

Novel Modeling of High-Frequency Stock Trading Data

Presentation

Yuying Huang¹, Ke Xu², Li Xing³, Xuekui Zhang¹

¹Dept. of Math & Stat, University of Victoria, ²Dept. of Economics, University of Victoria, ³Dept. of Math & Stat, University of Saskatchewan

Master's Project Defence 2021

- ① Introduction
- ② Research Problem
- ③ Strategies
- ④ Data Source
- ⑤ Results

1 Introduction

2 Research Problem

3 Strategies

4 Data Source

5 Results

Background Introduction

Research Background

High-frequency trading (HFT) rises from increased full electronic automation in stock exchanges, which features the use of extraordinarily high-speed and sophisticated computer programs for generating, routing, and executing orders. Investment banks, hedge funds and institutional investors design and implement the algorithmic trading strategy to identify the emerging surge in stock price and then speculate from the stock market.

What to Address

- It's impossible to analyze a whole dataset
- Information redundancy issue
- Most machine learning models assume independence among the data, while close observations in high-frequency stock price data are highly correlated

Innovation & Contribution

- Improve the stock price prediction accuracy by pre-processing raw data and constructing proper input variables (predictors) from high-frequency data
 - Propose novel strategy to decrease data redundancy while avoiding information loss in high-frequency data during the data thinning process
 - Propose novel strategy to break correlation between highly correlated data points while maintain robust model performance
 - Model long-term stock price input and evaluate its impact

① Introduction

② Research Problem

③ Strategies

④ Data Source

⑤ Results

Definition and Notification

1. Limit Order Book (LOB) : timestamp, best ask price, best bid price
2. Stock mid-price, as our dependent variable, is defined as:

$$MP = \frac{Ask_{best} + Bid_{best}}{2}$$

3. three-class classification problem: mid-price movement $Y = \{\text{Upward, Downward, Stationary}\}$

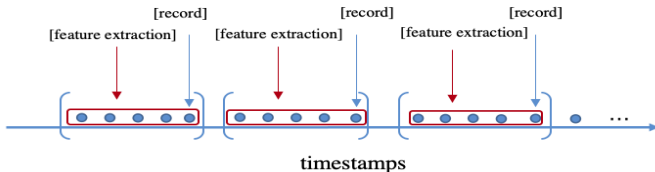


Figure 1: Illustration of event-based inflow framework

Definition and Notification

5. We regard the last observation in every Δt window as the "record" events, based on which, the stock mid-price movement can be labelled and features can be extracted accordingly. In our case, $\Delta t=5$ events, blocks are independent with each other

7.

$$\left\{ \begin{array}{ll} \text{Upwards} & \text{if } \frac{MP_{t+1}}{MP_t} > 1 + 10^{-5} \\ \text{Downwards} & \text{if } \frac{MP_{t+1}}{MP_t} < 1 - 10^{-5} \\ \text{Stationary} & \text{otherwise} \end{array} \right. \quad (1)$$

- ① Introduction
- ② Research Problem
- ③ Strategies**
- ④ Data Source
- ⑤ Results

Novel Strategies

Strategy I: data thinning with summary features–proposed feature set

To make amends for the lost information, during the data pre-processing, we involve new variables to extract and generalize summary features within each Δt window

Table 1: novel data thinning strategy with lately proposed “within-window” feature set

Attributes (k=5)	Category	Description
$V_4 = \{(P_{t-i}^{ask} - P_{t-i}^{bid}) / MP_{t-i}\}$	within window	bid-ask spread return
$V_5 = (P_{t-1}^{bid} - P_{t-k}^{bid}) / P_{t-k}^{bid}$	within window	best bid price difference return
$V_6 = (P_{t-1}^{ask} - P_{t-k}^{ask}) / P_{t-k}^{ask}$	within window	best ask price difference return
$V_7 = (P_{t-1}^{bid} - P_{t-k}^{ask}) / P_{t-k}^{ask}$	within window	bid-ask spread crossing return
$V_8 = (P_{t-1}^{ask} + \dots + P_{t-k}^{ask}) / k$	within window	mean best ask price in nanosecond
$V_9 = (P_{t-1}^{bid} + \dots + P_{t-k}^{bid}) / k$	within window	mean best bid price in nanosecond
$V_{10} = (MP_{t-1} + \dots + MP_{t-k}) / k$	within window	mean mid-price in nanosecond
$V_{11} = \sum_{t-k}^{t-1} V^{ask}$	within window	best ask price market depth
$V_{12} = \sum_{t-k}^{t-1} V^{bid}$	within window	best bid price market depth
$V_{13} = \{SD_{t-i}^{sec}\}$	within window	volatility in a second
$V_{14} = 1 / \{\text{time difference within the window}\}$	within window	window slope

