Algorithmic Trading

With Deep Learning

Ritesh Kumar

2024SP_MS_DSP_498-DL_SEC61:  Capstone

Module 6

Status Report

Kimberly Chulis and Srabashi Basu

July 26, 2024

Table of Contents

**Abstract**

In the dynamic landscape of financial markets, deep learning techniques have revolutionized algorithmic trading strategies. This project aims to develop an advanced algorithmic trading system for BankNifty, a prominent Indian banking index. The system will leverage historical data, including Cumulative Open Interest (COI), Price, Volume, India Volatility Index, and Technical Indicators such as Moving Averages, RSI, and MACD, to predict market trends and make informed trading decisions.

The primary objective is to design a sophisticated trading system that analyzes historical data and informs trading decisions using deep learning models. Key components of the project include Feature Engineering, Backtesting, Risk Management, Performance Metrics, Scalability, and Adaptation. High-frequency data collection at 1-minute intervals will enable the system to capture short-term market movements and trends.

Data preprocessing steps will include handling missing data, feature engineering to create new features, data normalization, and splitting the dataset into training, validation, and testing sets. The deep learning models employed will include Long Short-Term Memory (LSTM) networks, Convolutional Neural Networks (CNN), and Transformer Models, each chosen for their ability to capture sequential patterns, identify complex features, and handle long-range dependencies in time-series data.

Effective feature engineering will be crucial, with techniques like recursive feature elimination (RFE) and feature importance ranking identifying the most relevant features. Technical indicators will enhance the feature set.

Risk management strategies, including stop-loss orders and position sizing, will protect capital and minimize losses. Performance metrics such as the Win-Loss Ratio and Profit Factor will provide a comprehensive assessment of the system's effectiveness.

**Project Plan**

The project plan is designed to ensure a systematic approach to development and deployment. It begins with the publication of the project charter, which clearly defines the scope, objectives, and deliverables. Setting up a live-data extraction system is crucial for gathering the required historical data, including Cumulative Open Interest (COI), Price, Volume, and various Technical Indicators.

The next phase involves conducting exploratory data analysis to understand data distributions, identify trends, and uncover any quality issues that need addressing. Preprocessing the data for modeling is critical, involving feature engineering and data normalization to prepare the dataset for effective model training. We have completed the preprocessing and data engineering.

Exploration of different deep learning techniques, such as LSTM, CNN, and Transformer models, forms the core of the project. Evaluating these models and comparing their performance will help identify the most suitable approach for our trading system. Fine-tuning hyperparameters and validating the models on a separate validation dataset ensures robustness and reliability. We are currently experimenting with different models and their respective hyperparameters.

Once the best-performing model is selected, the deployment phase will involve integrating the model into a live trading environment. Continuous monitoring and adjustments based on real-time data feedback will be essential to maintain optimal performance.

The final stages of the project involve compiling progress and results and presenting the outcomes. This structured plan aims to deliver a powerful tool for high-frequency trading, tailored to the complexities of BankNifty and potentially extendable to other indices.

## Status

**Algorithmic TradingWith Deep Learning**

| MENTORS | Kimberly Chulis and Srabashi Basu |
|---|---|
| TEAM | Pooja, Ezhilarasu, Shreenivas, and Ritesh |
| COURSE | 2024SP_MS_DSP_498-DL_SEC61 |
| PROJECT START DATE | 18-Jun-24 |

**MILESTONE 1:** Project Charter Sign-off
**MILESTONE 2:** Exploratory Data Analysis
**MILESTONE 3:** Modelling
**MILESTONE 4:**
**MILESTONE 5:** Final Submission

