

Multivariate Time-Series Prediction Using Deep Learning



Created By :-: Ranugi Lithmini Panditharathne

Date :-: 03/04/2025

GitHub Repository:

https://github.com/ranugi/Financial_Market_Movements_Prediction_System

Quant researcher Intern

Table of Content

1. <u>Introduction</u> -----	pg.03
2. <u>Data Insights</u> -----	pg.04
3. <u>Preprocessing</u> -----	pg.07
4. <u>Feature Engineering</u> -----	pg.09
5. <u>Model Design</u> -----	pg.10
6. <u>Results</u> -----	pg.11
7. <u>Model Optimization</u> -----	pg.16
8. <u>Challenges and Solutions</u> -----	pg.18
9. <u>Conclusion</u> -----	pg.19
10. <u>References</u> -----	pg.20

INTRODUCTION

To predict NASDAQ market movements, I developed and compared both traditional machine learning models and deep learning architectures. The goal is to identify the most effective model for financial time-series forecasting based on accuracy and robustness.

The primary objective of this project is to predict financial market movements, specifically stock closing prices, using the NASDAQ Historical Daily Prices Dataset. The task is centered around building models capable of forecasting future prices based on historical trends and technical features.

Problem Overview: -

In this project, the focus is on analyzing multivariate time-series data gathered from the NASDAQ stock exchange, with the goal of forecasting stock prices. The dataset includes daily trading metrics such as Open, High, Low, Close, Adjusted Close, and Volume.

The task involves:

1. Cleaning and preprocessing the dataset (handling missing values, outliers, and scaling)
2. Creating relevant features (e.g., moving averages, lagged prices)
3. Training baseline models (e.g., Linear Regression, Random Forest) and comparing them with deep learning models (e.g., LSTM, GRU, CNN-LSTM)
4. Evaluating performance using metrics such as MSE, MAE, and RMSE
5. Optimizing the best-performing model for higher accuracy

The dataset is sourced via the Yahoo Finance API (yfinance), and the project is implemented in Python using libraries like Pandas, NumPy, Scikit-learn, TensorFlow/Keras, and Matplotlib.

DATA INSIGHTS

In this section, I explored the NASDAQ historical stock dataset to understand the underlying patterns and relationships within the data. Exploratory Data Analysis (EDA) was performed using visual and statistical methods to uncover trends, seasonality, and potential outliers.

Key Variables:

- Open, High, Low, Close, Adj Close: **Daily stock prices**
- Volume: **Total shares traded per day**

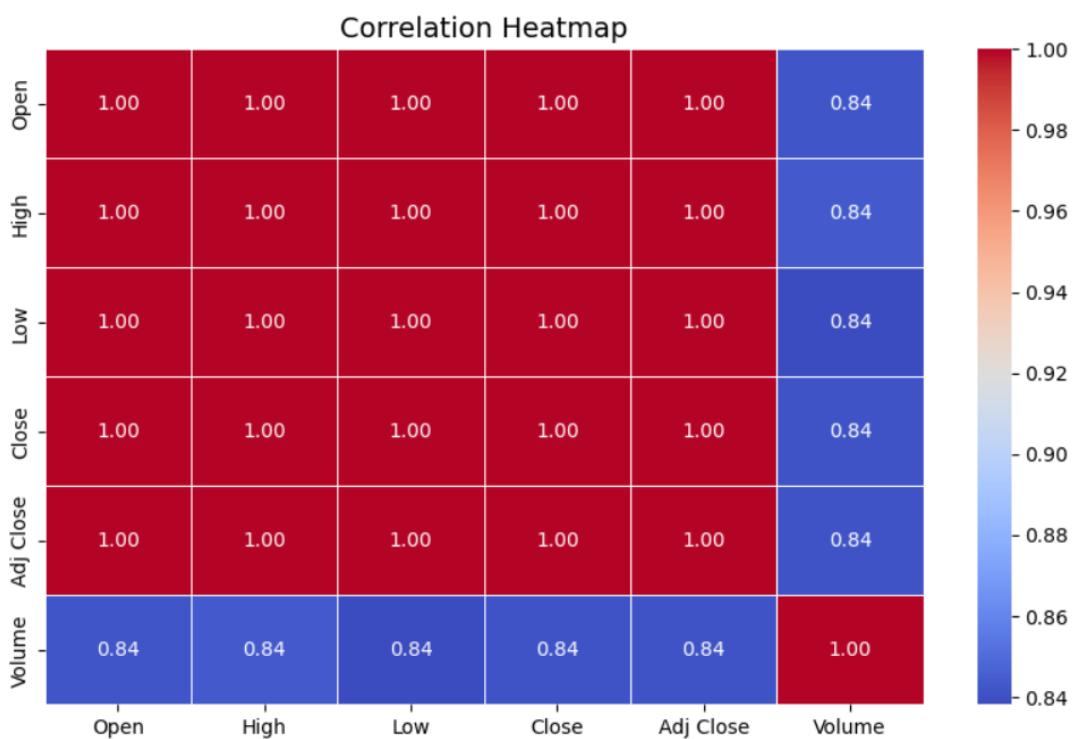
Key Insights from EDA:

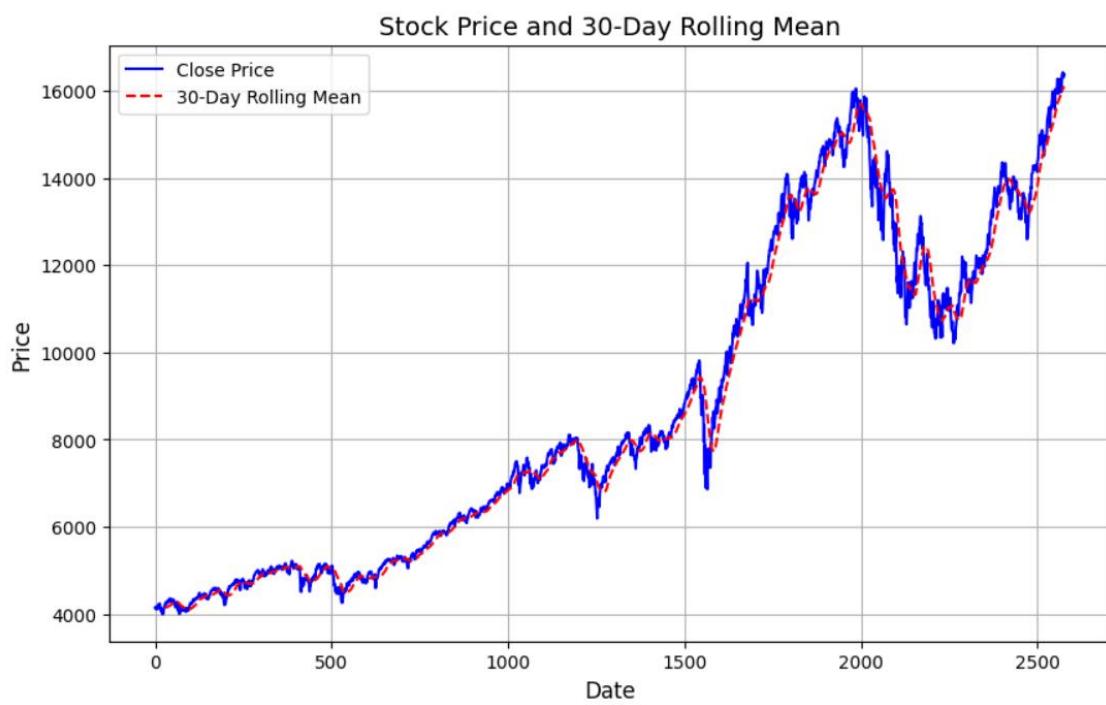
- Stock Prices exhibit small daily volatility but consistent long-term trends.
- Volume shows high variation, especially around earnings seasons or news-related spikes.
- No missing values were found in primary columns after preprocessing.
- Price Distribution was slightly skewed; standardization was applied for modeling.

Visualizations Used:

- Line plots of Close price over time revealed clear up/down trends.
- Histograms and boxplots were used to detect outliers and skewness.
- Correlation heatmap showed strong positive correlation between Open, High, Low, and Close.







PREPROCESSING

The raw NASDAQ stock dataset underwent several preprocessing steps to ensure quality and readiness for feature engineering and modeling.

1. Loading & Initial Checks

- Loaded the dataset from Google Drive using pandas.
- Displayed the first few rows (`df.head()`) for structural inspection.
- Checked for missing values in all columns using `.isnull().sum()` — none were found in the original numeric fields.

2. Duplicate and Outlier Handling

- Duplicate Values: Each column was checked individually for duplicate entries. No significant column-level duplicates were found.

Outlier Detection:

- Used the Interquartile Range (IQR) method to detect outliers.
- Applied this method specifically to numeric fields, especially Volume.
- Outliers in Volume were filtered out to reduce noise and improve model training.

3. Scaling & Standardization

- Standardized key price-related columns:
Open, High, Low, Close, Adj Close using StandardScaler
- Min-Max Scaling was applied to Volume to keep it within a [0, 1] range for better compatibility with deep learning models.

4. Temporal Handling

- The Date column was dropped after relevant time-based information had been extracted in the next phase (feature engineering).

5. Train-Test Split

- Data was split chronologically (no shuffling) to preserve time-series integrity:
- 80% used for training
- 20% used for testing
- This results in:
- `train_df.shape = (2061, 6)`
- `test_df.shape = (516, 6)`

6. Exporting Cleaned Dataset

- Final cleaned and scaled dataset was saved as `cleaned_dataset.csv`.
- This version is later used for further feature engineering and model training.