Novel Strategies

Strategy II: "Sampling + Ensemble" Model

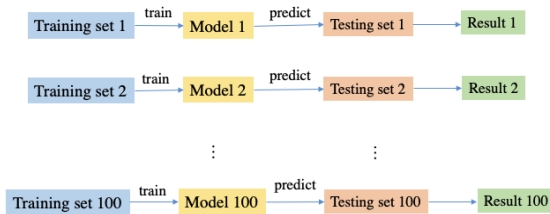


Figure 2: In an original framework, data sets are randomly split, trained and performed separately

Novel Strategies

Strategy II: "Sampling + Ensemble" Model

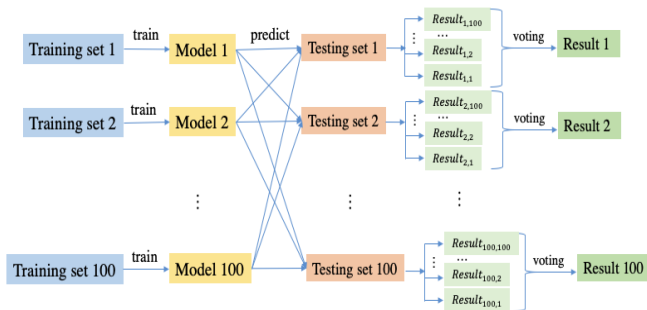


Figure 3: In a "Sampling + Ensemble" framework, data sets are randomly split. The same testing data set will be predicted through different trained models and a voting scheme will be applied to select the final prediction results

Novel Strategies

Strategy III: Combination of Long-term and Short-term Resolution

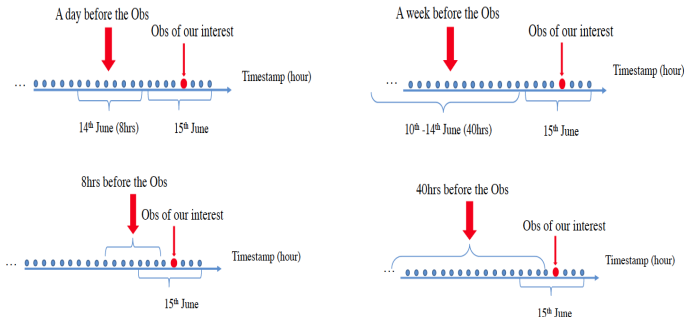


Figure 4: Illustration of two definitions of long-term curves. The upper two graphs represent the recurrence of daily pattern (definition 1) while the bottom two graphs represent the recurrence of weekly pattern (definition 2)

Novel Strategies

Strategy III: Combination of Long-term and Short-term Resolution

Table 2: Demonstration of 4 kinds of proposed FPCA scores including their signs, the number of variables used to represent them and the input data's time scope

Name	Definition	Prediction input
FPCA1	definition 1	8hrs (the day before today)
FPCA2	definition 1	40hrs (5 days before today)
FPCA3	definition 2	8hrs before this hour
FPCA4	definition 2	40hrs before this hour

- ① Introduction
- ② Research Problem
- ③ Strategies
- ④ Data Source
- ⑤ Results

Data Source

A glimpse at the data

- Stock high frequency intra day order-level data with nanosecond timestamps(e.g. HHMMSSxxxxxxxxxx)
- NYSE database
 - Provides best bid- and ask prices and their volume of the limit order book
 - Transaction prices and quantities for each trade
- Sample period: Jun 1, 2017 to Aug 31, 2017 (NYSE launched on May 31, 2019)
- Massive records from 30 most popular high-tech companies' stocks by market capitalization, such as APPLE, Microsoft

Data Source

Data Manipulation

- Filter
Initially remove all transactions recorded outside official trading time and clearly misrecorded transactions
- Winsorization & normalization
We first compute the 1st and 3rd quantiles $Q1$ and $Q3$ of our train sample, get the difference length $IQR=Q3-Q1$. And then, we replace the observations falling outside $[Q1-1.5IQR, Q3+1.5IQR]$ with the lower bound $Q1-1.5IQR$ and upper bound $Q3+1.5IQR$, respectively.
- Experiment and examine strategies
Repeated each of the experiment 100 times with different random seeds to subsample our data set and calculate the Recall, Precision, and F1 score. Compare the model performance between model with and without our novel strategies