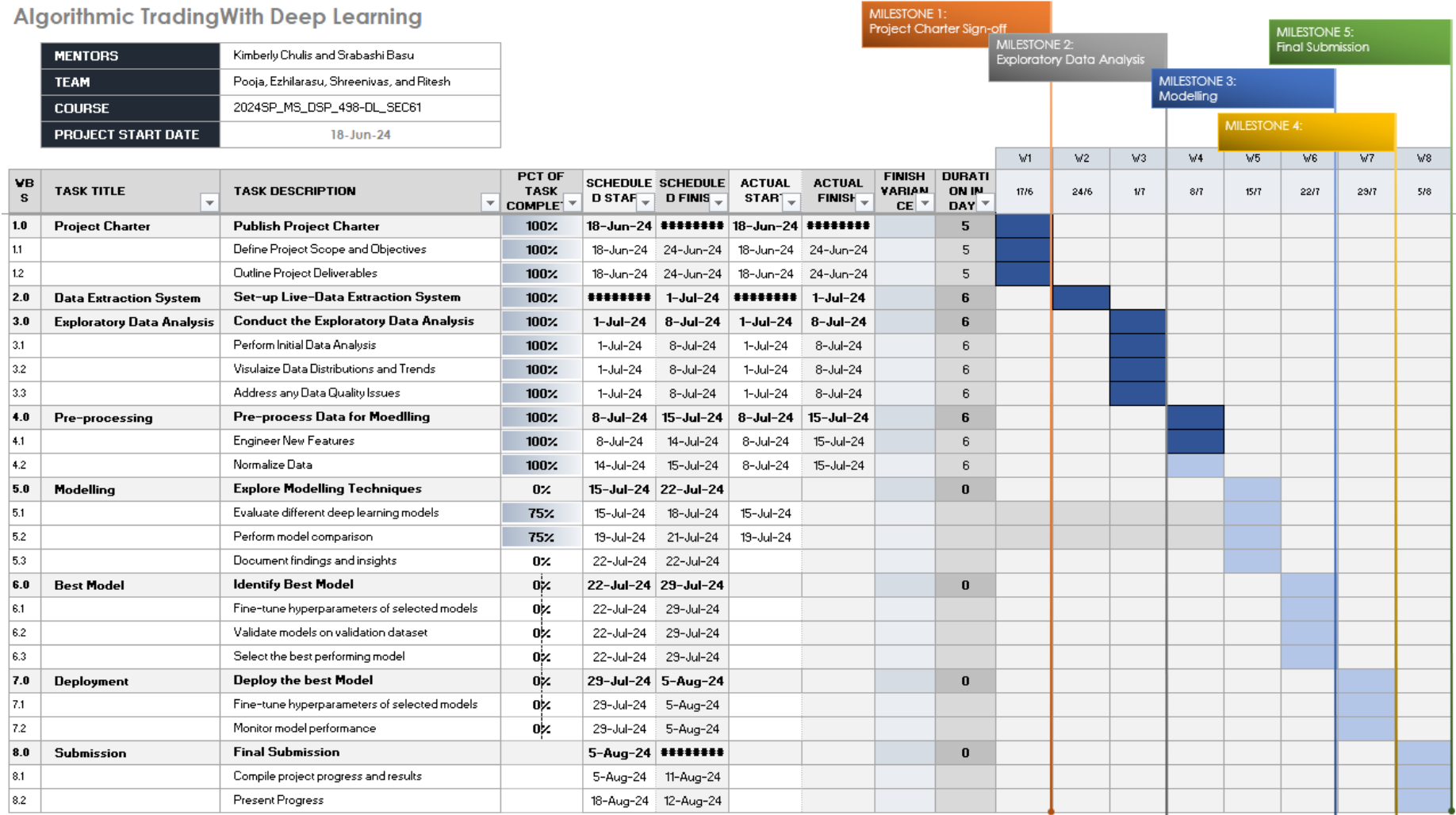| WBS | TASK TITLE | TASK DESCRIPTION | PCT OF TASK COMPLETE | SCHEDULED START | SCHEDULED FINISH | ACTUAL START | ACTUAL FINISH | FINISH VARIANCE | DURATION IN DAYS |
|---|---|---|---|---|---|---|---|---|---|
| 1.0 | Project Charter | Publish Project Charter | 100% | 18-Jun-24 | ######## | 18-Jun-24 | ######## | | 5 |
| 1.1 | | Define Project Scope and Objectives | 100% | 18-Jun-24 | 24-Jun-24 | 18-Jun-24 | 24-Jun-24 | | 5 |
| 1.2 | | Outline Project Deliverables | 100% | 18-Jun-24 | 24-Jun-24 | 18-Jun-24 | 24-Jun-24 | | 5 |
| 2.0 | Data Extraction System | Set-up Live-Data Extraction System | 100% | ######## | 1-Jul-24 | ######## | 1-Jul-24 | | 6 |
| 3.0 | Exploratory Data Analysis | Conduct the Exploratory Data Analysis | 100% | 1-Jul-24 | 8-Jul-24 | 1-Jul-24 | 8-Jul-24 | | 6 |
| 3.1 | | Perform Initial Data Analysis | 100% | 1-Jul-24 | 8-Jul-24 | 1-Jul-24 | 8-Jul-24 | | 6 |
| 3.2 | | Visulaize Data Distributions and Trends | 100% | 1-Jul-24 | 8-Jul-24 | 1-Jul-24 | 8-Jul-24 | | 6 |
| 3.3 | | Address any Data Quality Issues | 100% | 1-Jul-24 | 8-Jul-24 | 1-Jul-24 | 8-Jul-24 | | 6 |
| 4.0 | Pre-processing | Pre-process Data for Moedlling | 100% | 8-Jul-24 | 15-Jul-24 | 8-Jul-24 | 15-Jul-24 | | 6 |
| 4.1 | | Engineer New Features | 100% | 8-Jul-24 | 14-Jul-24 | 8-Jul-24 | 15-Jul-24 | | 6 |
| 4.2 | | Normalize Data | 100% | 14-Jul-24 | 15-Jul-24 | 8-Jul-24 | 15-Jul-24 | | 6 |
| 5.0 | Modelling | Explore Modelling Techniques | 0% | 15-Jul-24 | 22-Jul-24 | | | | 0 |
| 5.1 | | Evaluate different deep learning models | 75% | 15-Jul-24 | 18-Jul-24 | 15-Jul-24 | | | |
| 5.2 | | Perform model comparison | 75% | 19-Jul-24 | 21-Jul-24 | 19-Jul-24 | | | |
| 5.3 | | Document findings and insights | 0% | 22-Jul-24 | 22-Jul-24 | | | | |
| 6.0 | Best Model | Identify Best Model | 0% | 22-Jul-24 | 29-Jul-24 | | | | 0 |
| 6.1 | | Fine-tune hyperparameters of selected models | 0% | 22-Jul-24 | 29-Jul-24 | | | | |
| 6.2 | | Validate models on validation dataset | 0% | 22-Jul-24 | 29-Jul-24 | | | | |
| 6.3 | | Select the best performing model | 0% | 22-Jul-24 | 29-Jul-24 | | | | |
| 7.0 | Deployment | Deploy the best Model | 0% | 29-Jul-24 | 5-Aug-24 | | | | 0 |
| 7.1 | | Fine-tune hyperparameters of selected models | 0% | 29-Jul-24 | 5-Aug-24 | | | | |
| 7.2 | | Monitor model performance | 0% | 29-Jul-24 | 5-Aug-24 | | | | |
| 8.0 | Submission | Final Submission | | 5-Aug-24 | ######## | | | | 0 |
| 8.1 | | Compile project progress and results | | 5-Aug-24 | 11-Aug-24 | | | | |
| 8.2 | | Present Progress | | 18-Aug-24 | 12-Aug-24 | | | | |