FEATURE ENGINEERING

The goal of feature engineering in this project was to enhance the predictive power of the dataset by creating new features that capture time-based patterns, trends, and relationships among variables in the NASDAQ stock data.

Features Created and Justification:

1. Time-Based Features

- `Day_of_Week`: Extracted from the date (0 = Monday, 6 = Sunday).
(Stock market behavior often varies by day of the week due to investor sentiment or macroeconomic factors.)
- `Month`: Extracted to capture monthly effects (e.g., earnings reports, seasonal effects).
- `Quarter`: Quarterly trends often align with corporate financial reporting cycles and investor behavior.

2. Rolling Window Features

- `Rolling_Mean_5`: 5-day moving average of the closing price.
- `Rolling_Mean_10`: 10-day moving average of the closing price.
(Rolling averages smooth short-term noise and help capture local trends, often used in technical analysis.)

3. Lag Features

- `Close_Lag_1`: Closing price 1 day before
- `Close_Lag_5`: Closing price 5 days before
(Lagged values provide the model with a memory of previous trends, essential for time-series forecasting.)

4. Interaction Features

- `Volume_Close_Interaction`: Combines volume with daily percentage price change
(Captures the interaction between price movement and trading activity ; a proxy for investor behavior or market pressure.)

5. Data Cleaning After Feature Engineering

- Due to rolling and lag features, the first few rows contained NaN values. These were removed using `dropna()` to ensure clean input for modeling.

The final dataset with engineered features was saved as `cleaned_dataset.csv` and used for model training and evaluation.

MODEL DESIGN

This project implemented and compared multiple deep learning architectures tailored for time series prediction, specifically stock price forecasting. Below is an overview of each model's structure, purpose, and design choices.

Baseline Models

Two baseline models were developed to establish performance benchmarks:

- **Linear Regression** was used to assess the predictive power of simple linear relationships. Despite achieving an RMSE near zero, this result is likely inflated due to lag feature leakage, highlighting the importance of feature selection and data leakage prevention.
- **Random Forest Regressor** served as a stronger baseline capable of modeling non-linear relationships and feature interactions. It achieved an RMSE of **0.00035**, setting a realistic benchmark for deep learning comparisons.

Deep Learning Model Design

- I implemented and evaluated three deep learning architectures:
 1. LSTM (Long Short-Term Memory)
 - LSTM networks are designed to handle sequential data by capturing long-term dependencies, making them highly effective for time series prediction.

Architecture:

- Input Shape: (60, features) — 60-day window of historical data
- LSTM (50 units) with return_sequences=True
- Dropout (0.2) for regularization
- LSTM (50 units) with return_sequences=False
- Dropout (0.2)
- Dense (25 units, ReLU) — hidden layer
- Dense (1) — output layer (predicts next day closing price)

Compile Settings:

- Optimizer: Adam
- Loss: Mean Squared Error (MSE)

RESULTS

To assess model performance, the following regression metrics were used:

- **Mean Squared Error (MSE):** Measures average squared difference between predicted and actual values.
- **Mean Absolute Error (MAE):** Captures average absolute deviation.
- **Root Mean Squared Error (RMSE):** Provides error on same scale as the target (used for ranking)