Evaluation Indicator

Evaluate single strategy prediction accuracy

- Classification Indicator: Precision, Recall and F_β .
 - Precision $P = (\text{True Positive})/(\text{Predicted Positive})$
 - Recall $R = (\text{True Positive})/(\text{Positive})$
 - $F_\beta = (1 + \beta^2)PR/(\beta^2P + R)$, the weighted harmonic mean of P and R . When $\beta = 1$, meaning that both Precision and Recall are equally important, we obtain the $F_1 = 2PR/(P + R)$

Evaluate the difference between novel strategies

Two-sided Wilcoxon sign rank test (non-parametric test)

① Introduction

② Research Problem

③ Strategies

④ Data Source

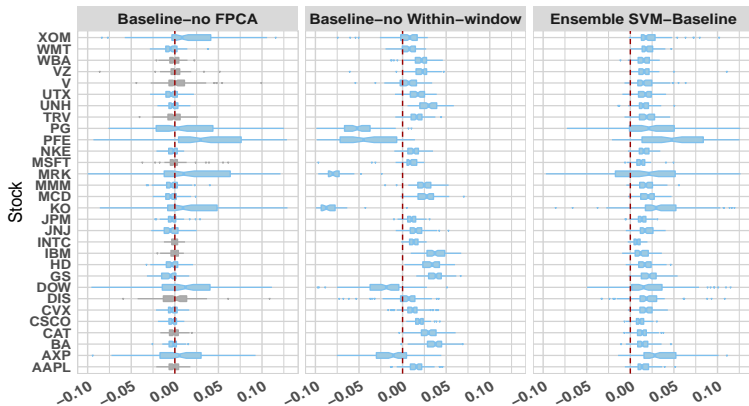
⑤ Results

Result

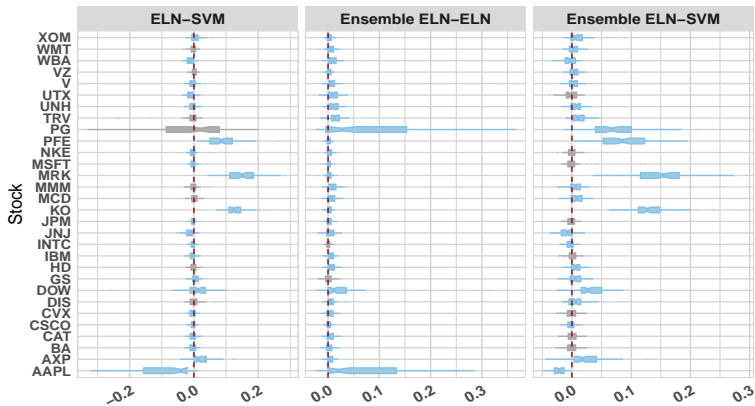
Comparisons

- baseline model: Strategy I and Strategy III (SVM and ELN)
- comparison between baseline model vs ensemble baseline model (SVM and ELN)
- comparison between baseline model vs baseline model without long-term variables (SVM)
- comparison between baseline model vs baseline model without within-window variables (SVM)
- variables selection (ELN)

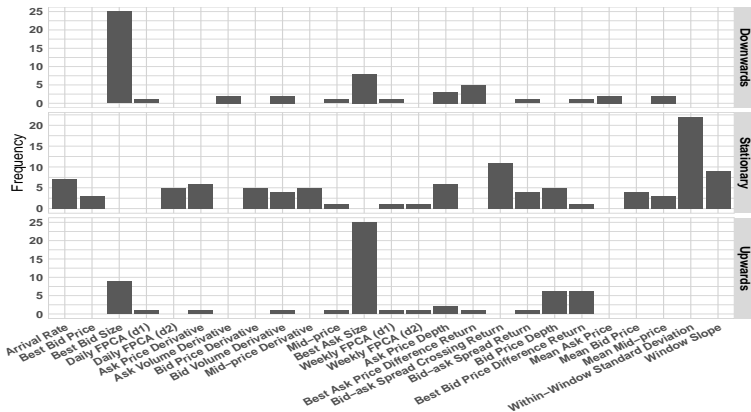
Result



Result



Result



Thanks!