Table 1

**Status Report**

As of July 26, 2024, the project has achieved substantial milestones and progressed significantly through several critical phases. The following provides a detailed overview of the status of various tasks and milestones:

1. Project Charter

   - Publish Project Charter: 100% Complete (18-Jun-24 to 24-Jun-24)

   - Define Project Scope and Objectives: 100% Complete (18-Jun-24 to 24-Jun-24)

   - Outline Project Deliverables: 100% Complete (18-Jun-24 to 24-Jun-24)

2. Data Extraction System

   - Set-up Live-Data Extraction System: 100% Complete (24-Jun-24 to 1-Jul-24)

3. Exploratory Data Analysis

   - Conduct the Exploratory Data Analysis: 100% Complete (1-Jul-24 to 8-Jul-24)

   - Perform Initial Data Analysis: 100% Complete (1-Jul-24 to 8-Jul-24)

   - Visualize Data Distributions and Trends: 100% Complete (1-Jul-24 to 8-Jul-24)

   - Address any Data Quality Issues: 100% Complete (1-Jul-24 to 8-Jul-24)

4. Pre-processing

   - Pre-process Data for Modelling: 100% Complete (8-Jul-24 to 15-Jul-24)

   - Engineer New Features: 100% Complete (14-Jul-24 to 15-Jul-24)

