

# Model Selection Documentation

## 1. Comparison of Different Modeling Approaches Tested

For stock price prediction, various machine learning models can be tested, ranging from simple linear models to complex non-linear models like tree-based algorithms. Below, we compare the performance of the following models:

(a)

### a. Linear Regression (Baseline Model)

- **Approach:** A basic linear regression model assumes a linear relationship between the features and the target variable (stock price).
- **Pros:** Simple, interpretable, and fast to train.
- **Cons:** It assumes that the relationship between features and target is linear, which may not hold true for stock prices, which are often non-linear.

### b. Random Forest Regressor

- **Approach:** A Random Forest is an ensemble of decision trees that uses multiple decision trees to predict the output. It is more flexible than linear regression because it can model non-linear relationships.
- **Pros:** Handles non-linearity well, reduces overfitting compared to a single decision tree, and can model complex relationships between features and target.
- **Cons:** More computationally expensive, less interpretable than linear models.

### c. XGBoost (Gradient Boosting)

- **Approach:** XGBoost is an optimized implementation of gradient boosting, a technique that builds an ensemble of trees in a sequential manner. Each subsequent tree corrects the errors made by the previous one.
- **Pros:** Often performs very well for structured data, can handle non-linearities, and is robust to overfitting with proper hyperparameter tuning.
- **Cons:** Can be computationally intensive, requires careful tuning of hyperparameters to avoid overfitting.

#### d. LSTM (Long Short-Term Memory)

- **Approach:** LSTM is a type of recurrent neural network (RNN) specifically designed to handle sequences of data, such as time series. It can capture temporal dependencies, making it suitable for stock price prediction.
- **Pros:** Good at learning long-term dependencies in time series data, which is crucial for stock prices.
- **Cons:** Requires more computational resources and time for training. LSTMs can be prone to overfitting if not carefully regularized.

#### (b) Model Performance Comparison

To compare the models, we evaluate their performance using the following metrics:

- **Root Mean Squared Error (RMSE):** This metric measures the average magnitude of the error between predicted and actual values. It penalizes larger errors more heavily than smaller ones, making it a good measure for stock price prediction.

Formula for RMSE:

$$RMSE = \sqrt{\sum_{i=1}^n \frac{(\hat{y}_i - y_i)^2}{n}}$$

$\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n$  are predicted values  
 $y_1, y_2, \dots, y_n$  are observed values  
 $n$  is the number of observations

- **Mean Absolute Error (MAE):** Measures the average of the absolute errors between predicted and actual values. It is less sensitive to outliers than RMSE but still provides useful insights into prediction accuracy.
- **Directional Accuracy:** In stock price prediction, directional accuracy measures how often the model correctly predicts the direction of price change (up or down) rather than the exact price. This is particularly useful in financial applications where knowing the direction of price movement can be more valuable than the exact value.

Here's an example of model performance comparison (you would replace the data with your actual results):

Model	RMSE	MAE	Directional Accuracy (%)
Linear Regression	5.23	4.01	60%
Random Forest	4.88	3.85	65%
XGBoost	4.12	3.21	70%
LSTM	4.56	3.62	68%

## 2. Explanation of Evaluation Metrics Used for Selection

- RMSE (Root Mean Squared Error):** RMSE is a popular metric for regression tasks. It penalizes larger errors more severely than smaller ones, which is useful when we want to ensure the model does not make large mistakes in predicting stock prices. A lower RMSE indicates a better fit.
- MAE (Mean Absolute Error):** MAE provides a simple average of the absolute prediction errors. While it doesn't penalize large errors as much as RMSE, it gives a more straightforward sense of how close the predicted values are to the actual values on average.
- Directional Accuracy:** This metric is particularly important for stock price prediction because even if the model's exact prediction is slightly off, getting the direction right (whether the price will go up or down) can have significant value in practical trading strategies. This metric is crucial for assessing the model's practical usefulness in real-world applications, such as automated trading systems.

## 3. Justification for Final Model Choice

### Final Model Chosen: XGBoost

After testing various models, **XGBoost** was chosen as the final model based on the following reasons:

- Best Performance:** XGBoost showed the lowest RMSE and highest directional accuracy among all models tested. Its ability to capture complex, non-linear relationships in the data made it the most suitable for stock price prediction.

- **Versatility:** XGBoost can handle a mix of numerical and categorical data and is robust to overfitting, especially with proper regularization. Its performance is consistent across various datasets.
- **Efficiency:** XGBoost is computationally efficient, especially with large datasets, and can be trained relatively faster compared to neural network-based models like LSTM.
- **Practical Trading Value:** Given that XGBoost's directional accuracy was higher than the other models, it provides the most practical trading value, as predicting price direction is often more valuable than exact price prediction in financial markets.

#### 4. Analysis of Model Limitations and Potential Improvements with Additional Time/Data

##### Limitations of the Chosen Model (XGBoost):

1. **Lack of Temporal Context:** XGBoost is not specifically designed for time-series data and does not inherently capture temporal dependencies in the data. While lagged features and moving averages were used to mitigate this, the model could benefit from incorporating deeper time-series features.
2. **Feature Engineering:** The model heavily relies on engineered features like lagged variables and moving averages. These features are manually selected and may not capture all the complexities of the stock market data.
3. **Overfitting on Noisy Data:** Stock prices are influenced by many unpredictable factors, such as news and market sentiment, which may not be reflected in the dataset. The model might overfit to historical patterns that don't necessarily predict future trends effectively.

##### Potential Improvements:

1. **Incorporating More Features:** Adding external factors, such as market sentiment, social media mentions, or news sentiment, could significantly improve the model's predictive power. Feature selection could also be improved by using advanced techniques such as feature importance from models or autoencoders for dimensionality reduction.
2. **Advanced Time-Series Models:** Using models specifically designed for time-series data, such as **Long Short-Term Memory (LSTM)** or **ARIMA** models, could help capture temporal dependencies more effectively. This would allow the model to better understand how past prices influence future prices over time.

3. **Hyperparameter Tuning:** XGBoost performance can be further improved by performing exhaustive hyperparameter tuning (using GridSearchCV or RandomizedSearchCV) to find the optimal set of parameters for better accuracy.
4. **Ensemble Methods:** Combining multiple models, such as XGBoost with an LSTM or ARIMA model, could potentially lead to better performance by taking advantage of the strengths of each model.
5. **Increasing Data Size:** The performance of the model could improve if more data is available. More extensive datasets would allow the model to capture more patterns, seasonality, and volatility in the stock market. Additionally, using rolling windows or more granular data (e.g., hourly or minute-level data) may help.