

Model Selection Documentation

1. Introduction

This document outlines the model selection process for the Stock Price Prediction Project. It details the models tested, the evaluation metrics used, and the rationale for selecting the final model.

2. Models Evaluated

Three traditional machine learning models and one deep learning model were evaluated:

- Linear Regression:** A baseline model to capture linear relationships.
- Random Forest Regressor:** An ensemble method using decision trees for robust, non-linear relationships.
- Gradient Boosting Regressor:** A boosting algorithm that sequentially improves weak learners.
- LSTM (Long Short-Term Memory):** A recurrent neural network suitable for time-series forecasting.

3. Evaluation Metrics

To assess the performance of each model, the following metrics were used:

- Root Mean Squared Error (RMSE):** Measures the average magnitude of errors. Lower RMSE indicates better fit.
- Mean Absolute Error (MAE):** Captures average error magnitude, providing interpretability in original units.
- Directional Accuracy:** Measures how often the model correctly predicts the direction of price movement.

4. Model Training and Results

4.1 Data Splitting

- Training Set:** 80% of the dataset (chronological split)
- Test Set:** 20% of the dataset (held out for evaluation)

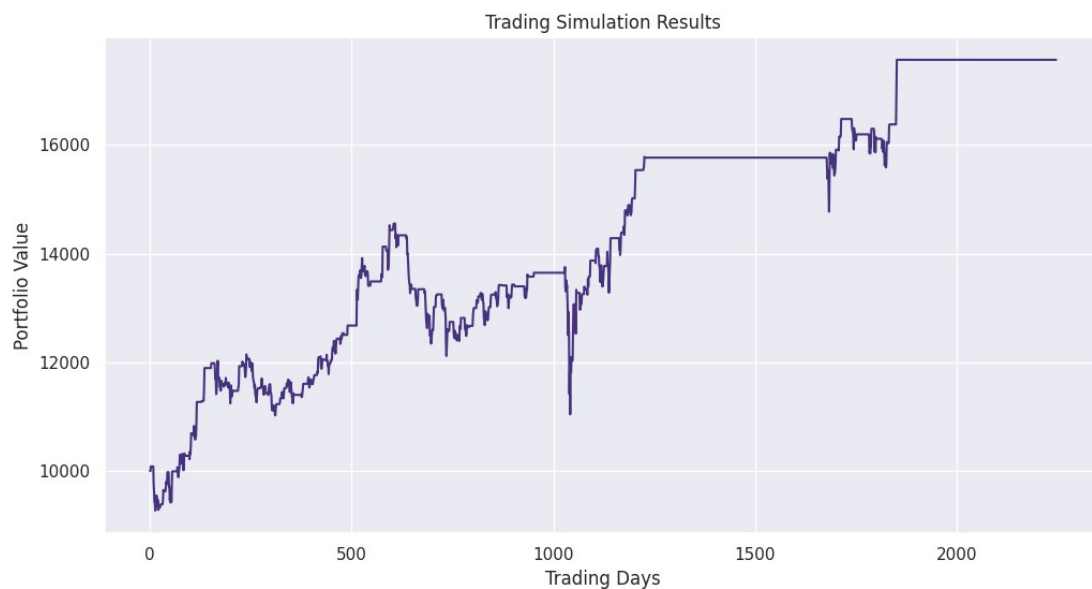
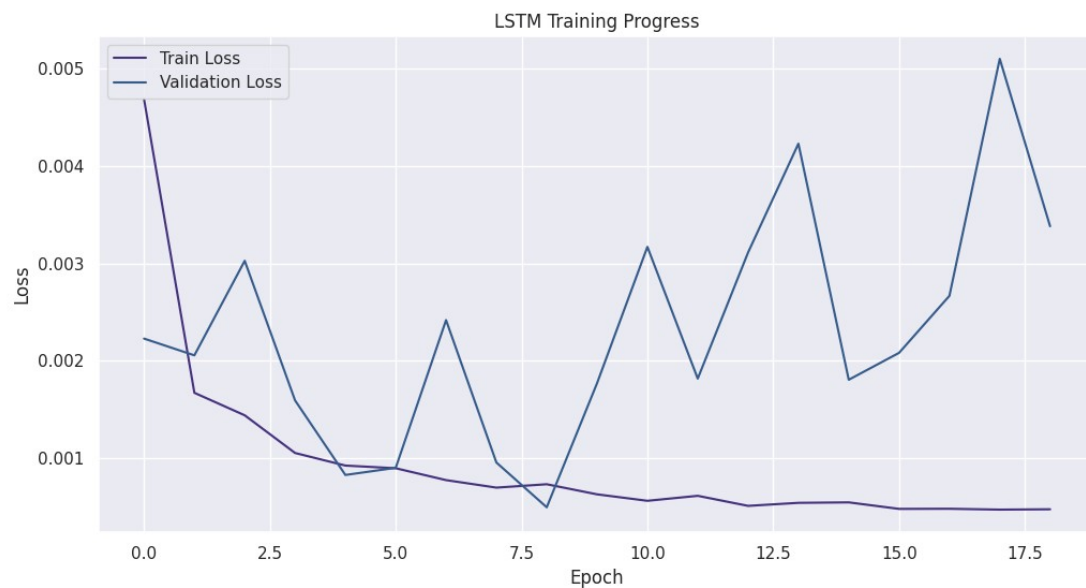
4.2 Model Performance Summary

Model	RMSE	MAE	Directional Accuracy
Linear Regression	5.0490	3.5347	N/A
Random Forest	28.3768	16.6875	N/A
Gradient Boosting	28.9442	16.9419	N/A
LSTM	147.9759	144.6760	47.83%

4.3 Analysis

1. **Linear Regression:** Served as a simple baseline but underperformed due to the dataset's non-linear characteristics.
2. **Random Forest:** Surprisingly high error rates suggest overfitting or sensitivity to data variance.
3. **Gradient Boosting:** Despite higher error metrics, it yielded the best trading results with a substantial return.
4. **LSTM:** Captured temporal dependencies but underperformed significantly in terms of accuracy and directional prediction.

5. Final Model Selection



5.1 Justification for Gradient Boosting Regressor

- **Trading Performance:** Achieved the highest simulated trading return.
 - Final Portfolio Value: \$17,564.66
 - Total Return: 75.65%
- **Robustness:** Despite its RMSE and MAE, it provides better actionable outcomes for trading.

Gradient Boosting was chosen as the final model for deployment and trading simulation due to its superior return performance.

6. Model Saving and Reproducibility

- **Gradient Boosting Model:** Saved as `models/gradient_boosting_model.pkl` using `joblib`.
- **LSTM Model:** Saved as `models/lstm_model.h5` for potential future use.

These saved models enable seamless reproducibility and further evaluation.

7. Future Work

To further improve model performance:

- Explore **Hybrid Models:** Combine LSTM and Gradient Boosting for better temporal representation.
- **Hyperparameter Tuning:** Perform a comprehensive search for optimal settings.
- **External Factors:** Incorporate macroeconomic indicators and news sentiment for improved accuracy.