   - Normalize Data: 100% Complete (15-Jul-24 to 15-Jul-24)

5. Modelling

   - Explore Modelling Techniques: 75% Complete (15-Jul-24 to 22-Jul-24)

   - Evaluate different deep learning models: 75% Complete (15-Jul-24 to 21-Jul-24)

   - Perform model comparison: 75% Complete (Scheduled: 21-Jul-24 to 22-Jul-24)

   - Document findings and insights: 0% Complete (Scheduled: 22-Jul-24 to 24-Jul-24)

6. Best Model

   - Identify Best Model: 0% Complete (Scheduled: 22-Jul-24 to 29-Jul-24)

   - Fine-tune hyperparameters of selected models: 0% Complete (Scheduled: 22-Jul-24 to 29-Jul-24)

   - Validate models on validation dataset: 0% Complete (Scheduled: 22-Jul-24)

   - Select the best performing model: 0% Complete (Scheduled: 22-Jul-24 to 29-Jul-24)

7. Deployment

   - Deploy the Best Model: 0% Complete (Scheduled: 29-Jul-24 to 5-Aug-24)

   - Fine-tune hyperparameters of selected models: 0% Complete (Scheduled: 29-Jul-24 to 5-Aug-24)

   - Monitor model performance: 0% Complete (Scheduled: 5-Aug-24 to 12-Aug-24)

8. Submission

   - Final Submission: 0% Complete (Scheduled: 5-Aug-24 to 12-Aug-24)

   - Compile project progress and results: 0% Complete (Scheduled: 5-Aug-24 to 12-Aug-24)

   - Present Progress: 0% Complete (Scheduled: 12-Aug-24)

Risks and Issues

1. Internet Outage: On 3rd July, an internet outage occurred between 11:22 AM and 1:07 PM, resulting in a gap in our data collection efforts. To prevent future disruptions, starting from 8th July, we have migrated our data collection process to Sagemaker Studio on AWS, ensuring high availability and reducing the risk of data loss due to connectivity issues.

2. Model Performance Issues: The current RMSE values are above the ideal benchmark of 11, and the direction agreement is below the desired 60%. We are addressing this by:

- Hyperparameter Tuning: Experimenting with various hyperparameter configurations.

- Feature Engineering: Refining existing features and creating new ones.

- Model Ensemble: Combining predictions from multiple models to improve accuracy.

- Data Augmentation: Enhancing the dataset by generating synthetic data or augmenting existing data.

At this stage, no additional risks or issues are anticipated. The project is progressing as planned, and the team is prepared to address any potential challenges that may arise.

**Data Description**

In the rapidly changing landscape of financial markets, deep learning techniques have revolutionized algorithmic trading strategies. This dataset is specifically curated to aid in the development of a sophisticated algorithmic trading system for BankNifty. It aims to harness historical data to accurately predict market trends and make informed trading decisions by utilizing high-frequency data collection and advanced technical indicators. The comprehensive data provided will support the implementation of cutting-edge deep learning models, ultimately enhancing trading performance and strategy.

Date and Time Variables:

- date: Represents the date of the observation.

- time: Represents the specific time of the observation. Together, 'date' and 'time' provide a precise timestamp for each data point.

Price and Volume Variables:

- open: The price of the asset at the start of the trading interval.

- high: The highest price of the asset during the trading interval.

- low: The lowest price of the asset during the trading interval.

- close: The price of the asset at the end of the trading interval.

- vix: The Volatility Index, which measures market expectation of near-term volatility.

- ce_vol: The call option volume, representing the number of call options traded.

- pe_vol: The put option volume, representing the number of put options traded.

- total_vol: The total trading volume, combining both call and put options.

Open Interest Variables:

- ce_oi: Call option open interest, indicating the total number of call options outstanding.

- pe_oi: Put option open interest, indicating the total number of put options outstanding.

- ce_oi_chg: Change in call option open interest.

- pe_oi_chg: Change in put option open interest.

- tot_oi_chg: Total change in open interest.

Put/Call Ratio Variables:

- pcr_vol: Put/Call ratio based on volume.

- pcr_oi: Put/Call ratio based on open interest.

Moving Averages:

- SMA_5: 5-period Simple Moving Average.

