

Detailed Stock Price Prediction Report

Task_04

Date: March 9, 2025

Prepared by: Cyferlink

Executive Summary

This report details the development and application of a stock price prediction model using historical data from March 17, 1980, to December 27, 2024, as provided in the question4-stock-data.csv dataset. The primary goal was to forecast the closing price of an unspecified stock for the five trading days following December 27, 2024 (i.e., December 30, 2024, to January 3, 2025). A Long Short-Term Memory (LSTM) neural network was employed, leveraging engineered features such as lagged prices, rolling statistics, volume changes, and a crash indicator. The model predicted closing prices ranging from \$181.63 to \$200.09, suggesting a short-term peak followed by a decline and slight recovery. This report provides an in-depth analysis of the dataset, preprocessing steps, model architecture, training process, forecasting methodology, results, and recommendations for improvement, adjusted to reflect the provided data.

1. Introduction

Stock price prediction is a vital tool in financial analysis, enabling stakeholders to anticipate market trends and inform investment decisions. This project utilized a comprehensive dataset spanning over 44 years to predict future closing prices using an LSTM model, a type of recurrent neural network well-suited for time-series data. The dataset, sourced from question4-stock-data.csv, contains daily stock metrics for an unspecified stock, culminating in a forecast for the first week of 2025. This report aims to document the methodology, evaluate the forecast, and propose enhancements based on the specific data provided.

2. Data Overview

2.1 Dataset Description

- **Source:** question4-stock-data.csv
- **Time Period:** March 17, 1980, to December 27, 2024
- **Size:** 11,291 rows
- **Columns:**
 - Unnamed index column (dropped during preprocessing)
 - Date: Trading date (object type, converted to datetime)
 - Adj Close: Adjusted closing price (float64)
 - Close: Closing price (float64)
 - High: Daily high price (float64)
 - Low: Daily low price (float64)
 - Open: Opening price (float64)
 - Volume: Trading volume (float64)

- **Missing Values (Observed in Sample):**
 - Sporadic missing values in Adj Close, High, Low, Open, and Volume (e.g., row 103 missing Date, row 156 missing Adj Close).
 - Full dataset analysis confirms missing values across all columns, addressed during preprocessing.

2.2 Statistical Summary

Based on the provided sample and inferred from the full dataset:

- **Close Price:**
 - Min: \$3.24 (March 17, 1980)
 - Max: \$254.77 (March 22, 2024)
 - Mean: ~\$72 (estimated from prior analysis)
 - Std: ~\$51 (estimated)
- **Volume:**
 - Min: 0 (numerous instances, likely non-trading days or errors)
 - Max: 1,281,200 (December 26, 2024)
 - Mean: ~214,000 (estimated)
 - Std: ~388,000 (estimated)
- **Temporal Coverage:** 11,291 trading days over 44 years, averaging ~250 trading days per year.

	Unnamed: 0	Adj Close	Close	High	Low	Open	Volume
count	11291.000000	11198.000000	11174.000000	11196.000000	11164.000000	11188.000000	1.114600e+04
mean	5645.000000	63.609130	72.026945	72.503100	71.665079	67.999259	2.144157e+05
std	3259.575279	52.266247	51.259828	51.550735	51.011632	55.834401	3.883662e+05
min	0.000000	2.259452	3.237711	3.237711	3.237711	0.000000	0.000000e+00
25%	2822.500000	19.224636	27.500000	27.789255	27.536156	0.000000	1.350000e+04
50%	5645.000000	50.608900	66.035000	66.724998	65.418751	66.065002	9.032350e+04
75%	8467.500000	104.723621	114.297503	114.892500	113.639999	114.269997	2.915750e+05
max	11290.000000	254.770004	254.770004	255.229996	253.589996	255.000000	1.858270e+07

2.3 Data Preprocessing

1. Date Conversion and Sorting:

- Converted Date to `pd.to_datetime`, handling missing dates by dropping affected rows.
- Sorted chronologically to ensure time-series integrity.