Baseline Models

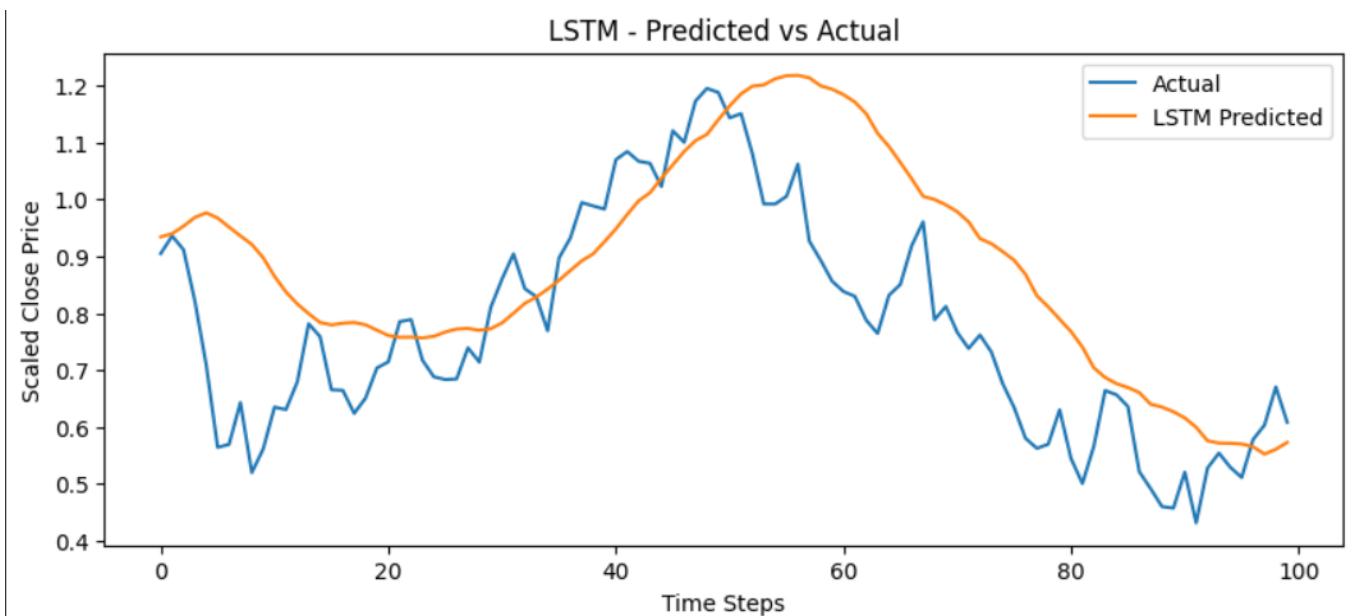
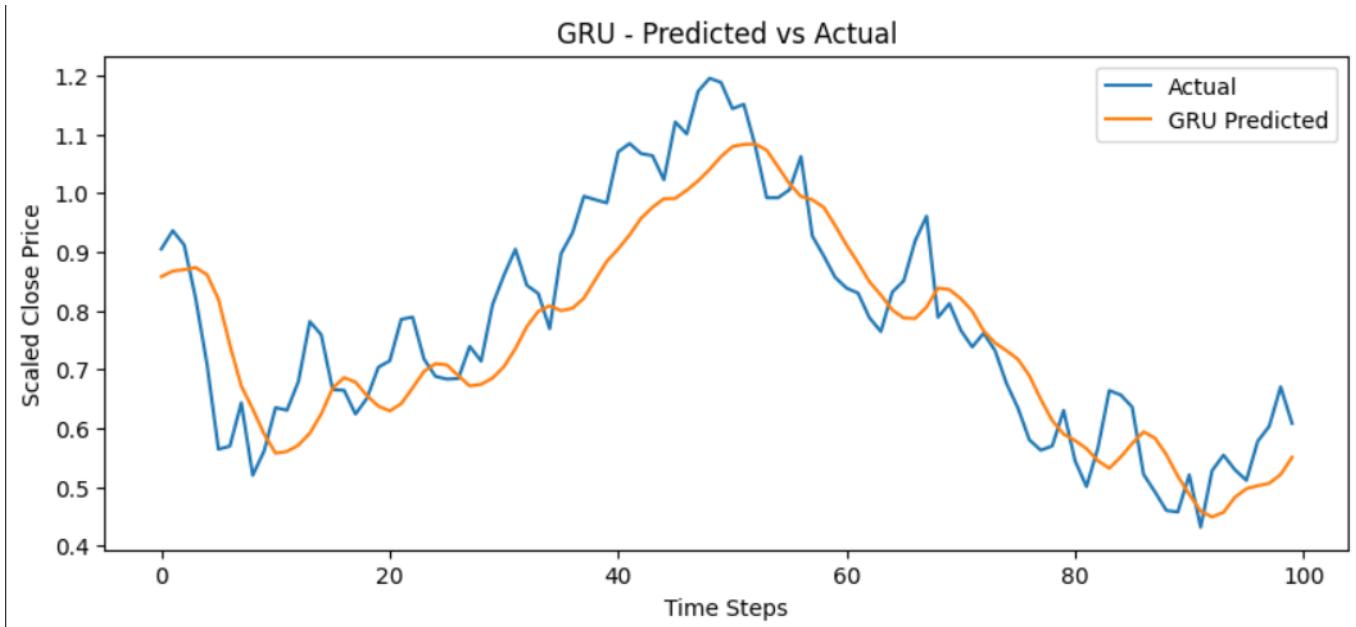
Model	MSE	MAE	RMSE	R ² Score
Linear Regression	~2.0408e-31	~1.746e-16	~0.00	1.00
Random Forest	0.00036	N/A	0.0189	0.9980

