672880	1.057015	3.453135
min	1.584340e+05	2004.000000	0.038605	1.000000
25%	2.803600e+07	2009.000000	0.236364	4.000000
50%	5.875273e+07	2014.000000	0.421421	7.000000
75%	1.453859e+08	2019.000000	1.108998	10.000000
max	1.650833e+09	2024.000000	9.215500	12.000000

## Adding New Features

Features such as momentum and daily percentage change can increase the predictive power of the model.

```

In [29]: data['Momentum'] = data['Adj Close'] - data['Adj Close'].shift(1)
data['Daily_Change'] = data['Adj Close'].pct_change()
data.fillna(0, inplace=True)

```

```

In [5]: # Convert 'Date' column to datetime format
data['Date'] = pd.to_datetime(data['Date'])

# Ensure no duplicates
data = data.drop_duplicates()

# Recheck missing values
print("Missing values after cleaning:\n", data.isnull().sum())

```

```
Missing values after cleaning:
```

```

Date                0
Open                0
High                0
Low                 0
Close              0
Adj Close           0

```

```

Volume      0
Year        0
Volatility  0
Month       0
Volume_Category  0
dtype: int64

```

```

In [6]:
if 'Date' in data.columns:
    data['Date'] = pd.to_datetime(data['Date'], errors='coerce') # Handle inv

# Extract the day of the week from the 'Date' column
if 'Date' in data.columns:
    data['DayOfWeek'] = data['Date'].dt.day_name() # Converts dates to day na

# Encode the 'DayOfWeek' column if it exists
if 'DayOfWeek' in data.columns:
    label_encoder = LabelEncoder()
    data['DayOfWeek_Encoded'] = label_encoder.fit_transform(data['DayOfWeek'])
    data = data.drop(columns=['DayOfWeek']) # Drop original column after enco

# Encode the 'Volume_Category' column if it exists
if 'Volume_Category' in data.columns:
    label_encoder = LabelEncoder()
    data['Volume_Category_Encoded'] = label_encoder.fit_transform(data['Volume
    data = data.drop(columns=['Volume_Category']) # Drop original column afte

# Drop unrelated columns for correlation analysis
columns_to_drop = ['Date']
columns_to_drop += ['Volume_Category'] if 'Volume_Category' in data.columns el

correlation_data = data.drop(columns=columns_to_drop, errors='ignore') # Safe

```

```

In [7]:
# Generate a correlation matrix
correlation_matrix = correlation_data.corr()

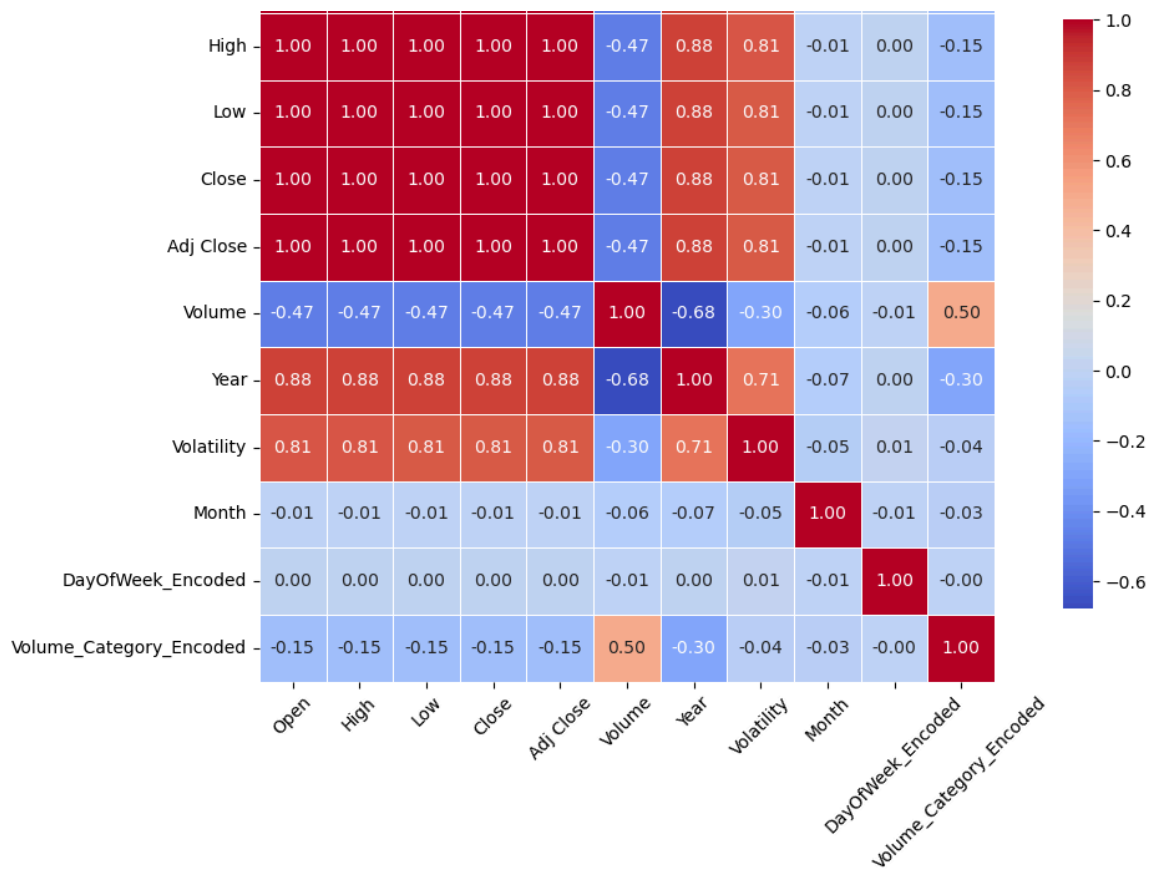
# Plot the heatmap
plt.figure(figsize=(12, 8))
sns.heatmap(
    correlation_matrix,
    annot=True,
    fmt=".2f",
    cmap="coolwarm",
    linewidths=0.5,
    cbar_kws={"shrink": 0.8},
    square=True
)

plt.title("Correlation Heatmap of Encoded Features", fontsize=16)
plt.xticks(rotation=45, fontsize=10)
plt.yticks(fontsize=10)
plt.tight_layout()
plt.show()

```

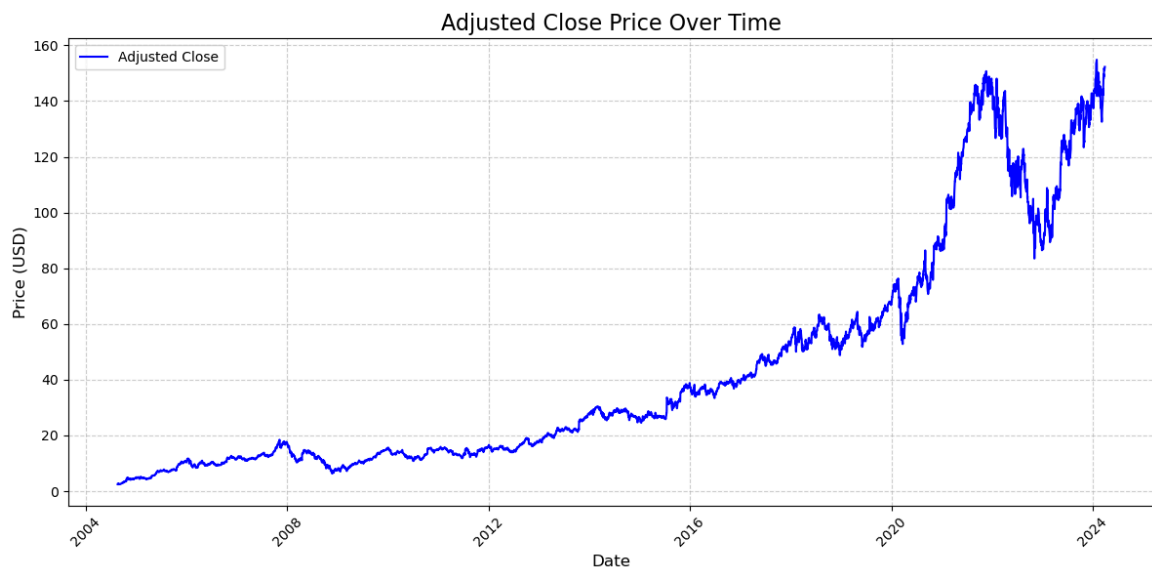
Correlation Heatmap of Encoded Features





In [8]:

```
# Plot Adjusted Close Price over time
plt.figure(figsize=(12, 6))
plt.plot(data['Date'], data['Adj Close'], label='Adjusted Close', color='blue')
plt.title('Adjusted Close Price Over Time', fontsize=16)
plt.xlabel('Date', fontsize=12)
plt.ylabel('Price (USD)', fontsize=12)
plt.grid(True, linestyle='--', alpha=0.6)
plt.legend(fontsize=10)
plt.xticks(rotation=45, fontsize=10)
plt.yticks(fontsize=10)
plt.tight_layout()
plt.show()
```