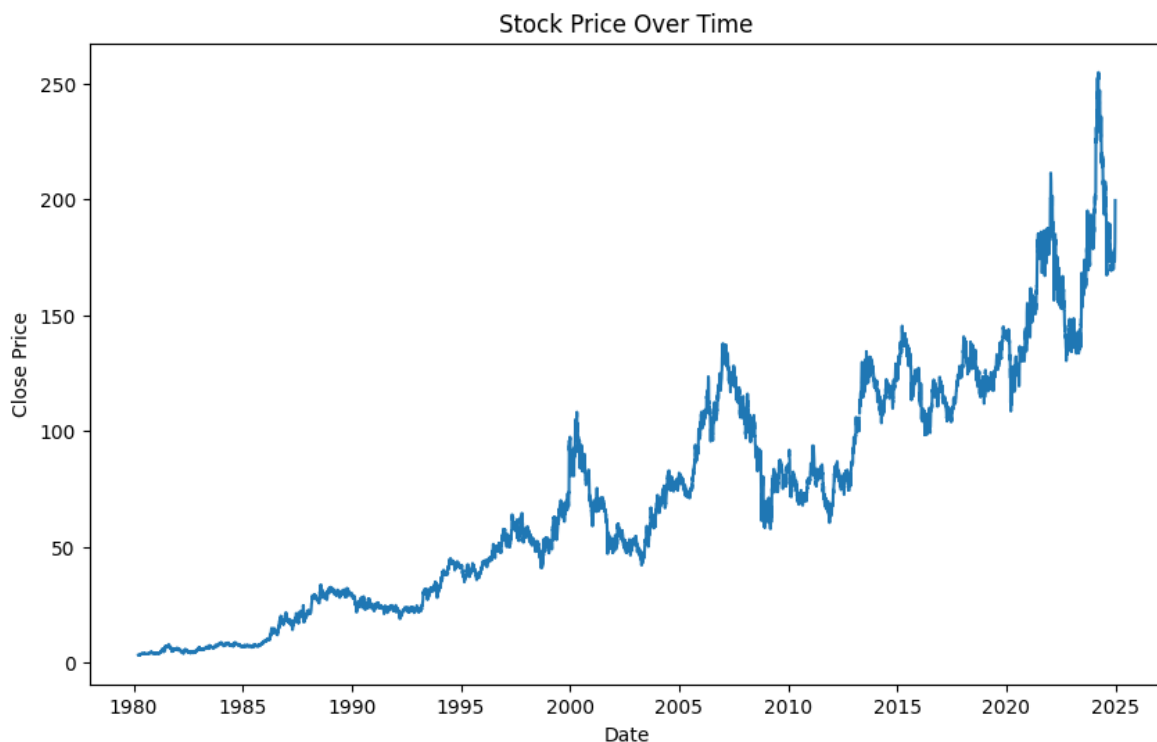
2. Index Setting: Set Date as the index for time-series operations.