- SMA_10: 10-period Simple Moving Average.

- SMA_20: 20-period Simple Moving Average.

- EMA_12: 12-period Exponential Moving Average.

- EMA_20: 20-period Exponential Moving Average.

- EMA_26: 26-period Exponential Moving Average.

Technical Indicators:

- RSI: Relative Strength Index, a momentum oscillator that measures the speed and change
  of price movements.

- MACD: Moving Average Convergence Divergence, calculated as the difference between the 12-period and 26-period EMAs.

- Signal_Line: 9-period EMA of the MACD.

- MACD_Histogram: The difference between MACD and Signal Line.

- Middle_Band: Middle band of Bollinger Bands, usually a 20-period SMA.

- Upper_Band: Upper band of Bollinger Bands.

- Lower_Band: Lower band of Bollinger Bands.

- %K: Stochastic Oscillator %K, indicating the current price relative to the high and low range over a set period.

- %D: Stochastic Oscillator %D, the 3-period SMA of %K.

- ATR: Average True Range, measuring market volatility.

- OBV: On-Balance Volume, using volume flow to predict changes in stock price.

- VWAP: Volume Weighted Average Price.

- ROC_10: Rate of Change over 10 periods.

- ROC_20: Rate of Change over 20 periods.

- CCI: Commodity Channel Index, measuring the variation of an asset's price from its statistical mean.

- plus_DM: Positive Directional Movement.

- minus_DM: Negative Directional Movement.

- TR: True Range, the greatest of the following: current high - current low, current high - previous close, or current low - previous close.

- plus_DI: Positive Directional Indicator.

- minus_DI: Negative Directional Indicator.

- DX: Directional Movement Index, indicating trend strength.

- ADX: Average Directional Index, smoothing the DX over a specified period.

Pivot Points and Support/Resistance Levels:

- exp_day: Expiration day, indicating the day options expire.

- pivot: Pivot point, calculated as the average of the high, low, and closing prices from the previous trading period.

- s1, s2, s3: First, second, and third support levels derived from the pivot point.

- r1, r2, r3: First, second, and third resistance levels derived from the pivot point.

Target Variable:

- nxt_move: The variable of interest, representing the predicted next movement in price or the number of points to the next close.

Data Summary:

- Minimum datetime: 2024-06-24 10:00:53

- Maximum datetime: 2024-07-02 15:31:04

- Total rows: 2561

This dataset provides an extensive and detailed view of market behavior, encompassing a wide range of technical indicators and price/volume metrics. It includes essential variables such as open, high, low, and close prices, along with trading volumes and open interest data for both call and put options. Advanced technical indicators like moving averages, RSI, MACD, Bollinger Bands, and stochastic oscillators offer insights into market trends, momentum, and volatility.

Additionally, the dataset features pivotal points, support and resistance levels, and volatility measures like the India Volatility Index (VIX) and Average True Range (ATR), which are crucial for understanding market dynamics. The inclusion of both simple and exponential

moving averages (SMA and EMA) over various periods allows for the analysis of short-term and long-term trends.

Metrics such as the Put/Call ratio based on volume and open interest further enhance the dataset's ability to gauge market sentiment. The dataset also incorporates crucial elements like the On-Balance Volume (OBV) and Volume Weighted Average Price (VWAP), providing a robust foundation for assessing the flow of funds and average trading price.

By offering a granular view of market data at 1-minute intervals, this dataset is exceptionally suited for high-frequency trading analysis. It supports the development and validation of advanced predictive models, including LSTM, CNN, and Transformer models, ensuring a comprehensive approach to algorithmic trading analysis. This dataset is instrumental in building and fine-tuning sophisticated trading strategies, ultimately contributing to more informed and effective decision-making in the financial markets.

**Target Variable**

The descriptive statistics of the target variable `nxt_move` are as follows:

- Count: There are 2,378 observations for the `nxt_move` variable.

- Mean: The average value of `nxt_move` is approximately 0.224.

- Standard Deviation (std): The standard deviation is 22.48, indicating a significant spread around the mean.

- Minimum (min): The minimum value is -151.55, suggesting a substantial negative movement.

- 25th Percentile (25%): The 25th percentile is -11.04, meaning 25% of the data points are less than -11.04.

- Median (50%): The median value is 0.00, indicating that half of the observations are below and half above 0.

