

# proposal

April 11, 2022

## 1 Domain Background

Investment firms, hedge funds and even individuals have been using financial models to better understand market behavior and make profitable investments and trades. A wealth of information is available in the form of historical stock prices and company performance data, suitable for machine learning algorithms to process. Accurate prediction of stock market returns is a very challenging task due to volatile and non-linear nature of the financial stock markets. Since a lot of investors use AI models to work with stocks it's become very important to predict price for making decision about future investment or portfolio rebalancing. For those who use Reinforcement Learning it's also quite important since agent should know when to hold or sell/buy new assets to expect better reward for the action in future.

Academic research prove that stock closing price can be predicted with ML algorithm. According to this research prediction is quite accurate for some stocks. Please check links of this research below. Even though that algorithms prove to be working we still should keep in mind, that they are unable to predict future and sometime unexpected events might happen. For example natural disasters, issue in supply chain, etc. can lead to stock price go down unexpectedly. Hence it's good to use different sources to take a good decision, because ML algorithm in most cases can find patterns in stock behavior rather than predict real stock price.

Common algorithms used to predict stock:

- RNN
- LSTM
- Random Forest
- Transformers

### 1.0.1 Reserches

- [Stock Closing Price Prediction using Machine Learning Techniques](#)
- [A Data Organization Method for LSTM and Transformer When Predicting Chinese Banking Stock Prices](#)

## 2 Problem Statement

Build stock price predictor which will use historical data to predict stock price for given date. It's quite clear for me that it's impossible to predict future stock movement, because it depends on a lot of different factors. For example:

- Social activity

- Company performance
- Natural disasters
- Etc.

Some of these factors can be predicted from experience but some of them are still hard to predict, like natural disasters. Hence the main goal of this project is to build a predictor which tries to find patterns in the past stock movement and help with the decision-making process together with some other information.

Ideally there should be a lot of such types of predictors which use different types of information, for example market news or satellite images, etc. Which combined in an ensemble model for better prediction, but this requires a lot of time and effort to be developed. Hence in this project we try to reduce scope and build a small part of a bigger system and did small research if it is worth to make this type of predictions at all.

### 3 Solution Statement

Here we have a regression problem since we try to predict stock price. According to different researches such ML architectures as RNN and Transformers perform better than for example Random Forest (XGBoost). It looks logical since stock price is a time series and RNN and Transformers work good with this type of data. For this type of project Transformers will be used since RNN networks have some following issues comparing to Transformers architectures.

RNN architecture issues: - parallel computing is difficult to implement; - for long sequences of data there is loss of information; - problem with vanishing gradient.

To test the idea 2 random stocks from well-known companies (for example: Apple, Google, Facebook etc.) which are part of S&P 500 index will be selected.

To solve this problem different technical indicators will be used.

Technical indicators:

- Simple moving average (SMA)
- Exponential moving average (EMA)
- Moving Average Convergence Divergence (MACD)
- Bollinger band
- Momentum
- Trading volume

#### 3.0.1 Technical indicator calculation

EMA and MA will be calculated for 3, 8 and 21 days

### Simple moving average

$$SMA = \frac{A_1 + A_2 + \dots + A_n}{n} \quad (1)$$

$$\textbf{where:} \quad (2)$$

$$A = \text{Average in period } n \quad (3)$$

$$n = \text{Number of time periods} \quad (4)$$

$$(5)$$

### Exponential moving average

$$EMA_t = \left[ V_t \times \left( \frac{s}{1+d} \right) \right] + EMA_y \times \left[ 1 - \left( \frac{s}{1+d} \right) \right] \quad (6)$$

$$\textbf{where:} \quad (7)$$

$$EMA_t = \text{EMA today} \quad (8)$$

$$V_t = \text{Value today} \quad (9)$$

$$EMA_y = \text{EMA yesterday} \quad (10)$$

$$s = \text{Smoothing} \quad (11)$$

$$d = \text{Number of days} \quad (12)$$

$$(13)$$

### Moving Average Convergence Divergence

$$BOLU = MA(TP, n) + m * \sigma[TP, n] \quad (14)$$

$$BOLD = MA(TP, n) - m * \sigma[TP, n] \quad (15)$$

$$\textbf{where:} \quad (16)$$

$$BOLU = \text{Upper Bollinger Band} \quad (17)$$

$$BOLD = \text{Lower Bollinger Band} \quad (18)$$

$$MA = \text{Moving average} \quad (19)$$

$$TP \text{ (typical price)} = (\text{High} + \text{Low} + \text{Close}) \div 3 \quad (20)$$

$$n = \text{Number of days in smoothing period} \quad (21)$$

$$m = \text{Number of standard deviations} \quad (22)$$

$$\sigma[TP, n] = \text{Standard Deviation over last } n \text{ periods of TP} \quad (23)$$

$$(24)$$

### Momentum

$$\text{Momentum} = V - Vx \quad (25)$$

$$\textbf{where:} \quad (26)$$

$$V = \text{Latest price} \quad (27)$$

$$Vx = \text{Closing price} \quad (28)$$

$$x = \text{Number of days ago} \quad (29)$$

## 4 Datasets and Inputs

To train model historical data for last 5 years will be used if it's available.

There are several open sources for historical stock price data which you are free to use:

- [Yahoo Finance API Specification](#)
- [Polygon Financial Market Data Platform](#)
- [TradingView](#)
- [Quandl](#).

Actual data will be loaded from Yahoo Finance. And has following structure:

- Open
- High
- Low
- Close
- Adj Close
- Volume
- company\_name
- Date

For feature engineering following data will be used: - Adjusted Close - Simple moving average (SMA) - Exponential moving average (EMA) - Moving Average Convergence Divergence (MACD) - Bollinger band - Momentum - Trading volume - Price date

Price date will be splitted further to month, year, day of month, day of year.

Data will be separated to 3 sets: - train - validation - test

Model will be deployed as SageMaker endpoint.

## 5 Benchmark Model

To measure model performance mean squared error (MSE) will be used, better model should have lower error. When best model selected it will be tested against real stock data with measuring error. To check how well it perform we have to check result against similar models, for that reason similar researches was found. Please check links below.

According to read researches it's possible to acheive quite good accuracy on some markets. For this project try to acheived  $MSE \leq 0.002$  for markets with low volatility since markets with high volatility more risky and difficult to predic.

### 5.0.1 Research articles

- [Stock Price Prediction Mobile Application](#)
- [Stock Market Analysis + Prediction using LSTM](#)
- [Stock predictions with state-of-the-art Transformer and Time Embeddings](#)

## 6 Evaluation Metrics

As evaluation metric Mean Squared Error will be used.

## 7 Project Design

SageMaker will be used this project. For model development on of the SageMaker built-in algorithms will be used.

Project stages:

- Feature engenering
- Training model
- Deploy endpoint