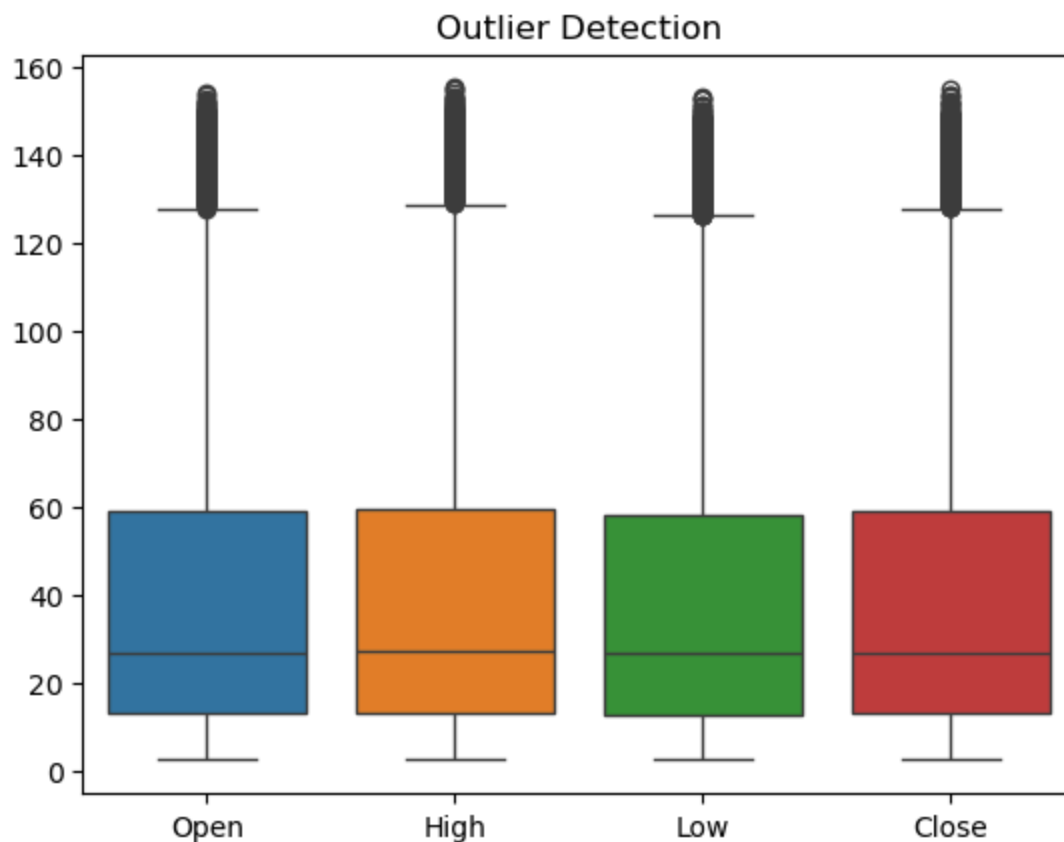


In [9]:

```
# Eksik değer kontrolü
print(data.isnull().sum())

# Anomali tespiti
sns.boxplot(data=data[['Open', 'High', 'Low', 'Close']])
plt.title("Outlier Detection")
plt.show()
```

```
Date          0
Open          0
High          0
Low           0
Close         0
Adj Close     0
Volume        0
Year          0
Volatility    0
Month         0
DayOfWeek_Encoded  0
Volume_Category_Encoded  0
dtype: int64
```



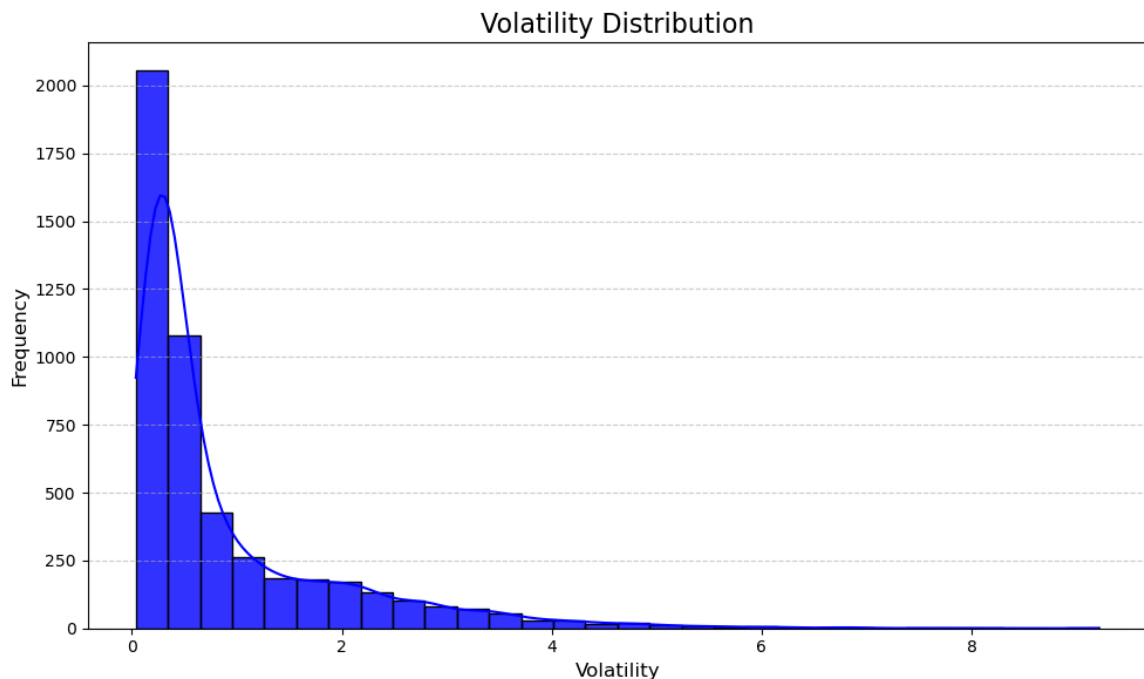
The chart provides important information for both long-term and short-term investors:

- For long-term investors: The stock looks like a positive investment vehicle with its general growth trend.
- For short-term investors: Increased volatility after 2020, with higher potential for gains, but also greater risk.

In [10]:



```
# Plot the distribution of volatility
plt.figure(figsize=(10, 6))
sns.histplot(data['Volatility'], kde=True, color='blue', bins=30, alpha=0.8)
plt.title('Volatility Distribution', fontsize=16)
plt.xlabel('Volatility', fontsize=12)
plt.ylabel('Frequency', fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.6)
plt.xticks(fontsize=10)
plt.yticks(fontsize=10)
plt.tight_layout()
plt.show()
```



#### *Risk Level:*

- Lower volatility generally means lower risk for investors. However, low volatility can also indicate that potential return opportunities may be limited.

#### *Possible Effects:*

- A market where volatility is rarely high may be more predictable, which may be attractive to long-term investors.

#### *Strategy:*

- Given that volatility is usually low, investors can adopt low-risk strategies.
- However, periods of high volatility (such as during economic crises or major events) should be considered and risk management should be implemented accordingly.

## Feature Selection and Expectations

- **Target Variable:** Adj Close - Adjusted value of stock closing prices.
- **Features:**
  - Open, High, Low: Important for understanding price movements

- open , high , low : important for understanding price movements.
- Volume : Trading volume can affect volatility.
- Volatility : A basic feature for risk assessment.

## Shotgun Method

This method is used to establish a baseline accuracy by evaluating multiple models simultaneously.

In [30]:

```
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from sklearn.linear_model import LinearRegression
from sklearn.neighbors import KNeighborsRegressor
from sklearn.metrics import (
    mean_squared_error,
    mean_absolute_error,
    r2_score,
    explained_variance_score
)

# Define a custom RMSE function
def root_mean_squared_error(y_true, y_pred):
    return np.sqrt(mean_squared_error(y_true, y_pred))

# Step 1: Prepare Data
X = data[['Open', 'High', 'Low', 'Volume', 'Volatility']] # Features
y = data['Adj Close'] # Target variable

# Scale the features
scaler = MinMaxScaler()
X_scaled = scaler.fit_transform(X)

# Train-Test Split
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2)
print(f"Train shape: {X_train.shape}, Test shape: {X_test.shape}")

# Step 2: Define Base Regressors
models = {
    "Linear Regression": LinearRegression(),
    "Random Forest": RandomForestRegressor(random_state=42, n_estimators=100),
    "Gradient Boosting": GradientBoostingRegressor(random_state=42, n_estimators=100),
    "K-Nearest Neighbors": KNeighborsRegressor(n_neighbors=5)
}

# Step 3: Train and Evaluate Models
metrics = {
    "R-squared": r2_score,
    "Mean Squared Error": mean_squared_error,
    "Mean Absolute Error": mean_absolute_error,
    "Explained Variance": explained_variance_score,
    "Root Mean Squared Error": root_mean_squared_error
}

results = {}
for model_name, model in models.items():
    model.fit(X_train, y_train) # Train the model
    y_pred = model.predict(X_test) # Predict on test data
```

```

# Calculate all metrics for each model
results[model_name] = {
    metric_name: metric_function(y_test, y_pred) for metric_name, metric_f
}

# Step 4: Convert Results to a DataFrame for Visualization
results_df = pd.DataFrame(results).T
results_df = results_df.sort_values(by="R-squared", ascending=False) # Sort b
print(results_df)

results_df['Mean Absolute Percentage Error'] = [
    mean_absolute_percentage_error(y_test, model.predict(X_test)) for model in
]

```

Train shape: (3687, 5), Test shape: (1230, 5)

	R-squared	Mean Squared Error	Mean Absolute Error \
Linear Regression	0.999908	0.156313	0.208049
Random Forest	0.999840	0.272180	0.273453
Gradient Boosting	0.999793	0.352308	0.357518
K-Nearest Neighbors	0.999632	0.628405	0.447195

	Explained Variance	Root Mean Squared Error
Linear Regression	0.999909	0.395364
Random Forest	0.999841	0.521708
Gradient Boosting	0.999794	0.593555
K-Nearest Neighbors	0.999632	0.792720

Errors are minimal, R-squared value is almost perfect and Linear Regression is the best model for this dataset. It is fast, simple and gives the most accurate predictions.