- 75th Percentile (75%): The 75th percentile is 11.05, meaning 25% of the data points are greater than 11.05.

- Maximum (max): The maximum value is 278.55, suggesting a substantial positive movement.
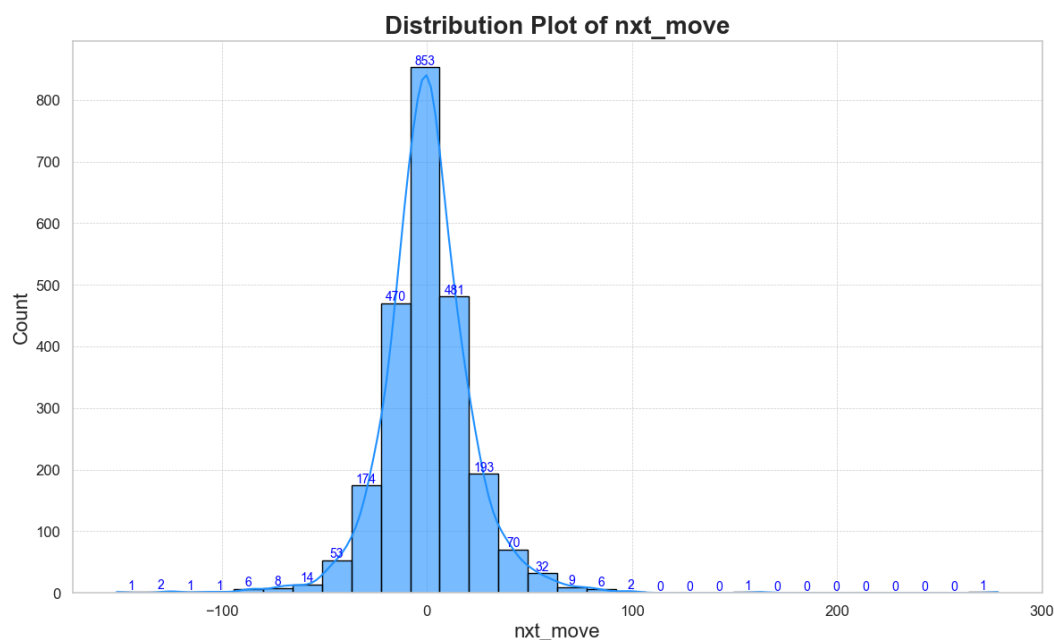
The target variable `nxt_move` has a mean close to zero with a high standard deviation, indicating a wide range of values. The distribution is centered around zero, as shown by the median. The data contains both significant negative and positive values, highlighting the potential for considerable market movements in either direction. This variability and spread in `nxt_move` are critical for developing predictive models, as they must account for the broad range of possible outcomes.



Plot 1

- Median and Quartiles: The boxplot shows that the median (50th percentile) of `nxt_move` is around 0, indicating that half of the observations are below and half above 0. The

interquartile range (IQR), represented by the box, is relatively narrow, highlighting that 50% of the data points lie between approximately -11 and +11.



Plot 2

- Shape and Central Tendency: The distribution plot is bell-shaped, indicating a normal-like distribution centered around 0. The highest count of observations (853) is near the mean value.

The `nxt_move` variable has a distribution centered around zero with a high concentration of data points within a narrow range. However, it also exhibits substantial outliers and a slight skewness towards the positive side. This variability and presence of outliers are critical for developing robust predictive models, as they must effectively handle and predict these extreme market movements.

**Exploratory Data Analysis**

The correlation matrix reveals that `nxt_move` has weak correlations with most variables, indicating that no single factor strongly predicts the next market movement. The variable

`pcr_oi` shows the strongest negative correlation, suggesting that the put/call ratio based on

open interest is inversely related to `nxt_move`. Other variables, including price metrics

(`open`, `high`, `low`, `close`) and moving averages (SMA, EMA), display very weak

correlations. This implies that while these factors are essential for understanding market

trends, predicting `nxt_move` effectively requires a combination of multiple indicators rather

than relying on individual ones.



Plot 3

Positive Correlations:

- Put/Call Ratios: `pcr_oi` (0.031451) and `pcr_vol` (0.006902) have slight positive correlations, suggesting that higher put/call ratios might be associated with increases in `nxt_move`.