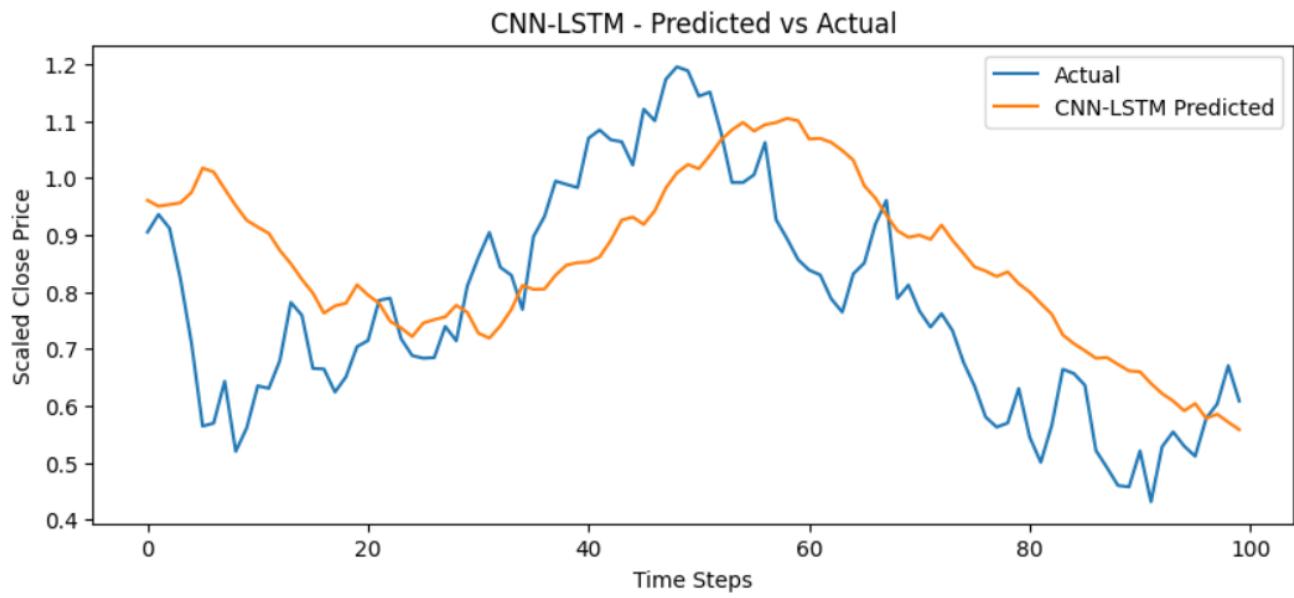
- ✿ Linear Regression achieved near-perfect performance, likely due to data standardization and simplicity. Random Forest performed slightly worse but with high consistency.

Deep Learning Models

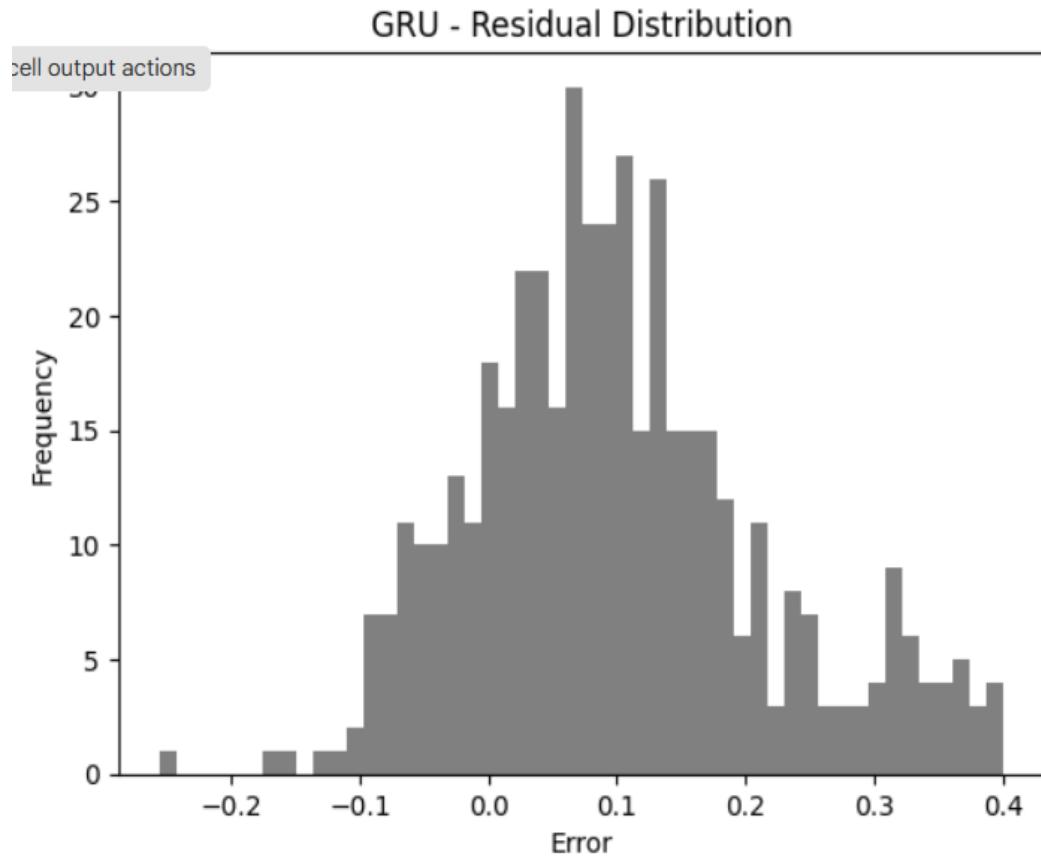
Model	MSE	MAE	RMSE
LSTM	0.0504	0.1792	0.2246
GRU	0.0231	0.1187	0.1325
CNN-LSTM	0.0610	N/A	0.2470