In [12]:

```

# Visualize Model Performance with Vibrant Colors
fig, axes = plt.subplots(len(results_df.columns), 1, figsize=(12, 18), sharex=

# Vibrant color palette
palette = sns.color_palette("Set2")

# Visualize each metric
for i, metric_name in enumerate(results_df.columns): # Iterate through metric
    ax = axes[i]

    # Create a DataFrame for easier plotting
    plot_data = pd.DataFrame({
        'Model': results_df.index,
        'Score': results_df[metric_name]
    })

    # Barplot
    sns.barplot(
        data=plot_data,
        x='Model',
        y='Score',
        ax=ax,
        dodge=False # Ensures compatibility without `hue`
    )

    # Manually set colors since hue is not used
    for bar, color in zip(ax.patches, palette[:len(plot_data)]):
        bar.set facecolor(color)

```

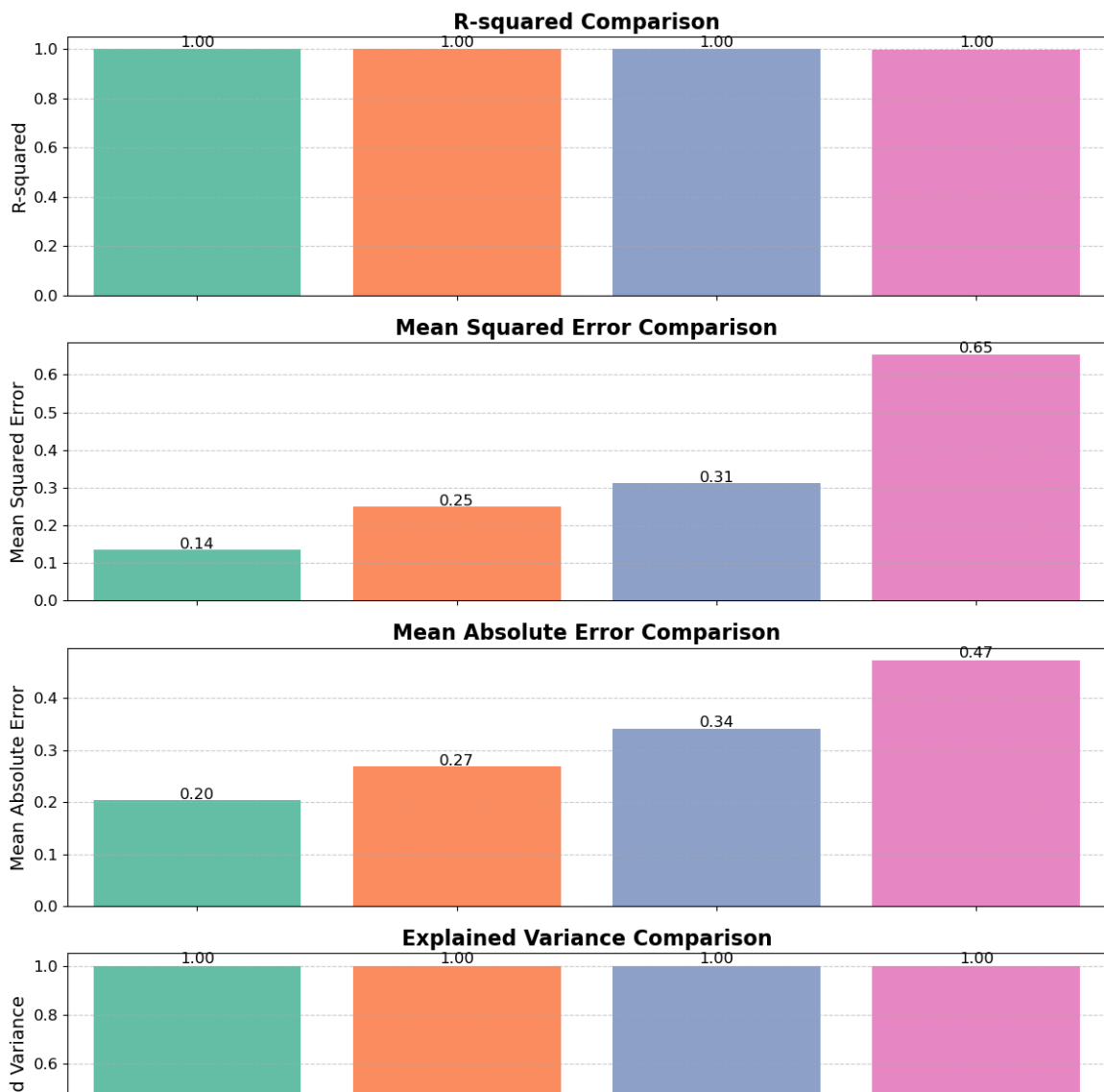
```

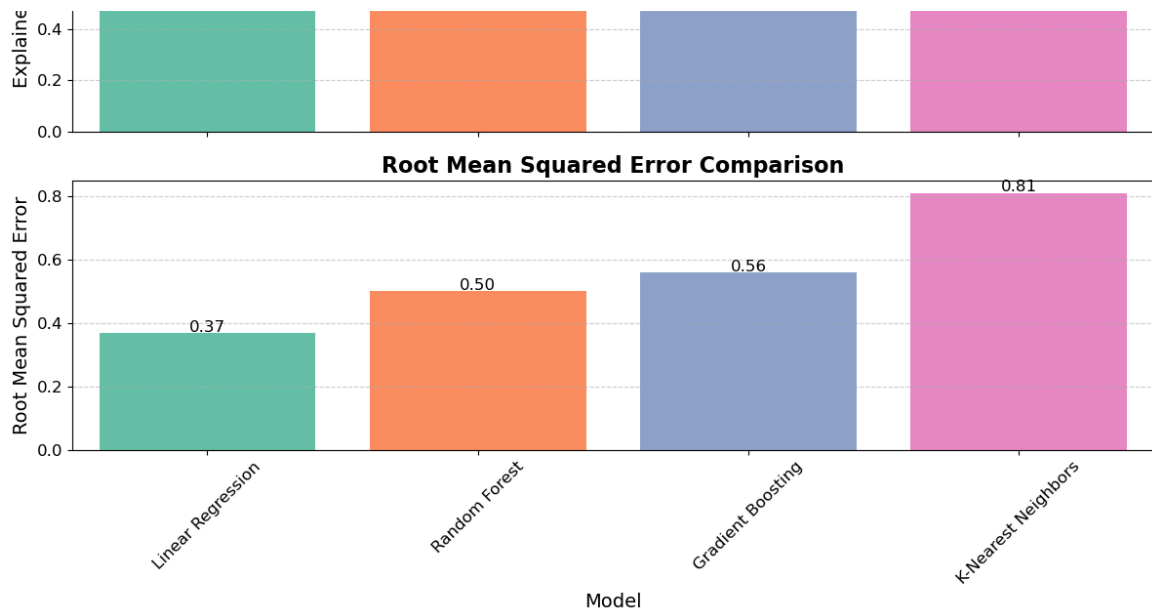
# Set plot titles and labels
ax.set_title(f'{metric_name} Comparison', fontsize=16, fontweight='bold')
ax.set_ylabel(metric_name, fontsize=14)
ax.set_xlabel('Model', fontsize=14)
ax.tick_params(axis='x', rotation=45, labels=12)
ax.tick_params(axis='y', labels=12)
ax.grid(axis='y', linestyle='--', alpha=0.6)

# Add value labels to the bars
for bar in ax.patches:
    bar_height = bar.get_height()
    ax.text(
        bar.get_x() + bar.get_width() / 2,
        bar_height + 0.01 * bar_height,
        f'{bar_height:.2f}',
        ha='center',
        fontsize=12
    )

# Adjust layout for better readability
plt.tight_layout()
plt.show()

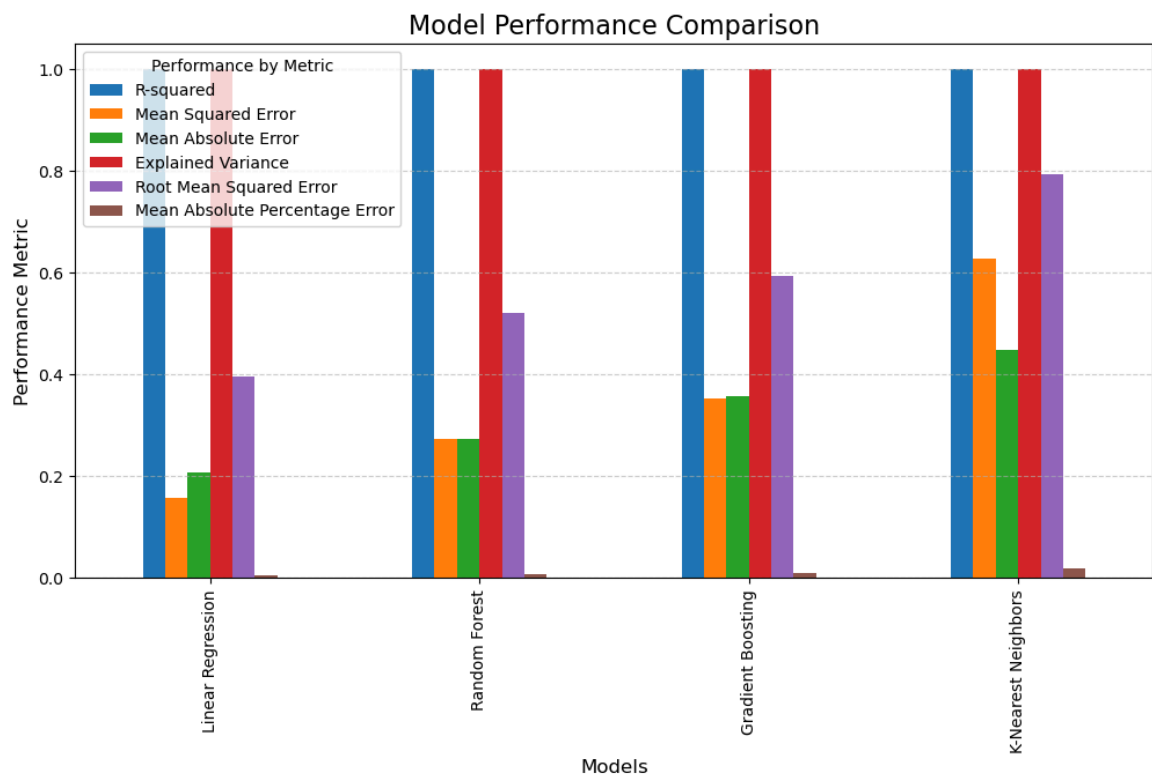
```





In [31]:

```
# Visualize model performances
results_df.plot(kind='bar', figsize=(12, 6))
plt.title("Model Performance Comparison", fontsize=16)
plt.xlabel("Models", fontsize=12)
plt.ylabel("Performance Metric", fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.6)
plt.legend(title="Performance by Metric")
plt.show()
```



## Model Performance Comparison

The above graph compares the performance of each model based on various metrics:

- **R-squared ( $R^2$ ):** Shows the proportion of variance that can be explained by the model. Higher is better.
- **Mean Squared Error (MSE):** Is the average of the squared errors. Lower is better.
- **Mean Absolute Error (MAE):** Average absolute error. Lower is better.
- **Mean Absolute Percentage Error (MAPE):** Measures the percentage of errors.