- Volatility Index: `vix` (0.023932) has a positive correlation, indicating that higher market volatility is associated with increases in `nxt_move`.

- Open Interest Changes: `pe_oi_chg` (0.023425) and `tot_oi_chg` (0.012114) show slight positive correlations, suggesting that changes in open interest, particularly in puts, might be linked to increases in `nxt_move`.

- Expiration Day: `exp_day` (0.012748) shows a slight positive correlation, indicating that the day of options expiration can have a slight impact on `nxt_move`.

Minimal Correlations:

- Volume Metrics: Variables such as `ce_vol`, `pe_vol`, and `total_vol` have very minimal negative correlations, suggesting that trading volume has little to no direct correlation with `nxt_move`.

- Directional Indicators and ADX: `minus_DI` and `ce_oi` have minimal positive correlations, while `ce_oi_chg` has a minimal negative correlation, indicating that these metrics have a negligible impact on `nxt_move`.

Overall, the data indicates that several price, volume, and technical indicators have varying degrees of correlation with `nxt_move`. Most correlations are relatively weak, suggesting a complex relationship between these indicators and market movements. The negative correlations generally suggest that lower values in these indicators are associated with decreases in `nxt_move`, while the few positive correlations indicate an association with increases.

Ten Strongest Correlations

```
plus_DM      -0.083949
minus_DM      0.058600
s3           -0.053277
s1           -0.051322
s2           -0.051155
low          -0.051004
close        -0.050778
VWAP         -0.050679
open         -0.049579
high         -0.049290
```

In the exploratory data analysis (EDA) phase, we examined a comprehensive dataset. Currently spanning seven trading days, this dataset includes detailed historical data on Cumulative Open Interest (COI), price metrics, trading volumes, the India Volatility Index (VIX), and numerous technical indicators such as moving averages, RSI, and MACD.

Despite the dataset's richness, our correlation analysis revealed that no single factor strongly predicts the target variable, `nxt_move`. The strongest correlations observed were still relatively low, indicating that predicting market movements effectively will require a multifactorial approach. Although the correlations are low, this can still be leveraged effectively using advanced modeling techniques.

**Accomplishments**

Data Collection: We have successfully collected over 7500 data points so far. This substantial amount of data provides a solid foundation for our analysis and modeling efforts. By the time we conclude modeling, we expect to have collected an additional ~2,500 data points, bringing the total to approximately 10,000 data points. This expanded dataset will enhance our ability to explore and develop robust modeling techniques.

Exploratory Data Analysis:

We have completed the initial data analysis phase, which involved visualizing data distributions and trends. This analysis has provided valuable insights into the dataset's characteristics and highlighted any data quality issues that needed addressing. Our team has

successfully addressed these issues, ensuring the data is clean and ready for further processing and modeling.

Pre-processing and Feature Engineering:

We have completed this phase of the project, which involved the creation of technical features essential for our modeling efforts. The pre-processing steps included data normalization, handling missing values, and transforming raw data into features that better represent the underlying patterns. Time of the day has been normalized, using sine and cosine conversion for time, and all other features have been normalized using standard scalar. This comprehensive approach has enabled us to build a robust dataset, primed for effective model training.

Model Development:

We are currently in the model development phase, experimenting with various deep learning techniques such as LSTM, CNN, and Transformer models. Our focus is on evaluating the performance of these models, fine-tuning hyperparameters, and validating their effectiveness using a separate validation dataset. This rigorous approach aims to identify the most suitable model for our high-frequency trading system.

| | CNN | | | LSTM | | | Transformers | | |
|---|---|---|---|---|---|---|---|---|---|
| seq_len | 5 | 10 | 30 | 5 | 10 | 30 | 5 | 10 | 30 |
| MSE | 430.19 | 349.08 | 478.82 | 372.30 | 359.22 | 362.05 | 461.26 | 413.18 | 367.22 |
| RMSE | 20.74 | 18.68 | 21.88 | 19.30 | 18.95 | 19.03 | 21.48 | 20.33 | 19.16 |
| dir_agr% | 50.27 | 52.34 | 49.27 | 51.36 | 47.93 | 49.85 | 50.00 | 51.24 | 51.9 |

Table 2

Modelling Progress Highlights

1.  Performance Metrics:

- Models have been evaluated based on Mean Squared Error (MSE) and Root Mean Squared Error (RMSE) across different sequence lengths.

