

---

# Machine Learning Project Proposal

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Introduction

This report will describe our proposal for the main project for the course 80245013 - *Machine Learning* held at Tsinghua University. The report is divided into Background, Definition, Related Work and the Proposed Method.

## 2 Background

Financial markets are a cornerstone of the global economy, influencing corporate valuations and individual investments alike. With the rise of advanced machine learning algorithms, the challenge of accurately forecasting market behavior has become increasingly relevant. Still, financial data is notoriously difficult to model due to its unpredictable patterns and sudden shifts that traditional methods struggle to capture [6].

In response, Jane Street, a market leader in automated trading[4], has launched a Kaggle competition to develop more accurate and robust prediction models [3]. Solving this challenge could have significant real-world impact by improving trading strategies, leading to more informed decisions and increased profitability. Additionally, the insights gained from applying machine learning to noisy, non-stationary data could enhance the field of time series forecasting, with benefits extending beyond the field of economics and finance to areas such as healthcare, astrology and climate science [5].

## 3 Definition

As our objective is to predict developments in the stock market, there isn't any mathematical definition of our problem. Still, to score highly in the competition we have to create a model which minimizes error. The kaggle competition uses the following error calculation to estimate the scoring in the competition.

$$R^2 = 1 - \frac{\sum w_i (y_i - \hat{y}_i)^2}{\sum w_i y_i^2}$$

Figure 1: The formula used to decide the model placements in the competition

21

22 Here is an explanation of the symbols used:

- 23 •  $w_i$ : A weight applied to each observation  $i$  in the dataset, allowing certain observations  
24 to have more or less influence on the error calculation, possibly based on relevance or  
25 importance.
- 26 •  $y_i$ : The actual observed value for observation  $i$  in the dataset, representing real stock market  
27 values that the model aims to predict.

28 •  $\hat{y}_i$ : The predicted value generated by the model for observation  $i$ .

29 This error calculation evaluates the weighted squared differences between actual and predicted values,  
30 adjusting the impact of errors by  $w_i$  to reflect competition scoring criteria. As of October 30, 2024,  
31 the 10 best scores in the competition are in the range 0.0068-0.0050.

## 32 4 Related Work

### 33 Hybrid Bidirectional LSTMs (H.BLSTMs)

- 34 • **Advantages:** Captures long-term dependencies, adapt to any changing market conditions.
- 35 • **Disadvantages:** Complexity and time computation, dependence on data quantity and quality.

### 36 Extended Kalman filter Non-linear Autoregressive Neural Network (EKF-NAR)

- 37 • **Advantages:** Computes with a big improved accuracy and can also handle complex patterns.
- 38 • **Disadvantages:** Can barely handle very complex models with some linearization error.

## 39 5 Proposed Method

40 Our project will implement Jamba, a novel Mamba-Transformer hybrid machine learning model. We  
41 hope that using Jamba we can leverage the strengths of transformer architectures and Mamba, which  
42 is optimized for high-dimensional and sequential data. Our choice of Jamba is driven by its ability to  
43 capture complex temporal dependencies and its adaptability to high-volume, real-time data, making  
44 it a suitable candidate for predicting stock market behavior in this competition. Furthermore it is a  
45 newly developed model, only being released this year, in 2024.

### 46 5.1 Dataset Selection

47 We will use the provided Kaggle dataset, containing real-world data derived from Jane Streets  
48 production systems, as the models primary data source. We anticipate that this dataset will provide  
49 sufficient information for training our model, and therefore don't believe we will need to explore  
50 supplementary datasets during the project.

### 51 5.2 Baseline Approaches

52 In the Kaggle competition, our model's performance will be evaluated against other participants using  
53 the scoring formula outlined in section 3. This will enable us to gauge our model's effectiveness  
54 relative to the alternative solutions employed by other teams. Consequently, our focus will be on  
55 optimizing our model to achieve the highest possible score, rather than conducting comparative  
56 analyses against other models independently.

### 57 5.3 Implementation

58 Our implementation of the Jamba model on the Kaggle dataset will follow these steps:

- 59 1. **Data Preprocessing:** Raw data will be cleaned, scaled, and organized to remove noise and  
60 manage missing values. For time-series data, we will create feature windows capturing  
61 recent past values as input to the model.
- 62 2. **Model Training:** Jamba will be trained on this data, with the training process involving  
63 cross-validation to optimize hyperparameters such as learning rate, sequence length, and  
64 transformer depth.
- 65 3. **Evaluation Metrics:** We will use the competition's scoring metric,  $R^2$ , as the primary  
66 metric, but may also use MAE (Mean Absolute Error) and MSE (Mean Squared Error) for  
67 additional insights.

68 Our approach may evolve based on preliminary results, but this outline provides a structured plan for  
69 the method.

70 **References**

- 71 [1] Science Direct. *Extended Kalman filter Non-linear Autoregressive Neural Network*. URL: <https://www.sciencedirect.com/science/article/abs/pii/S0957417423003809>. (ac-  
72 //www.sciencedirect.com/science/article/abs/pii/S0957417423003809. (ac-  
73 cessed: 30.10.2023).
- 74 [2] Science Direct. *Hybrid Bidirectional-LSTM (H.BLSTM) model*. URL: [https://maxiomwealth.com/askguru/2024/04/16/why-is-the-stock-market-so-  
75 maxiomwealth.com/askguru/2024/04/16/why-is-the-stock-market-so-  
76 difficult-to-predict/](https://maxiomwealth.com/askguru/2024/04/16/why-is-the-stock-market-so-difficult-to-predict/). (accessed: 30.10.2023).
- 77 [3] Jane Street Group. *Jane Street Real-Time Market Data Forecasting*. URL: [https://www.kaggle.com/competitions/jane-street-real-time-market-data-forecasting/  
78 kaggle.com/competitions/jane-street-real-time-market-data-forecasting/  
79 team](https://www.kaggle.com/competitions/jane-street-real-time-market-data-forecasting/team). (accessed: 23.10.2023).
- 80 [4] Jane Street Group. *Who We Are*. URL: <https://www.janestreet.com/who-we-are/>.  
81 (accessed: 23.10.2023).
- 82 [5] Analytics Steps. *5 Applications of Time Series Analysis*. URL: <https://www.analyticssteps.com/blogs/5-applications-time-series-analysis>. (accessed:  
83 analyticssteps.com/blogs/5-applications-time-series-analysis. (accessed:  
84 23.10.2023).
- 85 [6] MAXIOM Wealth. *Why is the stock market so difficult to predict?* URL: [https://maxiomwealth.com/askguru/2024/04/16/why-is-the-stock-market-so-  
86 maxiomwealth.com/askguru/2024/04/16/why-is-the-stock-market-so-  
87 difficult-to-predict/](https://maxiomwealth.com/askguru/2024/04/16/why-is-the-stock-market-so-difficult-to-predict/). (accessed: 23.10.2023).