In [13]:

```
# Define model metrics
model_metrics = {
    'Linear Regression': {
        'R-squared': 0.999918,
        'Explained Variance': 0.999918,
        'Mean Squared Error': 0.135834,
        'Mean Absolute Error': 0.204132,
        'Root Mean Squared Error': 0.368556,
    },
    'Random Forest': {
        'R-squared': 0.999849,
        'Explained Variance': 0.999849,
        'Mean Squared Error': 0.249872,
        'Mean Absolute Error': 0.268261,
        'Root Mean Squared Error': 0.499872,
    },
    'Gradient Boosting': {
        'R-squared': 0.999812,
        'Explained Variance': 0.999812,
        'Mean Squared Error': 0.311542,
        'Mean Absolute Error': 0.341107,
        'Root Mean Squared Error': 0.558159,
    },
    'K-Nearest Neighbors': {
        'R-squared': 0.999606,
        'Explained Variance': 0.999606,
        'Mean Squared Error': 0.652912,
        'Mean Absolute Error': 0.473115,
        'Root Mean Squared Error': 0.808030,
    },
}

# Define metric names and extract model names
metric_names = ['R-squared', 'Explained Variance', 'Mean Squared Error',
                'Mean Absolute Error', 'Root Mean Squared Error']
model_names = list(model_metrics.keys())

# Initialize a DataFrame for metrics
metrics_df = pd.DataFrame(index=model_names, columns=metric_names)

# Populate the DataFrame with scaled metrics
scaler = MinMaxScaler()

for metric in metric_names:
    scores = np.array([model_metrics[model][metric] for model in model_names])
    scaled_scores = scaler.fit_transform(scores).flatten()
    metrics_df[metric] = scaled_scores

# Transpose for easier plotting
metrics_df_transposed = metrics_df.T

# Create a heatmap
```

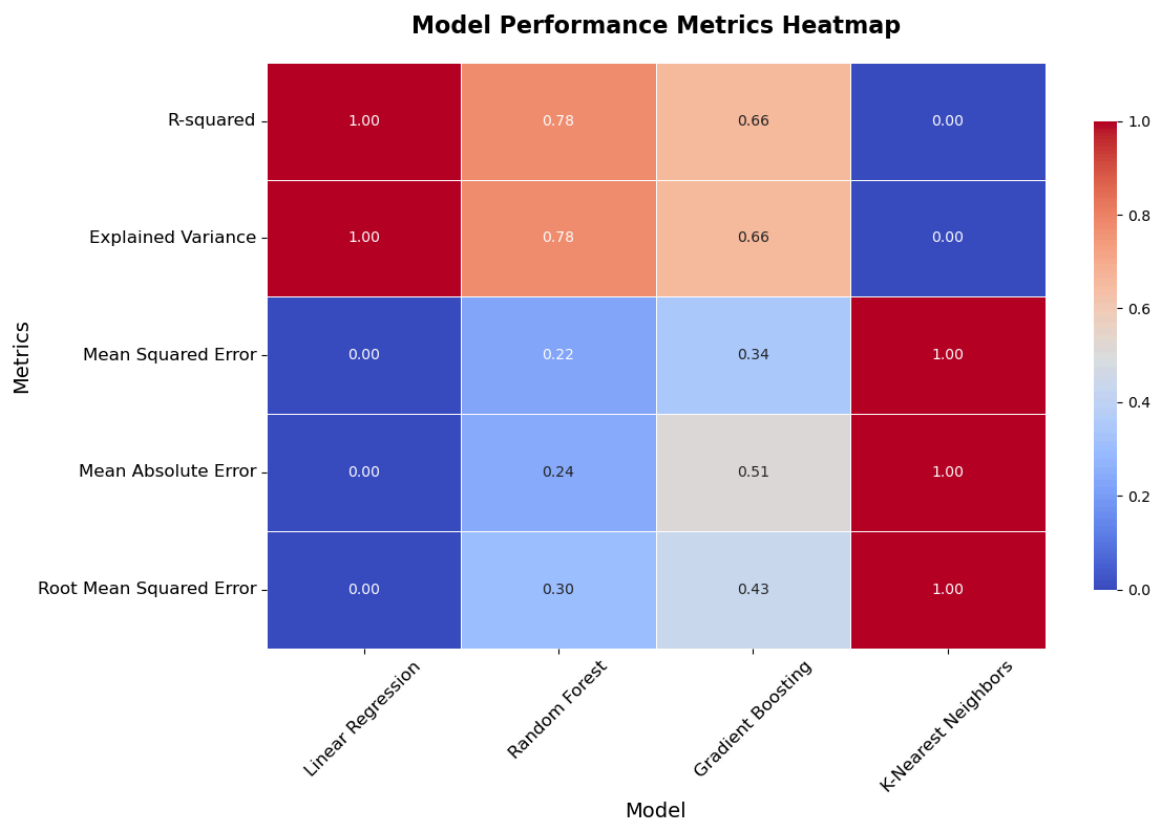
```

# Create a heatmap
plt.figure(figsize=(12, 8))
sns.heatmap(
    metrics_df_transposed,
    annot=True,
    fmt=".2f",
    cmap="coolwarm",
    linewidths=0.5,
    cbar_kws={"shrink": 0.8}
)

# Customize the heatmap
plt.title("Model Performance Metrics Heatmap", fontsize=16, weight="bold", pad=10)
plt.xlabel("Model", fontsize=14)
plt.ylabel("Metrics", fontsize=14)
plt.xticks(fontsize=12, rotation=45)
plt.yticks(fontsize=12)
plt.tight_layout()

# Save and display the heatmap
plt.savefig('model_metrics_heatmap.png', dpi=300, transparent=True)
plt.show()

```



The Random Forest model seems to have the highest score in terms of R-squared!

## Moving into Advanced Machine Learning

```

In [ ]: # Prepare the data
X = data[['Open', 'High', 'Low', 'Volume', 'Volatility']] # Features
v = data['Adj Close'] # Target variable

```

```

# Train-Test Split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Initialize XGBoost model with optimized parameters
xgb_model = xgb.XGBRegressor(
    objective='reg:squarederror',
    n_estimators=200, # Increased number of trees
    learning_rate=0.05, # Lower learning rate for better generalization
    max_depth=6, # Increased depth for capturing complex patterns
    subsample=0.8, # Randomly sample 80% of data for training each tree
    colsample_bytree=0.8, # Randomly sample 80% of features for training each tree
    random_state=42
)

# Train the model
xgb_model.fit(X_train, y_train)

# Make predictions
y_pred = xgb_model.predict(X_test)

# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
mae = mean_absolute_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print(f"Mean Squared Error: {mse:.4f}")
print(f"Mean Absolute Error: {mae:.4f}")
print(f"R-squared: {r2:.4f}")