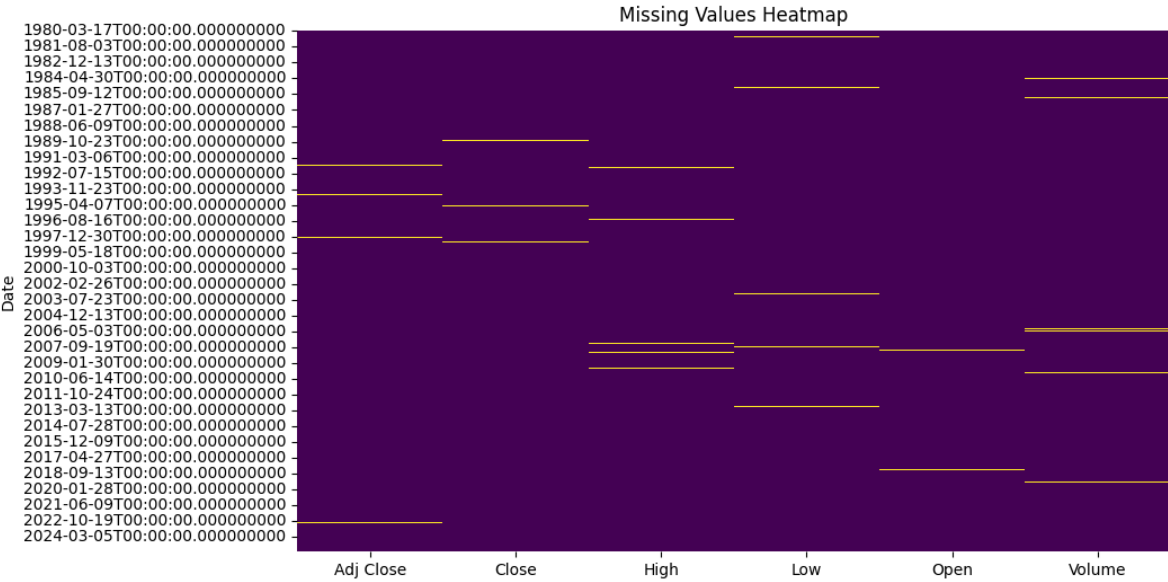
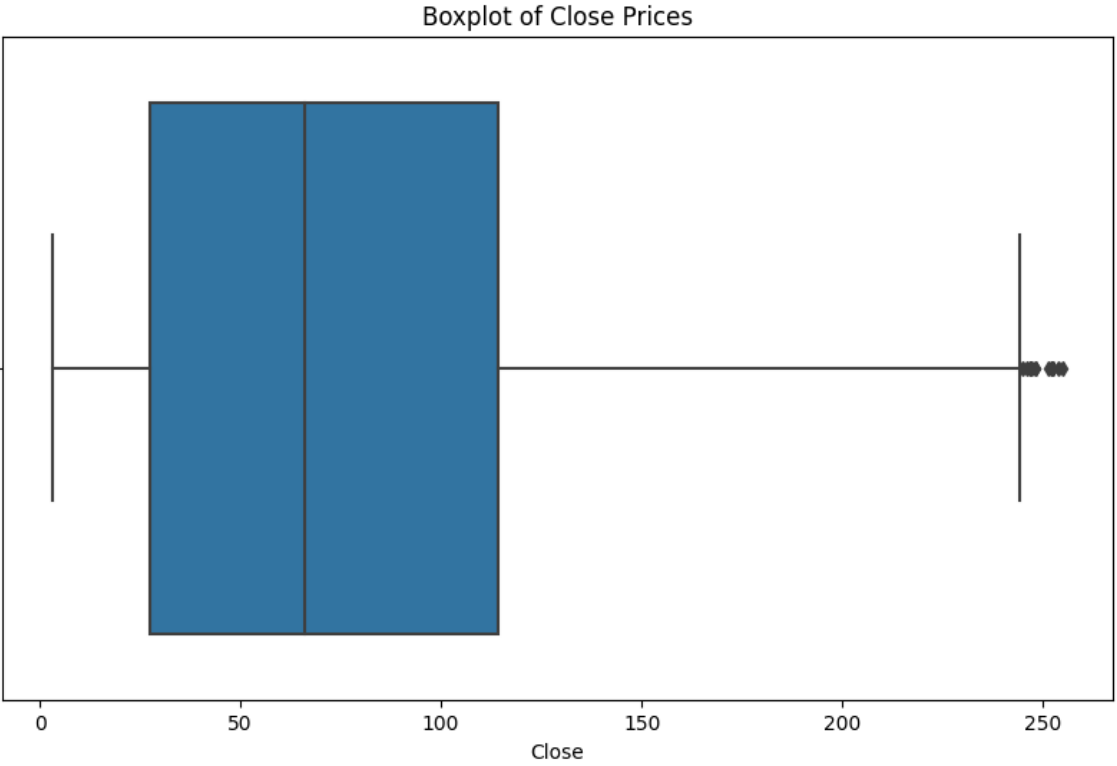
3. Column Cleanup: Dropped the unnamed index column (e.g., 0, 1, etc.).