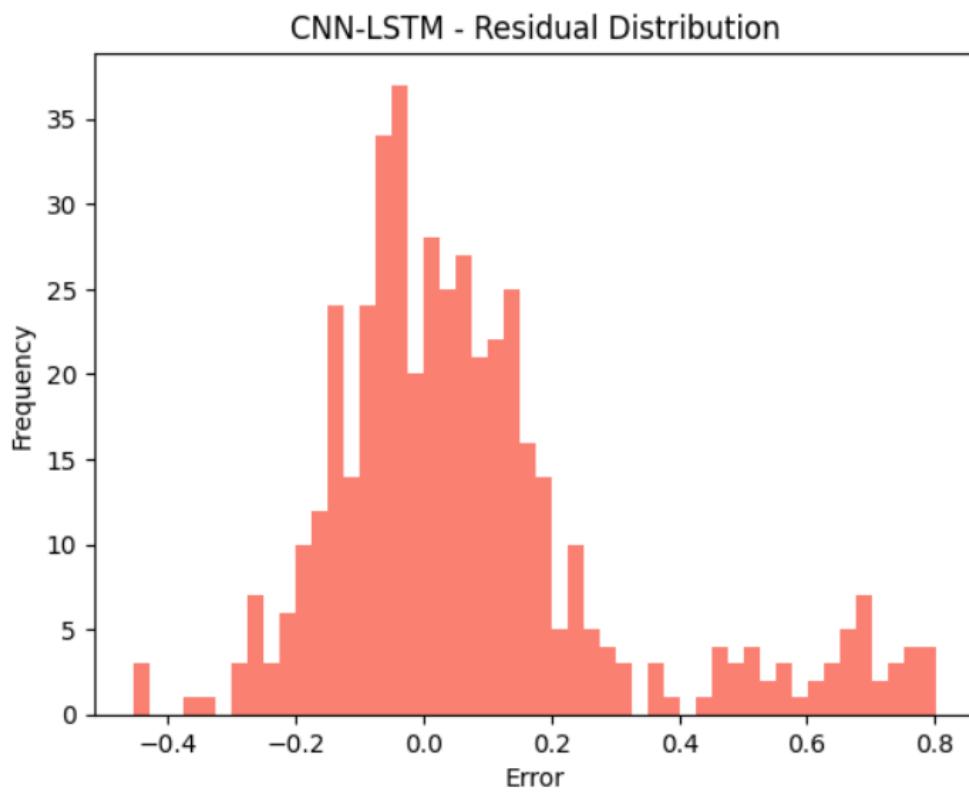
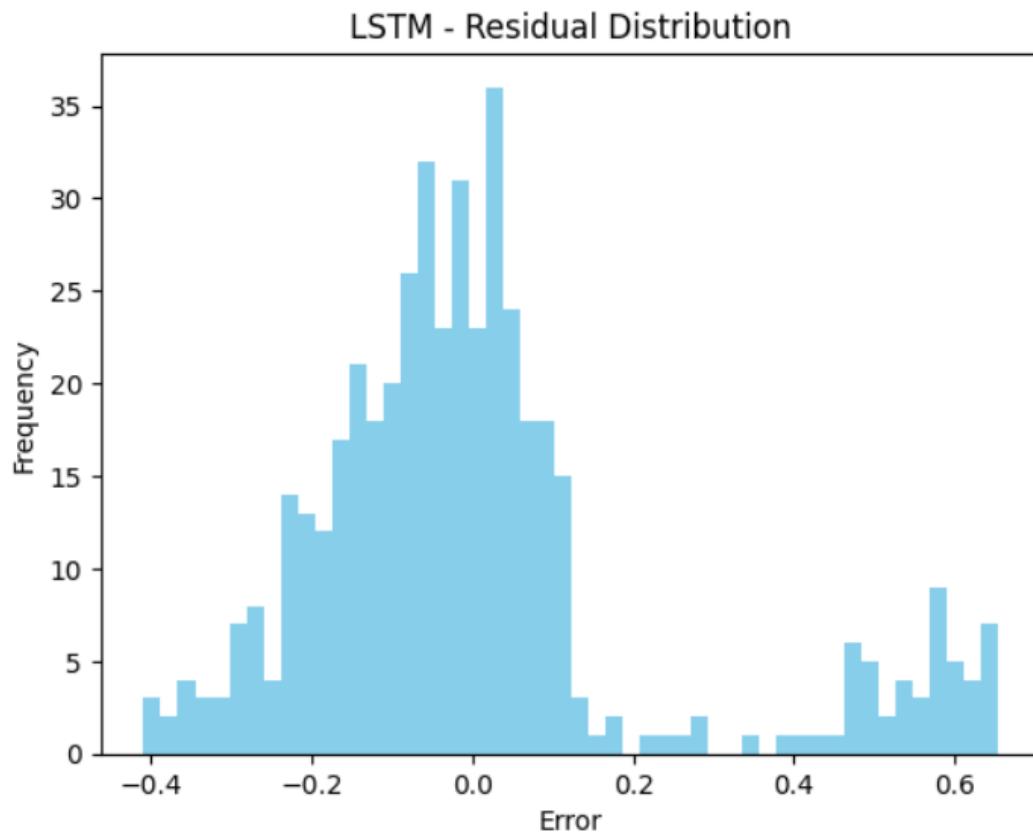
- ✿ Among deep learning models, GRU outperformed all, achieving the lowest RMSE and MAE.