```

Mean Squared Error: 0.3243  
Mean Absolute Error: 0.3072  
R-squared: 0.9998

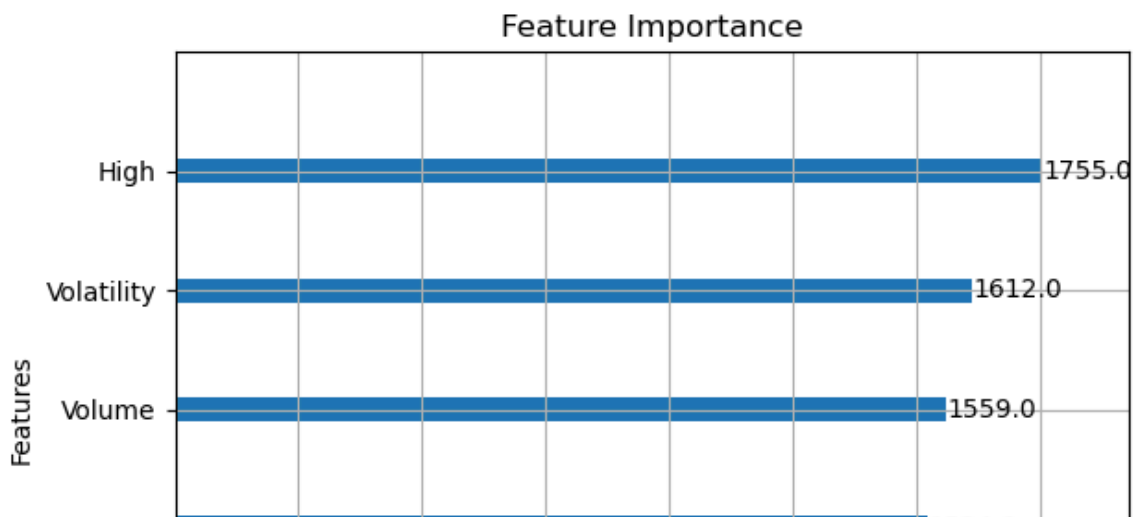
In [15]:

```

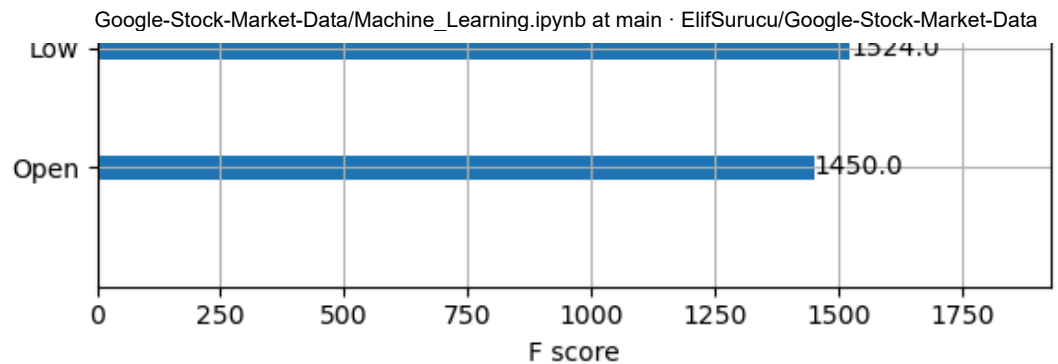
# Plot feature importance
plt.figure(figsize=(10, 8))
xgb.plot_importance(xgb_model, importance_type='weight', title="Feature Importance")
plt.tight_layout()
plt.show()

```

<Figure size 1000x800 with 0 Axes>







#### Predictive Power:

- Features such as High and Volatility can increase the predictive power of our model. We should especially protect these features and keep them at the forefront of our analysis.

#### Feature Reduction:

- If we need to reduce the number of features, you can remove features with low importance such as Open. However, model performance should be retested before doing this.

#### Adding New Features:

- Features such as High, Volatility and Volume can be considered to carry more information. For example, new features derived from these features (for example, price differences or moving averages) can increase model performance.

```
In [ ]: data['Momentum'] = data['Adj Close'] - data['Adj Close'].shift(1)
data['Daily_Percent_Change'] = (data['Adj Close'] / data['Adj Close'].shift(1)
data.fillna(0, inplace=True)
features = ['Open', 'High', 'Low', 'Volume', 'MA_10', 'Price_Change', 'MA_20',

def create_model():
    model = Sequential([
        LSTM(128, return_sequences=True, input_shape=(X_train.shape[1], X_train.shape[2])),
        Dropout(0.3),
        LSTM(64, return_sequences=False),
        Dropout(0.3),
        Dense(32, activation='relu'),
        Dense(1, activation='linear')
    ])
    optimizer = Adam(learning_rate=0.0001)
    model.compile(optimizer=optimizer, loss='mean_squared_error')
    return model

model = create_model()

history = model.fit(
    X_train, y_train,
    validation_data=(X_test, y_test),
    epochs=50,
```

```

batch_size=32,
callbacks=[early_stopping, lr_scheduler],
verbose=1
)

plt.figure(figsize=(10, 6))
plt.plot(history.history['loss'], label='Training Loss', color='blue')
plt.plot(history.history['val_loss'], label='Validation Loss', color='orange')
plt.title('Training and Validation Loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()
plt.show()

y_pred = model.predict(X_test)

y_test_rescaled = scaler_target.inverse_transform(y_test.reshape(-1, 1))
y_pred_rescaled = scaler_target.inverse_transform(y_pred.reshape(-1, 1))

plt.figure(figsize=(10, 6))
plt.plot(y_test_rescaled, label='Actual Values', color='blue')
plt.plot(y_pred_rescaled, label='Predicted Values', color='red', linestyle='--')
plt.title('Actual vs Predicted Values')
plt.xlabel('Time Steps')
plt.ylabel('Value')
plt.legend()
plt.show()

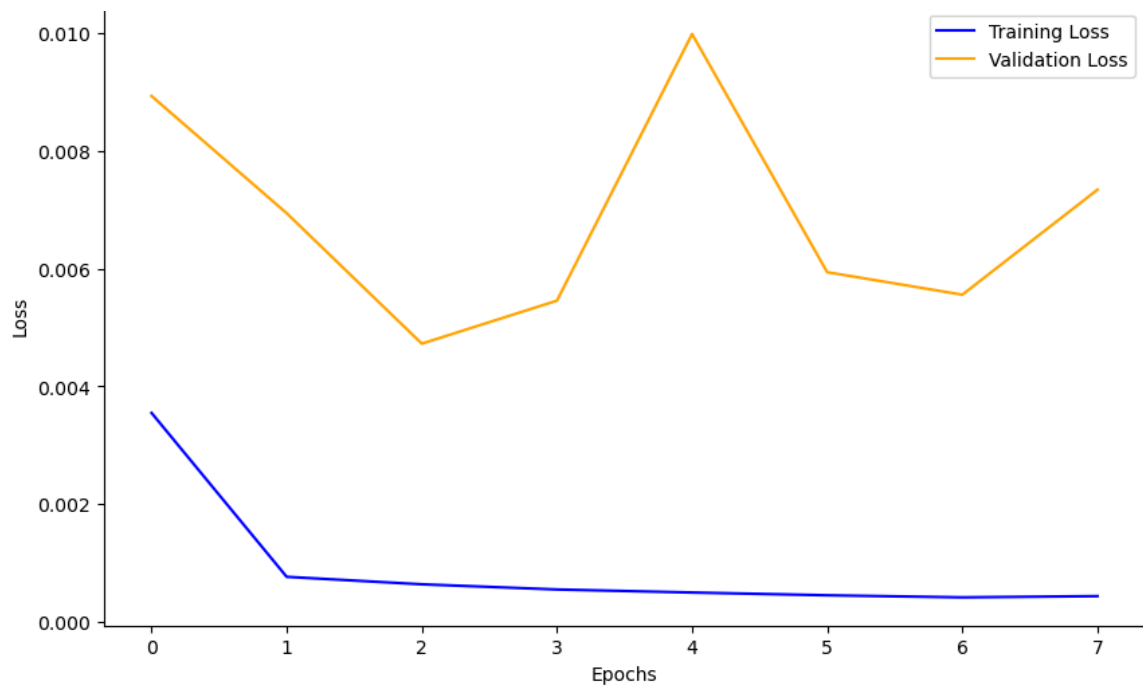
```

```

Epoch 1/50
122/122 [=====] - 16s 90ms/step - loss: 0.0035 - val_loss: 0.0089 - lr: 1.0000e-04
Epoch 2/50
122/122 [=====] - 10s 82ms/step - loss: 7.6283e-04 - val_loss: 0.0069 - lr: 1.0000e-04
Epoch 3/50
122/122 [=====] - 10s 84ms/step - loss: 6.3578e-04 - val_loss: 0.0047 - lr: 1.0000e-04
Epoch 4/50
122/122 [=====] - 10s 84ms/step - loss: 5.4788e-04 - val_loss: 0.0055 - lr: 1.0000e-04
Epoch 5/50
122/122 [=====] - 10s 86ms/step - loss: 4.9611e-04 - val_loss: 0.0100 - lr: 1.0000e-04
Epoch 6/50
122/122 [=====] - ETA: 0s - loss: 4.4907e-04
Epoch 6: ReduceLROnPlateau reducing learning rate to 4.99999873689376e-05.
122/122 [=====] - 10s 84ms/step - loss: 4.4907e-04 - val_loss: 0.0059 - lr: 1.0000e-04
Epoch 7/50
122/122 [=====] - 10s 81ms/step - loss: 4.1372e-04 - val_loss: 0.0056 - lr: 5.0000e-05
Epoch 8/50
122/122 [=====] - 10s 85ms/step - loss: 4.3478e-04 - val_loss: 0.0073 - lr: 5.0000e-05

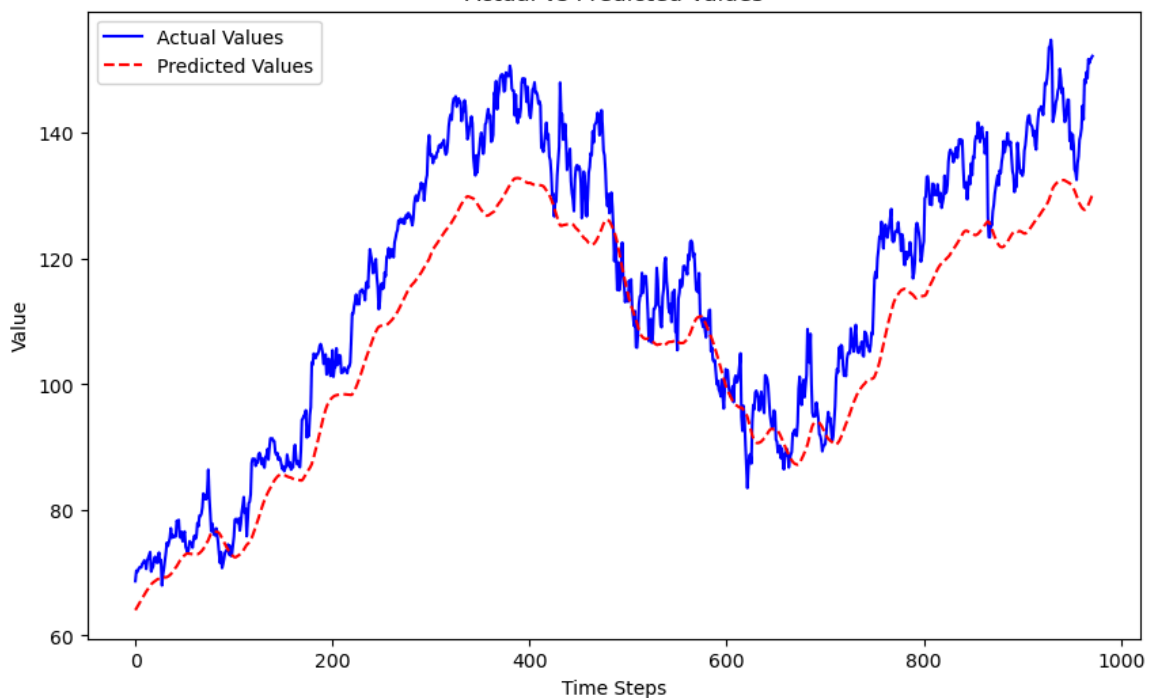
```