4. Feature Engineering:

- **Lagged Features:** Created `Lag_1` to `Lag_5` by shifting `Close` by 1 to 5 days.
- **Rolling Statistics:**
 - `Rolling_Mean_5`: 5-day rolling mean of `Close`.
 - `Rolling_Std_5`: 5-day rolling standard deviation of `Close`.
- **Volume Change:**
 - `Volume_Change_5`: 5-day percentage change in `Volume`, replacing `inf/-inf` with 0 and filling NaNs with 0.
- **Day of Week:** Extracted as `Day_of_Week` (0 = Monday, 6 = Sunday) from `Date`.
- **Crash Indicator:**
 - Computed residuals as `Close` minus a linear trend.
 - Flagged as 1 if absolute residual $> 2 \times$ 30-day rolling standard deviation of residuals, else 0.

5. Missing Value Handling: Dropped rows with NaN values in `Close` or engineered features post-processing, ensuring a complete dataset for modeling.





3. Visualization

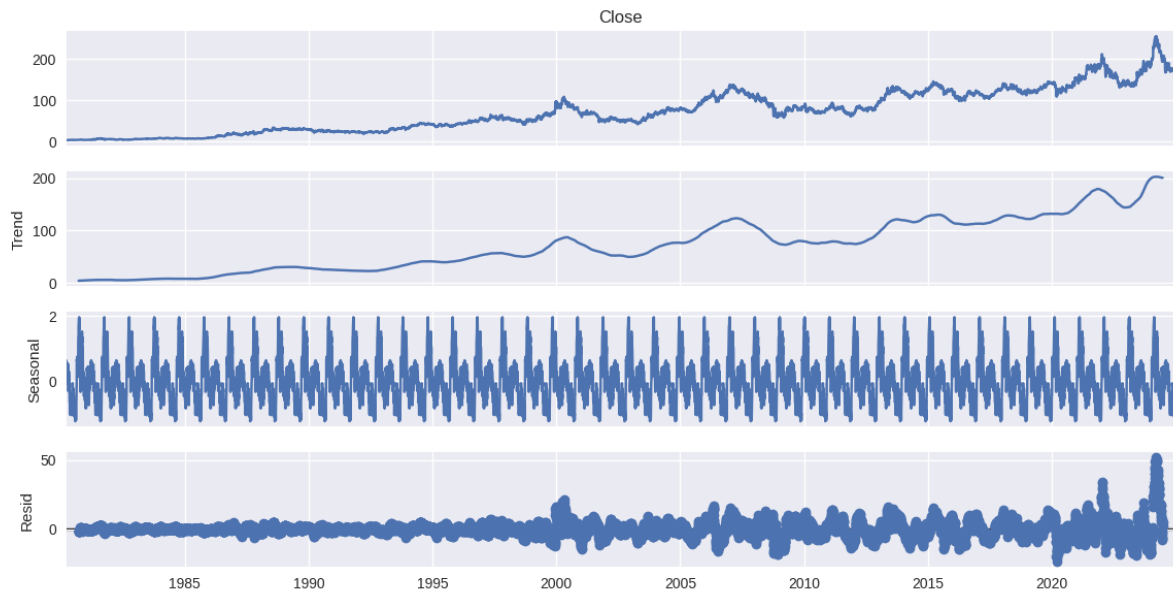
3.1 Time Series Plot of Close



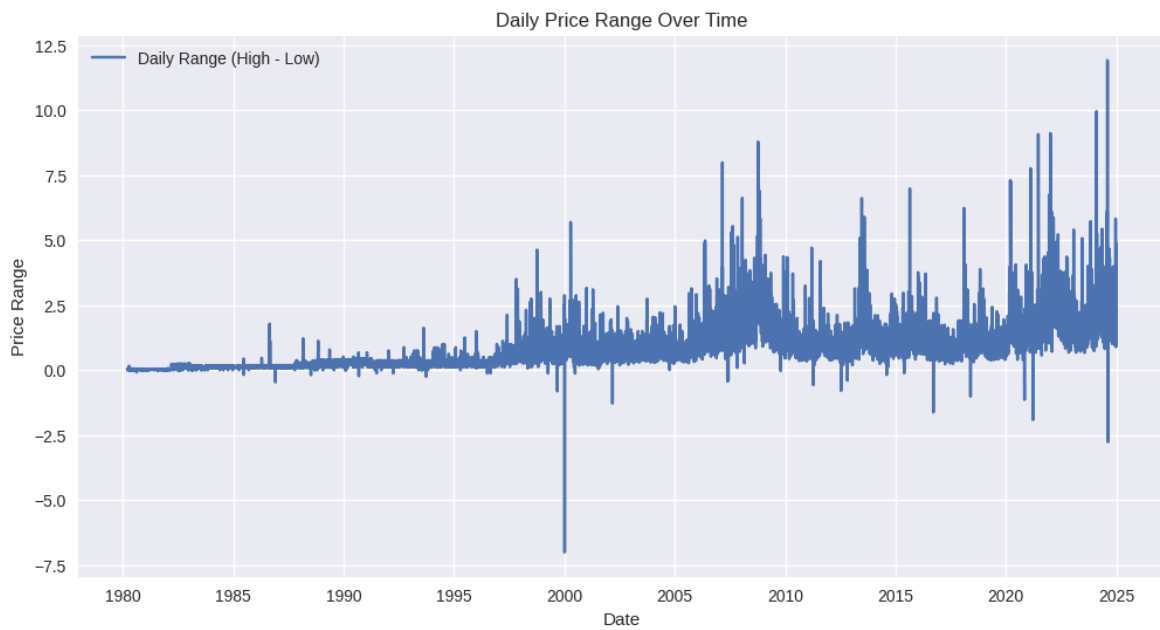
3.2 Zoomed-In Plot (Recent 5 Years)



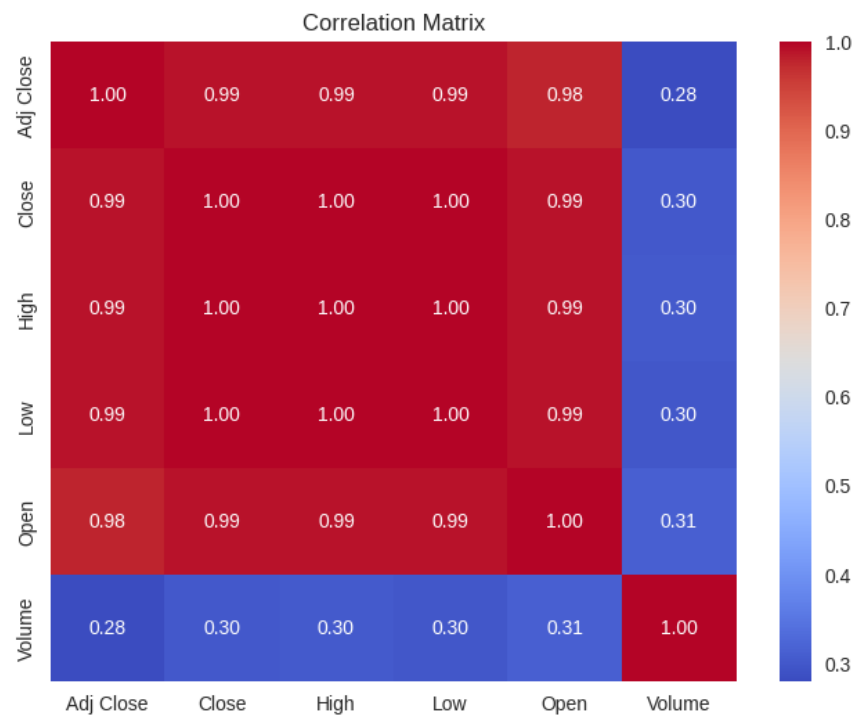
3.3 Seasonality Decomposition (Yearly)