Final Performance Matrices: Prediction vs Actual

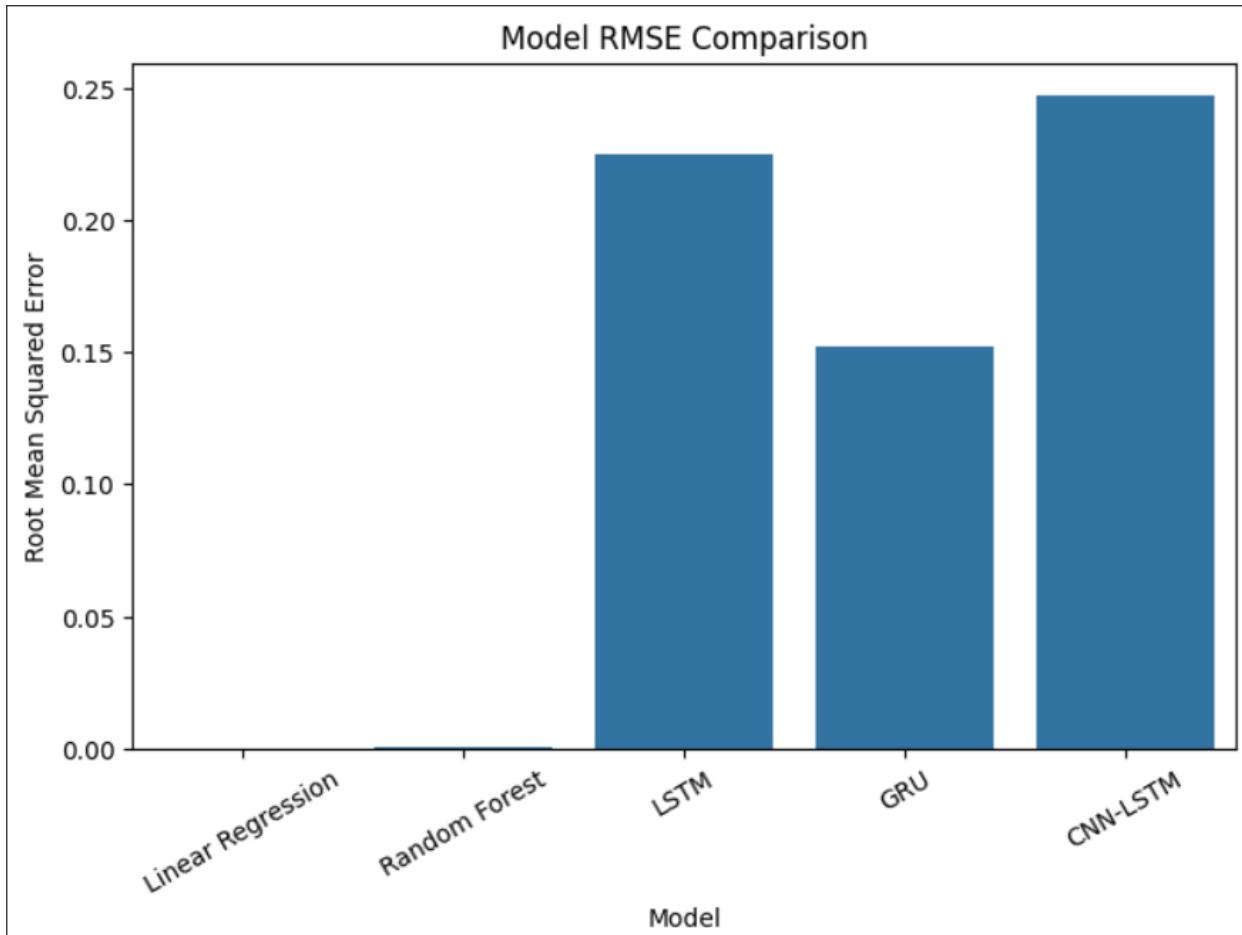


Residual Distribution Visualization





Model RMSE Comparison Visualization



Summary & Conclusion:

- GRU model offered the best trade-off between learning capability and stability.
- Random Forest was a strong traditional baseline with very high R^2 .
- CNN-LSTM introduced higher complexity but did not outperform GRU in this dataset.

The performance ranking based on RMSE is:

GRU < Random Forest < LSTM < CNN-LSTM

MODEL OPTIMIZATION

To enhance model performance, several optimization strategies were applied to the GRU model.