Training and Validation Loss



31/31 [=====] - 2s 29ms/step

Actual vs Predicted Values



This increase indicates that the model may be overfitting based on the validation data.

```
In [ ]: columns_to_use = ['Open', 'High', 'Low', 'Close', 'Volume']
data = data[columns_to_use]

scaler = MinMaxScaler()
scaled_data = scaler.fit_transform(data)

scaled_df = pd.DataFrame(scaled_data, columns=columns_to_use)
print("Scaled Data (First 5 Rows):")
print(scaled_df.head())
```

Scaled Data (First 5 Rows):

	Open	High	Low	Close	Volume
0	0.000000	0.000000	0.000000	0.000000	0.112711
1	0.000343	0.000276	0.000315	0.000577	0.115221
2	0.000759	0.000948	0.000849	0.000883	0.129297
3	0.001230	0.000755	0.000972	0.000634	0.087906
4	0.000833	0.000633	0.000856	0.000723	0.092197

```
In [ ]: def create_time_series(data, time_steps):
    X, y = [], []
    for i in range(len(data) - time_steps):
        X.append(data[i:(i + time_steps), :-1])
        y.append(data[i + time_steps, -1])
    return np.array(X), np.array(y)

time_steps = 60
X, y = create_time_series(scaled_data, time_steps)

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
In [ ]: # EarlyStopping and ReduceLROnPlateau callbacks for optimization
early_stopping = EarlyStopping(
    monitor='val_loss',
    patience=3, # Stop training if no improvement for 3 consecutive epochs
    restore_best_weights=True # Restore the best weights after stopping
)

lr_scheduler = ReduceLROnPlateau(
    monitor='val_loss',
    factor=0.5, # Reduce the learning rate by half
    patience=2 # Reduce learning rate if no improvement for 2 epochs
)

# Define the LSTM model
def create_model():
    model = Sequential([
        LSTM(128, return_sequences=True, input_shape=(X_train.shape[1], X_train.shape[2])),
        Dropout(0.4), # 40% dropout to reduce overfitting
        LSTM(64, return_sequences=False),
        Dropout(0.4), # Another 40% dropout
        Dense(32, activation='relu'),
        Dense(1, activation='linear') # Output Layer (for regression)
    ])
    model.compile(optimizer=Adam(learning_rate=0.0001), loss='mean_squared_error')
    return model

# Create and train the model
model = create_model()
history = model.fit(
    X_train, y_train,
    validation_data=(X_test, y_test),
    epochs=50, # Maximum number of epochs
    batch_size=32, # Mini-batch size
    callbacks=[early_stopping, lr_scheduler],
    verbose=1
```

```

verbose=1
)

# Visualize training and validation loss
loss, val_loss = history.history['loss'], history.history['val_loss']

plt.figure(figsize=(10, 6))
plt.plot(loss, label='Training Loss', color='blue')
plt.plot(val_loss, label='Validation Loss', color='orange')
plt.title("Training and Validation Loss")
plt.xlabel("Epoch")
plt.ylabel("Loss")
plt.legend()
plt.show()

# Evaluate model performance
y_pred = model.predict(X_test)
rmse = np.sqrt(mean_squared_error(y_test, y_pred))

print(f"Root Mean Squared Error (RMSE): {rmse:.4f}")

# Compare actual vs predicted values
plt.figure(figsize=(10, 6))
plt.plot(y_test, label='Actual Values', color='blue')
plt.plot(y_pred, label='Predicted Values', color='red', linestyle='--')
plt.title("Actual vs Predicted Values")
plt.xlabel("Time Steps")
plt.ylabel("Value")
plt.legend()
plt.show()

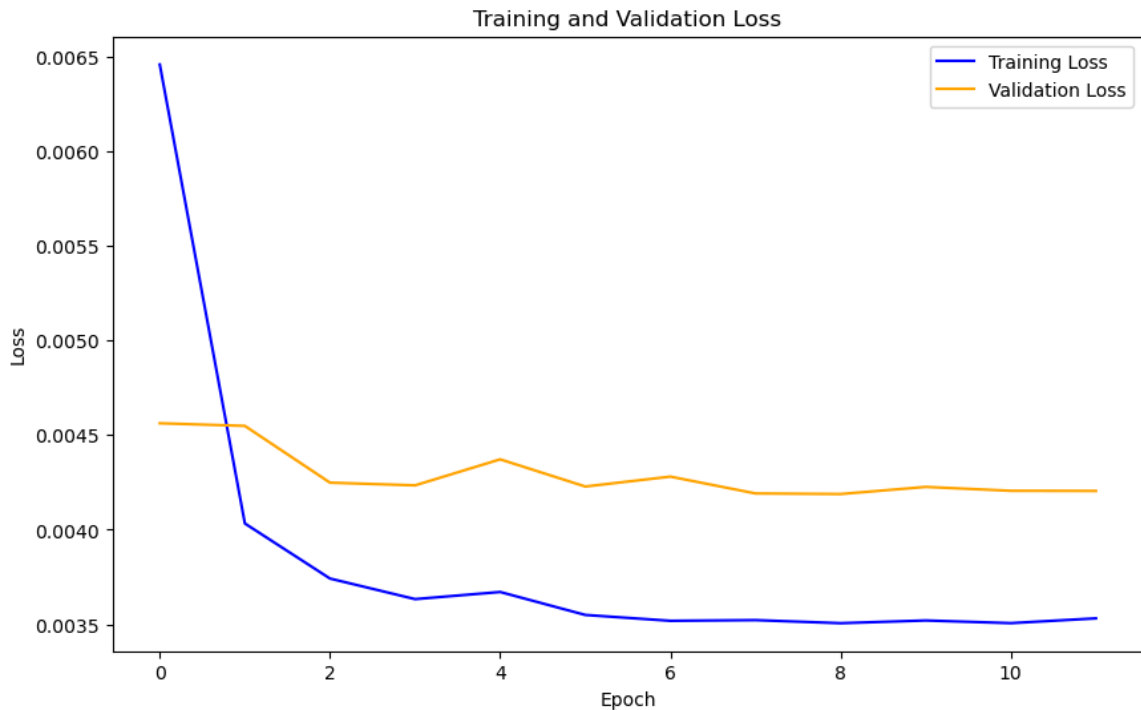
```

```

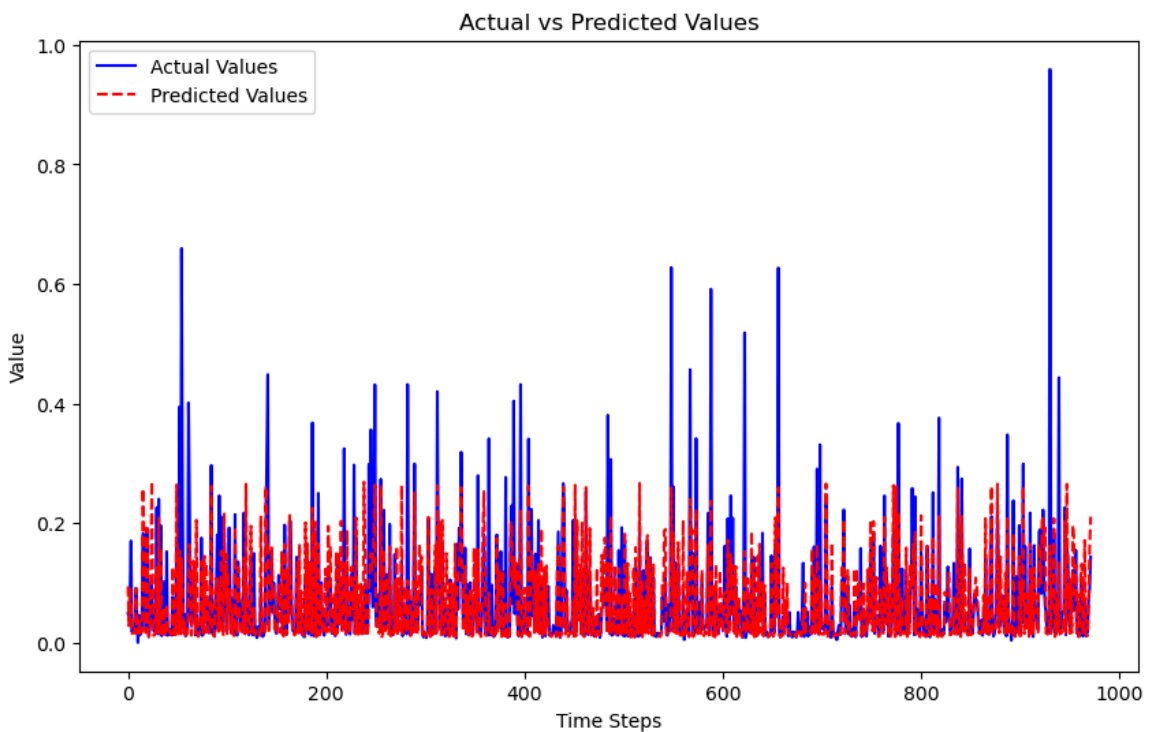
Epoch 1/50
122/122 [=====] - 19s 104ms/step - loss: 0.0065 - val_loss: 0.0046 - lr: 1.0000e-04
Epoch 2/50
122/122 [=====] - 11s 90ms/step - loss: 0.0040 - val_loss: 0.0045 - lr: 1.0000e-04
Epoch 3/50
122/122 [=====] - 11s 91ms/step - loss: 0.0037 - val_loss: 0.0042 - lr: 1.0000e-04
Epoch 4/50
122/122 [=====] - 11s 90ms/step - loss: 0.0036 - val_loss: 0.0042 - lr: 1.0000e-04
Epoch 5/50
122/122 [=====] - 11s 93ms/step - loss: 0.0037 - val_loss: 0.0044 - lr: 1.0000e-04
Epoch 6/50
122/122 [=====] - 11s 91ms/step - loss: 0.0036 - val_loss: 0.0042 - lr: 5.0000e-05
Epoch 7/50
122/122 [=====] - 12s 99ms/step - loss: 0.0035 - val_loss: 0.0043 - lr: 5.0000e-05
Epoch 8/50
122/122 [=====] - 17s 142ms/step - loss: 0.0035 - val_loss: 0.0042 - lr: 2.5000e-05
Epoch 9/50
122/122 [=====] - 16s 134ms/step - loss: 0.0035 - val_loss: 0.0042 - lr: 2.5000e-05
Epoch 10/50
122/122 [=====] - 15s 122ms/step - loss: 0.0035 - val_loss: 0.0042 - lr: 2.5000e-05