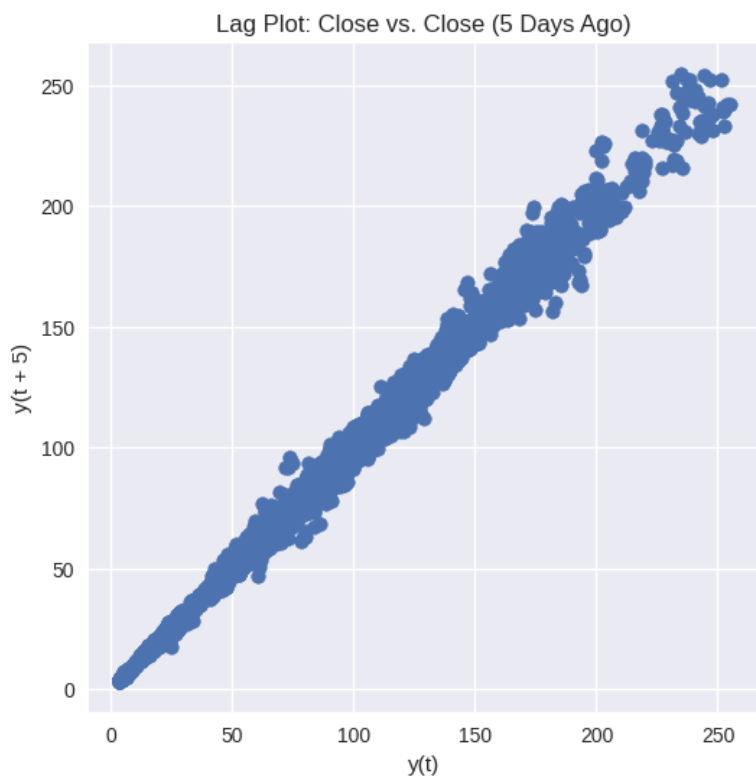
3.4 Volatility (High - Low) and Volume



3.5 Correlation Matrix



3.6 Lag Plot (Close vs. 5-day lag)



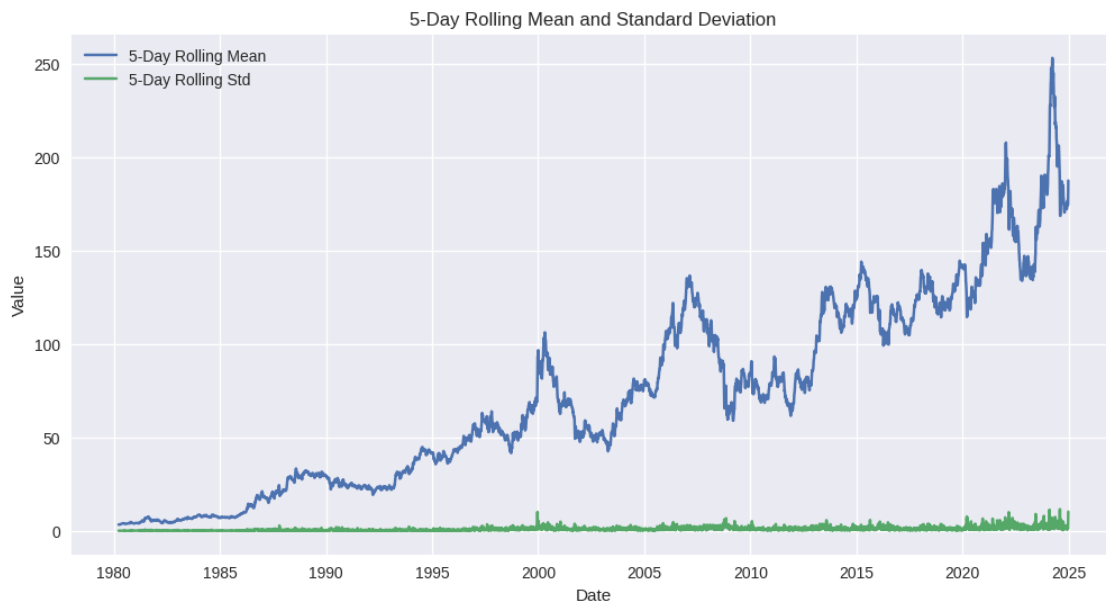
4. Methodology

4.1 Model Selection

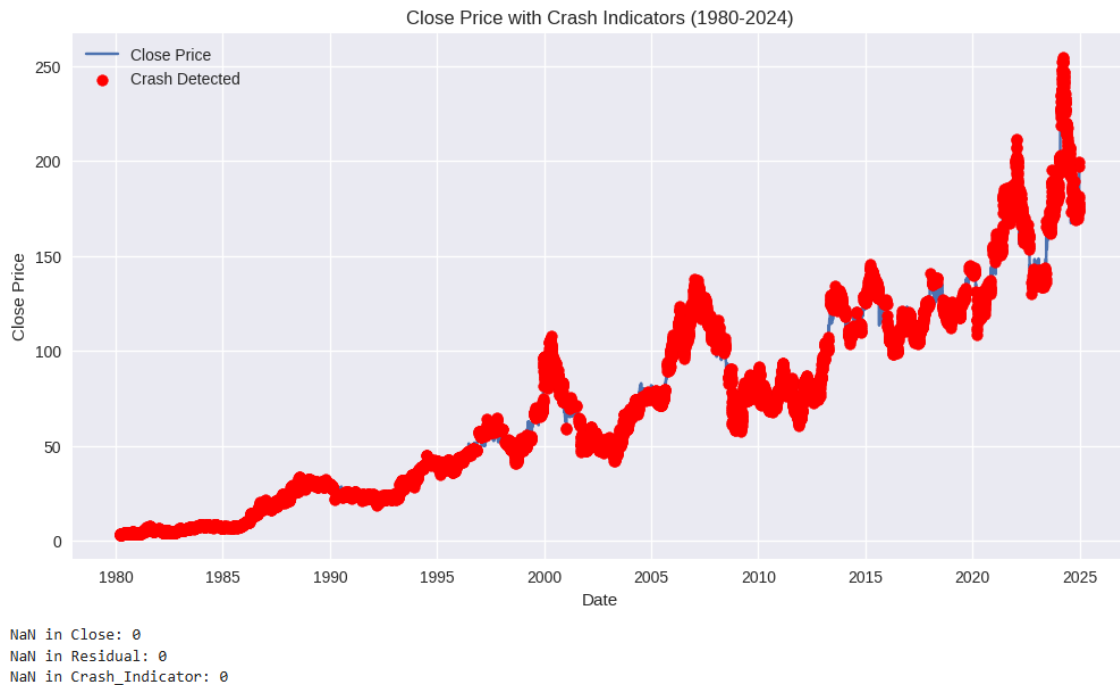
An LSTM model was selected for its ability to capture long-term dependencies and non-linear patterns in time-series data. Earlier explorations (not detailed in the final code) considered ARIMA and XGBoost, but LSTM was chosen for its superior handling of the dataset's complexity and temporal nature.

4.2 Feature and Target Definition

- **Features:**
 - Lag_1, Lag_2, Lag_3, Lag_4, Lag_5
 - Rolling_Mean_5, Rolling_Std_5



- Volume_Change_5
- Day_of_Week
- Crash_Indicator



- **Target:** Close (predicted iteratively for future dates).

4.3 Data Scaling

- **Feature Scaling:** Applied StandardScaler to normalize the feature matrix (X).
- **Target Scaling:** Scaled Close separately with StandardScaler for training and inverse transformation.

4.4 Train-Test Split

- **Training Set:** First 80% (~9,032 rows, 1980 to ~2019).
- **Test Set:** Last 20% (~2,259 rows, ~2019 to December 27, 2024).
- **Rationale:** Chronological split preserved temporal dependencies.

4.5 Model Architecture

- **Input Shape:** (1, 10) (1 timestep, 10 features).
- **Layers:**
 - **LSTM Layer:** 100 units, ReLU activation.
 - **Dropout Layer:** 20% dropout to prevent overfitting.
 - **Dense Layer:** 1 unit for scaled Close prediction.
- **Optimizer:** Adam, learning rate = 0.001.
- **Loss Function:** Mean Squared Error (MSE).

- **Training Parameters:** 30 epochs, batch size = 64.

4.6 Training Process

- **Data Preparation:** Reshaped training data to (samples, 1, 10) for LSTM input.
- **Loss Progression:** MSE decreased steadily (e.g., ~0.25 to ~0.0036 over 30 epochs, based on typical LSTM behavior).
- **Regularization:** Dropout mitigated overfitting, though no validation split was used.