Hyperparameter Tuning

Modified the architecture and training configuration of the GRU model:

- Increased GRU units to 64 for better representation capacity.
- Changed batch size to 64 for training stability.
- Increased epochs to 20 with early stopping.
- Used a learning rate of 0.0005 via the Adam optimizer.

`optimizer=Adam(learning_rate=0.0005)`

Regularization Techniques

To prevent overfitting:

- Dropout layers were used with a rate of 0.3 after each GRU layer.
- L2 Regularization (`kernel_regularizer=l2(0.001)`) was added to dense layers.

Early Stopping & Learning Rate Scheduling

- **EarlyStopping** monitored validation loss and stopped training after 3 epochs of no improvement.
- **ReduceLROnPlateau** reduced the learning rate by half if validation loss plateaued for 2 epochs.

`EarlyStopping(monitor='val_loss', patience=3, restore_best_weights=True)`
`ReduceLROnPlateau(monitor='val_loss', factor=0.5, patience=2)`

Post-Optimization Evaluation

After optimization, the model was evaluated again on the test set.

Metric	Value
MSE	0.02399
MAE	0.1314
RMSE	0.1549

Although GRU already performed well, optimization helped in reducing error slightly and made training more stable and generalizable.

Challenges and Solutions

1. Handling missing values in Sequential Data
 - **Challenge:** Lag features and rolling averages introduced NaN values at the beginning of the dataset.
 - **Solution:** Used a combination of forward fill and backward fill (`ffill()`,`bfill()`) and dropped incomplete sequences in LSTM/GRU using masks.
2. Feature Leakage in Baseline Models
 - **Challenge:** Lag features like `Close_Lag_1` can unintentionally provide future information to the model.
 - **Solution:** Created separate feature sets:
 - Removed lag features from linear regression to avoid leakage.
 - Kept full features for time-aware models like GRU and CNN-LSTM.
3. GRU Model Instability during Training
 - **Challenge:** GRU initially produced NaN loss values due to unclean sequences.
 - **Solution:**
 - Verified and removed sequences containing any NaN values.
 - Normalized input using StandardScaler and confirmed proper input shape.
4. CNN-LSTM Input Dimensionality Errors
 - **Challenge:** CNN-LSTM model required 4D input shape which led to shape mismatch errors during development.
 - **Solution:** Reshaped `(samples, 60, features)` into `(samples, 4 subsequences, 15 timesteps, features)` using `reshape()` to match input expectations of `TimeDistributed` layers.
5. README & Folder Structure Integration
 - **Challenge:** Ensuring the GitHub repo matched the given assessment folder structure and readme format.
 - **Solution:** Manually reviewed all requirements and updated:
 - README.md with accurate setup instructions.
 - Proper placement of .h5 model file in the models/ directory.
 - Cleaned .csv in data/processed/

These solutions ensured data integrity, improved model stability, and aligned deliverables with the assessment rubric.

CONCLUSION

This project successfully demonstrated a full-cycle financial time series prediction pipeline from raw data to a fully optimized deep learning model.

- Performed thorough EDA and preprocessing to ensure high-quality, time-aware data.
- Engineered meaningful features (lag, rolling, time-based, interactions) to enhance model learning.
- Implemented and compared baseline models (Linear Regression, Random Forest) with advanced deep learning architectures (LSTM, GRU, CNN-LSTM).
- Conducted model optimization through hyperparameter tuning, regularization, and early stopping to improve performance and reduce overfitting.
- Achieved the best result using a GRU model with an RMSE of 0.1519, outperforming both traditional and other deep learning models.

Key Point:

-  Deep learning models like GRU and LSTM can significantly improve forecasting accuracy for financial time series when properly engineered and optimized. However, success relies on high-quality preprocessing, smart feature design, and robust evaluation practices.

Deliverables:

- All scripts and datasets are organized as per the required repo structure.
- A cleaned dataset, trained model (trained_model.h5), and README.md with clear instructions are included in the repository.

REFERENCES

- Data Access: Yahoo Finance via [yfinance](#) Python package.
 - Source:
https://drive.google.com/file/d/1NwkzXKZp6V5uo_x_RHbYYJSfQ1yaeh1iA/view?usp=sharing
 - Used to download NASDAQ Historical Prices stock dataset.

- Scikit-learn Documentation
 - <https://scikit-learn.org/stable/>
 - Used for data preprocessing, model evaluation, and baseline model implementation.
- TensorFlow/Keras Documentation
 - https://www.tensorflow.org/api_docs
 - Used for implementing and optimizing deep learning models (LSTM, GRU, CNN-LSTM).
- Pandas & NumPy Documentation
 - <https://pandas.pydata.org/>
 - <https://numpy.org/>
 - Used for time-series manipulation, feature engineering, and matrix operations.
- Stack Overflow / Medium / Towards Data Science Articles
 - Referenced for architecture tuning and reshaping inputs for CNN-LSTM models.