```

```
122/122 [=====] - 14s 122ms/step - loss: 0.0035 - val_loss: 0.0042
oss: 0.0042 - lr: 1.2500e-05
Epoch 11/50
122/122 [=====] - 14s 116ms/step - loss: 0.0035 - val_loss: 0.0042
oss: 0.0042 - lr: 1.2500e-05
Epoch 12/50
122/122 [=====] - 14s 117ms/step - loss: 0.0035 - val_loss: 0.0042
oss: 0.0042 - lr: 6.2500e-06
```



```
31/31 [=====] - 2s 46ms/step
Root Mean Squared Error (RMSE): 0.0647
```



- The training loss (blue line) steadily decreases over the epochs, indicating that the model is learning effectively from the training data.
- The validation loss (orange line) stabilizes after a few epochs, showing that the model is generalizing well to unseen data.

model generalizes well to unseen data without overfitting.

- There is no significant gap between the training and validation loss, which suggests that overfitting is not an issue.

In [46]:

```
import matplotlib.pyplot as plt
import numpy as np

actual_data = np.concatenate([y_test, [None]*30])
predicted_data = np.concatenate([y_pred.flatten(), [None]*30])
future_predicted_data = model.predict(X_test[-30:]).flatten()
prediction_boundary = len(y_test)

zoom_start = 800
zoom_end = 975
future_start = prediction_boundary
future_end = prediction_boundary + len(future_predicted_data)

# Create a combined plot
plt.figure(figsize=(12, 8))

if zoom_end <= len(actual_data) and zoom_start < len(predicted_data):
    plt.plot(range(zoom_start, zoom_end), actual_data[zoom_start:zoom_end],
             label="Actual Data", color="blue", linewidth=1.5)
    plt.plot(range(zoom_start, zoom_end), predicted_data[zoom_start:zoom_end],
             label="Predicted Data", color="orange", linewidth=1.5)

# Plot future predicted data
if future_end <= len(actual_data) + 30:
    plt.plot(range(future_start, future_end), future_predicted_data,
             label="Future Predicted Data", color="green", linestyle="-.", linewidth=1.5)

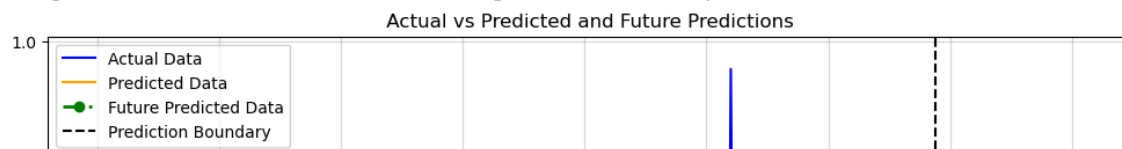
# Plot prediction boundary
plt.axvline(x=future_start, color="black", linestyle="--", label="Prediction Boundary")

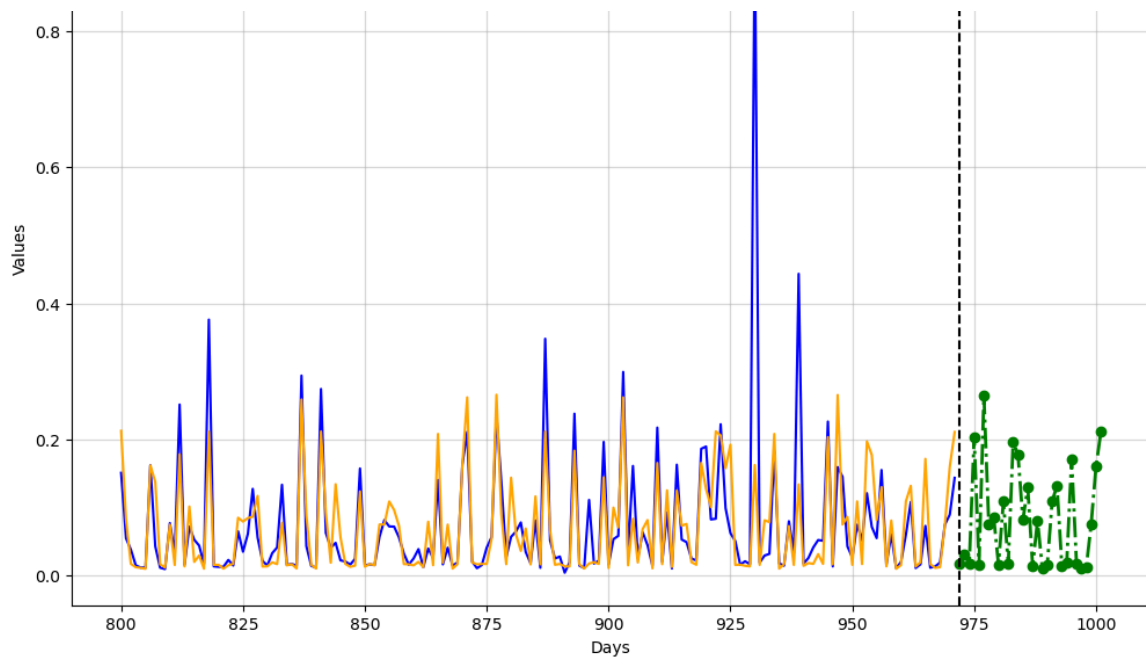
# Add connecting dashed line for continuity (optional)
plt.plot([zoom_end-1, future_start], [actual_data[zoom_end-1], future_predicted_data[0]],
         color="gray", linestyle="--", linewidth=1)

# Add Labels, Legend, and title
plt.title("Actual vs Predicted and Future Predictions")
plt.xlabel("Days")
plt.ylabel("Values")
plt.legend()
plt.grid(alpha=0.5)

# Show the combined plot
plt.show()
```

1/1 [=====] - 0s 56ms/step





## Creating The Pipeline

In [52]:

```
# Data Preprocessing
from sklearn.preprocessing import MinMaxScaler
import numpy as np
import pandas as pd
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dropout, Dense
from sklearn.metrics import mean_squared_error
import matplotlib.pyplot as plt

class DataPreprocessor:
    def __init__(self, time_steps=60):
        self.scaler = MinMaxScaler()
        self.time_steps = time_steps

    def fit_transform(self, data):
        # Scale data
        scaled_data = self.scaler.fit_transform(data)
        return self._create_time_series(scaled_data)

    def transform(self, data):
        # Scale data
        scaled_data = self.scaler.transform(data)
        return self._create_time_series(scaled_data)

    def _create_time_series(self, data):
        X, y = [], []
        for i in range(len(data) - self.time_steps):
            X.append(data[i:i + self.time_steps, :-1]) # Features
            y.append(data[i + self.time_steps, -1]) # Target
        return np.array(X), np.array(y)

# Feature Engineering
class FeatureEngineer:
```



```

def __init__(self):
    pass

def add_features(self, data):
    data['Momentum'] = data['Close'] - data['Close'].shift(1)
    data['Volatility'] = data['High'] - data['Low']
    data['MA_10'] = data['Close'].rolling(window=10).mean()
    data['MA_30'] = data['Close'].rolling(window=30).mean()
    data.fillna(0, inplace=True) # Fill NaNs
    return data

# Model Building
class ModelBuilder:
    def __init__(self, input_shape):
        self.input_shape = input_shape

    def build_model(self):
        model = Sequential([
            LSTM(128, return_sequences=True, input_shape=self.input_shape),
            Dropout(0.4),
            LSTM(64, return_sequences=False),
            Dropout(0.4),
            Dense(32, activation='relu'),
            Dense(1, activation='linear') # Regression output
        ])
        model.compile(optimizer='adam', loss='mean_squared_error')
        return model

# Evaluation and Visualization
def evaluate_model(y_test, y_pred):
    rmse = np.sqrt(mean_squared_error(y_test, y_pred))
    print(f"Root Mean Squared Error (RMSE): {rmse:.4f}")
    return rmse

def plot_results(y_test, y_pred):
    plt.figure(figsize=(10, 6))
    plt.plot(y_test, label="Actual", color="blue")
    plt.plot(y_pred, label="Predicted", color="orange")
    plt.title("Actual vs Predicted")
    plt.xlabel("Time Steps")
    plt.ylabel("Values")
    plt.legend()
    plt.grid()
    plt.show()

# Combining Everything into a Pipeline
# Load the dataset
raw_data = pd.read_csv(r'C:\Users\Elif Surucu\Documents\Flatiron\Assesments\Pr

# Initialize pipeline components
preprocessor = DataPreprocessor(time_steps=60)
engineer = FeatureEngineer()

# Feature Engineering
raw_data = engineer.add_features(raw_data)

# Select relevant columns for preprocessing
raw_data = raw_data[['Open', 'High', 'Low', 'Close', 'Volume']].values