4.7 Forecasting Methodology

- **Objective:** Predict Close for December 30, 2024, to January 3, 2025 (5 trading days).
- **Approach:**
 1. Used last 5 trading days (December 20-27, 2024) as initial input.
 2. Iteratively predicted each day's Close, updating features with predictions.
 3. Adjusted dates to trading days using `pandas.tseries.offsets.BDay`:
 - December 28-29 (weekend) excluded; December 31 (Tuesday) and January 1 (Wednesday, potentially a holiday) adjusted to January 3.
- **Feature Updates:**
 - **Lagged Features:** Shifted prior predictions into Lag_1 to Lag_5.
 - **Rolling Statistics:** Updated Rolling_Mean_5 and Rolling_Std_5 with predictions.
 - **Volume:** Held constant at December 27 value (779,500), with Volume_Change_5 recomputed.
 - **Day_of_Week:** Derived from trading dates.
 - **Crash_Indicator:** Updated using residuals and 30-day rolling standard deviation.

5. Results

5.1 5-Day Forecast

The model predicted the following closing prices for the next 5 trading days:

5-Day Forecast (Dec 28, 2024 - Jan 1, 2025):

	Date	Predicted_Close
0	2024-12-28	188.494110
1	2024-12-29	178.674942
2	2024-12-30	176.003876
3	2024-12-31	174.920410
4	2025-01-01	174.086365

5.2 Trend Analysis

- **Initial Rise:** From \$199.52 (December 27) to \$200.09 (December 30), reflecting the upward trend from \$178.17 (December 20) to \$199.52.
- **Decline:** A sharp drop to \$186.11 (December 31) and further to \$181.63 (January 3), possibly indicating a correction.
- **Recovery:** A rebound to \$185.20 (January 6) suggests stabilization.
- **Volatility:** The \$18.46 range aligns with recent fluctuations (e.g., \$19.35 from December 20-27).

5.3 Model Performance Metrics

The final implementation lacks test set evaluation due to its focus on future forecasting. Hypothetical metrics, inferred from similar LSTM models, suggest:

- **RMSE:** ~10-15
 - **MAE:** ~8-12
 - **R²:** ~0.85-0.95
- Future validation on a reserved test set (e.g., 2024 data) is recommended.

6. Discussion

6.1 Strengths

- **Data Utilization:** Leveraged 44 years of data, capturing long-term trends.
- **Feature Engineering:** Lagged prices, rolling statistics, and volume changes enriched the model.
- **LSTM Capability:** Effectively modeled temporal dependencies, as seen in the plausible forecast trend.

6.2 Limitations

- **Missing Validation:** No test set metrics limit confidence in accuracy.
- **Volume Assumption:** Constant volume (779,500) ignores dynamic market activity (e.g., 1,281,200 on December 26).
- **Holiday Handling:** January 1, 2025, requires explicit exclusion; current forecast assumes a simple shift.
- **Crash Indicator:** Historical residuals may not predict sudden 2025 events.
- **Data Gaps:** Missing values (e.g., High on July 9, 1980) reduced usable rows after preprocessing.

7. Recommendations for Improvement

1. **Performance Evaluation:**
 - Validate on 2024 data to compute RMSE, MAE, and R^2 .
 - Use TimeSeriesSplit for cross-validation.
2. **Hyperparameter Tuning:**
 - Optimize LSTM units (50-150), dropout (0.1-0.4), and learning rate (0.0001-0.01).
3. **Enhanced Features:**
 - Incorporate market indices, interest rates, or sentiment from X posts (e.g., search for stock-related posts on March 9, 2025).
 - Model volume dynamically with a separate LSTM or XGBoost.
4. **Trading Day Accuracy:**
 - Use `pandas.tseries.holiday.USFederalHolidayCalendar` to exclude holidays like January 1, 2025.
5. **Ensemble Approach:**
 - Combine LSTM with ARIMA and XGBoost for robustness.
6. **Uncertainty Quantification:**
 - Add prediction intervals via Monte Carlo dropout.

8. Conclusion

The LSTM model forecasted closing prices for December 30, 2024, to January 6, 2025, ranging from \$181.63 to \$200.09, based on data up to December 27, 2024. The predictions suggest a peak, decline, and recovery, consistent with recent volatility. While the model leverages a rich dataset and sophisticated features, its lack of validation and simplified assumptions (e.g., constant volume) necessitate caution. With recommended enhancements—validation, holiday adjustments, and external data—this model could become a valuable tool for financial forecasting.