- The goal is to achieve an RMSE below 11 for optimal performance, which remains a challenge.

2. Model Types:

- Various deep learning techniques are being explored, including CNNs, LSTMs, and Transformers.

- Each model type has been tested with sequence lengths of 5, 10, and 30 to assess performance variations.

3. Direction Prediction Accuracy:

- The direction in agreement metric indicates how accurately the models predict the direction of the next move.

- While some models show promising accuracy, reaching the 60% target remains a priority.

4. Key Observations:

- CNNs show varying performance with different sequence lengths, with the sequence length of 10 yielding the best results.

- LSTMs have shown consistent performance across sequence lengths, with the best accuracy at a sequence length of 5.

- Transformer models exhibit potential, with improvements needed to meet the desired benchmarks.

Next Steps:

- Model Deployment: Once the best-performing model is selected, we will proceed with the deployment phase. This involves integrating the model into a live trading environment

and ensuring seamless operation. Continuous monitoring and real-time data feedback will be crucial for maintaining optimal performance and making necessary adjustments.

- Final Reporting and Presentation: The final stages of the project involve compiling our progress and results into comprehensive reports. We will present our findings, highlighting the performance of our models and their potential applications in high-frequency trading. This structured plan aims to deliver a powerful tool for high-frequency trading, tailored to the complexities of BankNifty and potentially extendable to other indices.

**Risks and Issues**

Internet Outage: On July 3rd, an internet outage from 11:22 AM to 1:07 PM caused a disruption in our data collection process. This incident highlighted the need for a more reliable data collection system to ensure continuous and uninterrupted data flow. To mitigate this risk, we migrated our data collection process to Sagemaker Studio on AWS starting July 8th. This platform provides a robust and reliable environment for our data collection needs, ensuring high availability and significantly reducing the risk of data loss due to connectivity issues.

Model Performance Issues: Achieving the desired performance metrics for our models has proven to be challenging. The current Root Mean Squared Error (RMSE) values are above the ideal benchmark of 11, and the direction agreement metric is below the desired 60%. To address these issues, we are employing several strategies including hyperparameter tuning, refining existing features, and creating new ones. Additionally, we are experimenting with ensemble methods to combine predictions from multiple models, and employing data augmentation techniques to enhance our dataset, which can improve model robustness and accuracy.

Other Potential Risks: At this stage, no additional risks or issues are anticipated. The project is progressing as planned, and the team is well-prepared to address any potential challenges that may arise. Regular risk assessments and proactive mitigation strategies are in place to ensure the project's continued success.

**Conclusion**

In conclusion, the progress made in this phase of the project lays a strong foundation for the subsequent stages. Our successful data collection, amassing over 7,500 data points with an anticipated total of approximately 10,000, ensures a robust dataset for our modeling efforts. The initial exploratory data analysis has been pivotal, providing critical insights into data distributions and trends, and enabling us to address data quality issues effectively.

The completion of pre-processing and feature engineering has transformed raw data into valuable features, priming it for model training. Our current focus on model development, experimenting with LSTM, CNN, and Transformer models, is crucial for identifying the most suitable approach for our high-frequency trading system. Fine-tuning hyperparameters and validating these models ensure their robustness and reliability.

Looking ahead, the deployment phase will integrate the best-performing model into a live trading environment, with continuous monitoring and real-time adjustments to maintain optimal performance. Finally, compiling our progress and presenting our findings will highlight the efficacy of our models and their potential applications. This phase sets the stage for delivering a powerful, data-driven tool for high-frequency trading, tailored to the complexities of BankNifty and beyond.

**References**

1.  Manveer Kaur Mangat, Erhard Reschenhofer, Thomas Stark, Christian Zwatz. High-Frequency Trading with Machine Learning Algorithms and Limit Order Book Data. Data Science in Finance and Economics, 2022, 2(4): 437-463. ([Link](#))

2.  Arévalo, Andrés & Nino, Jaime & Hernandez, German & Sandoval, Javier. (2016). High-Frequency Trading Strategy Based on Deep Neural Networks. 9773. 424-436. 10.1007/978-3-319-42297-8_40. ([Link](#))