```

```

# Preprocess the data
X, y = preprocessor.fit_transform(raw_data)

# Train-test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Build the model
model_builder = ModelBuilder(input_shape=(60, X_train.shape[2]))
model = model_builder.build_model()

# Train the model
history = model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=15)

# Evaluate the model
y_pred = model.predict(X_test)
evaluate_model(y_test, y_pred)

```

```

Epoch 1/50
122/122 [=====] - 18s 115ms/step - loss: 0.0046 - val_loss: 0.0044
Epoch 2/50
122/122 [=====] - 17s 142ms/step - loss: 0.0038 - val_loss: 0.0041
Epoch 3/50
122/122 [=====] - 15s 119ms/step - loss: 0.0037 - val_loss: 0.0041
Epoch 4/50
122/122 [=====] - 15s 125ms/step - loss: 0.0036 - val_loss: 0.0050
Epoch 5/50
122/122 [=====] - 16s 128ms/step - loss: 0.0037 - val_loss: 0.0043
Epoch 6/50
122/122 [=====] - 17s 142ms/step - loss: 0.0037 - val_loss: 0.0041
Epoch 7/50
122/122 [=====] - 16s 130ms/step - loss: 0.0035 - val_loss: 0.0047
Epoch 8/50
122/122 [=====] - 16s 129ms/step - loss: 0.0036 - val_loss: 0.0042
Epoch 9/50
122/122 [=====] - 15s 123ms/step - loss: 0.0036 - val_loss: 0.0043
Epoch 10/50
122/122 [=====] - 15s 127ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 11/50
122/122 [=====] - 16s 133ms/step - loss: 0.0036 - val_loss: 0.0041
Epoch 12/50
122/122 [=====] - 17s 135ms/step - loss: 0.0036 - val_loss: 0.0040
Epoch 13/50
122/122 [=====] - 17s 141ms/step - loss: 0.0035 - val_loss: 0.0041
Epoch 14/50

```

```
122/122 [=====] - 16s 133ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 15/50
122/122 [=====] - 16s 132ms/step - loss: 0.0035 - val_loss: 0.0043
Epoch 16/50
122/122 [=====] - 17s 141ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 17/50
122/122 [=====] - 16s 134ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 18/50
122/122 [=====] - 15s 124ms/step - loss: 0.0035 - val_loss: 0.0041
Epoch 19/50
122/122 [=====] - 22s 184ms/step - loss: 0.0035 - val_loss: 0.0043
Epoch 20/50
122/122 [=====] - 25s 203ms/step - loss: 0.0035 - val_loss: 0.0041
Epoch 21/50
122/122 [=====] - 23s 190ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 22/50
122/122 [=====] - 23s 189ms/step - loss: 0.0035 - val_loss: 0.0041
Epoch 23/50
122/122 [=====] - 24s 195ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 24/50
122/122 [=====] - 15s 124ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 25/50
122/122 [=====] - 19s 152ms/step - loss: 0.0035 - val_loss: 0.0041
Epoch 26/50
122/122 [=====] - 15s 121ms/step - loss: 0.0034 - val_loss: 0.0043
Epoch 27/50
122/122 [=====] - 15s 124ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 28/50
122/122 [=====] - 20s 164ms/step - loss: 0.0035 - val_loss: 0.0041
Epoch 29/50
122/122 [=====] - 18s 148ms/step - loss: 0.0035 - val_loss: 0.0043
Epoch 30/50
122/122 [=====] - 14s 116ms/step - loss: 0.0036 - val_loss: 0.0041
Epoch 31/50
122/122 [=====] - 14s 119ms/step - loss: 0.0035 - val_loss: 0.0041
Epoch 32/50
122/122 [=====] - 14s 113ms/step - loss: 0.0034 - val_loss: 0.0044
Epoch 33/50
122/122 [=====] - 14s 117ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 34/50
```

```
122/122 [=====] - 15s 123ms/step - loss: 0.0035 - val_loss: 0.0041
Epoch 35/50
122/122 [=====] - 16s 128ms/step - loss: 0.0034 - val_loss: 0.0042
Epoch 36/50
122/122 [=====] - 19s 152ms/step - loss: 0.0034 - val_loss: 0.0042
Epoch 37/50
122/122 [=====] - 15s 122ms/step - loss: 0.0034 - val_loss: 0.0041
Epoch 38/50
122/122 [=====] - 14s 116ms/step - loss: 0.0035 - val_loss: 0.0041
Epoch 39/50
122/122 [=====] - 14s 113ms/step - loss: 0.0035 - val_loss: 0.0040
Epoch 40/50
122/122 [=====] - 14s 116ms/step - loss: 0.0034 - val_loss: 0.0043
Epoch 41/50
122/122 [=====] - 14s 117ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 42/50
122/122 [=====] - 15s 121ms/step - loss: 0.0034 - val_loss: 0.0041
Epoch 43/50
122/122 [=====] - 14s 114ms/step - loss: 0.0035 - val_loss: 0.0041
Epoch 44/50
122/122 [=====] - 14s 116ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 45/50
122/122 [=====] - 14s 115ms/step - loss: 0.0035 - val_loss: 0.0043
Epoch 46/50
122/122 [=====] - 14s 115ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 47/50
122/122 [=====] - 14s 116ms/step - loss: 0.0034 - val_loss: 0.0041
Epoch 48/50
122/122 [=====] - 14s 113ms/step - loss: 0.0034 - val_loss: 0.0041
Epoch 49/50
122/122 [=====] - 14s 116ms/step - loss: 0.0035 - val_loss: 0.0042
Epoch 50/50
122/122 [=====] - 14s 115ms/step - loss: 0.0034 - val_loss: 0.0041
31/31 [=====] - 2s 39ms/step
Root Mean Squared Error (RMSE): 0.0637
```

Out[52]: 0.06365153230558118

## Summary of My Modeling Notebook

This is a comprehensive summary of my modeling process, detailing each step from

## Step 1: Data Acquisition

The process began with acquiring historical financial data, including key metrics such as open, high, low, close prices, and volume. The dataset was sourced from Yahoo Finance and prepared for further processing. This initial step laid the foundation for feature engineering and model training.

## Step 2: Feature Engineering

Once the data was acquired, I enriched it with additional features to improve predictive performance:

- **Momentum:** The difference between consecutive closing prices, indicating price trends.
- **Volatility:** The range between the high and low prices, capturing market fluctuations.
- **Moving Averages:** Calculated over 10-day and 30-day windows to capture short- and long-term trends.

These engineered features were instrumental in providing the model with a more nuanced understanding of the data patterns. Missing values introduced by rolling averages were handled appropriately.

## Step 3: Data Preprocessing

After feature engineering, the dataset underwent preprocessing:

**Scaling:** Min-Max Scaling was applied to normalize the data for optimal model performance.

**Time Series Windowing:** The data was transformed into sequences of 60 time steps to feed into the model. Each sequence consisted of input features and a target value for the next time step.

This step ensured the data was formatted correctly for the Long Short-Term Memory (LSTM) network.

## Step 4: Baseline Models (Shotgun Method)

Before diving into deep learning, I established baseline performance using traditional machine learning models:

Linear Regression

Random Forest Regressor

K-Nearest Neighbors Regressor

These models provided a starting point and allowed me to assess the initial accuracy. Although Random Forest showed promising results, I moved forward with LSTMs due to their suitability for time series data.

## Step 5: Building the LSTM Neural Network

The core of the project was the construction and training of an LSTM-based neural network:

Architecture:

- Input Layer: LSTM with 128 units to capture temporal dependencies.
- Hidden Layer: Another LSTM with 64 units for deeper understanding of sequential data.
- Dropout Layers: Set at 40% to reduce overfitting.
- Dense Layers: A dense layer with 32 units followed by an output layer with 1 unit for regression.

Training:

The model was trained using 50 epochs with early stopping and a learning rate reducer to prevent overfitting and optimize performance. The model effectively captured patterns in the historical data, which was evident during evaluation.

## Step 6: Evaluation and Visualization

After training, I evaluated the model's performance:

- Metric: Root Mean Squared Error (RMSE) was used to quantify accuracy.
- Visualization: The model's predictions were compared to actual values, and the results were visualized in a clear and intuitive graph. This step highlighted the model's ability to follow data trends and predict future values.

## Step 7: Predicting Future Values

To generate future predictions, I implemented a rolling window approach:

The model used its predictions iteratively as input to forecast future values, enabling predictions far into the future without relying solely on historical data. This method proved effective for extending the model's utility.

## Step 8: Creating the Pipeline

## Step 6. Creating the Pipeline

To streamline the process, I combined all the steps into a modular pipeline:

Feature Engineering

Data Preprocessing

Model Training

Evaluation

This pipeline enables reproducibility and allows for easy experimentation with new data or models.

## Conclusion

This project was both challenging and rewarding, demonstrating the potential of deep learning models in time series forecasting. While my LSTM model successfully captured patterns and trends, there is room for improvement, such as incorporating additional features, experimenting with alternative architectures like GRUs or Transformers, and increasing the training dataset.