

Forecasting high-frequency stock market returns using embedded limit order book data

Marius Sterling
Niels Wesselhöfft



International Research Training Group 1792
Ladislaus von Bortkiewicz Chair of Statistics
Humboldt-Universität zu Berlin



<http://irtg1792.hu-berlin.de>
<http://lvb.wiwi.hu-berlin.de>



Motivation

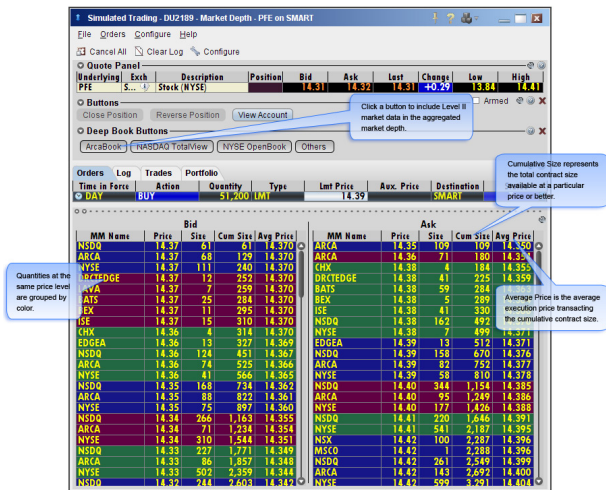


Figure: Limit Order Book (LOB) Level 2



Data Description

- Time span and aggregation
 - ▶ Time span from 01/07/2015 to 31/12/2015 (first iteration)
 - ▶ Tick data aggregated to second (last tick)
- Source
 - ▶ LOB tick data for NASDAQ stocks including Amazon from Lobster
 - ▶ Trading from 9:30 a.m. to 04:00 p.m. representing 23,400 seconds per day, 3m for half a year



Raw Data

- Limit prices and volumes (Bid | Ask) over
 - ▶ time $t = 1, \dots, T$,
 - ▶ depth $d = 1, \dots, D$ with $D = 200$
- Limit ask and bid price $P_{t,d}^a, P_{t,d}^b$
- Volume $V_{t,d}^a, V_{t,d}^b$
- Log price $p_{t,d}^a = \log P_{t,d}^a$

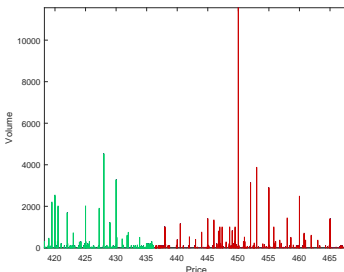


Figure: Limit Order Book in t



Cumulated Volumes

- Limit price vector and volume vector

$$P_t^a = [P_{t,1}^a, \dots, P_{t,D}^a]^\top$$

$$V_t^a = [V_{t,1}^a, \dots, V_{t,D}^a]^\top$$

- Cumulated Limit prices and volumes (Bid | Ask)

$$\bar{V}_{t,d}^b = \sum_{i=1}^d V_{t,i}^b$$

$$\bar{V}_{t,d}^a = \sum_{i=1}^d V_{t,i}^a$$

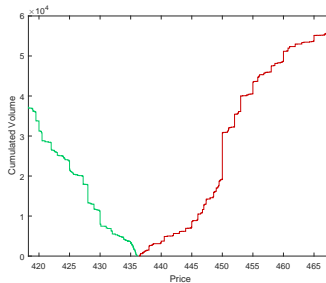


Figure: Limit Order Book with cumulated volumes in t



LOB over time

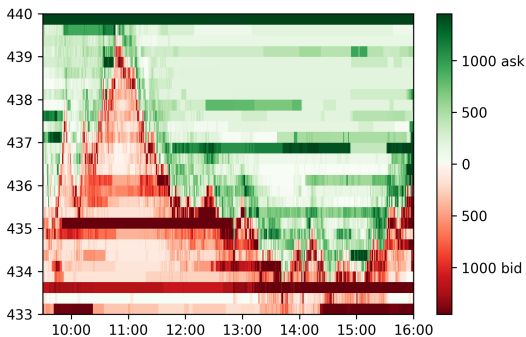


Figure: Limit Order Book (LOB) over one day



Embedding

- Embeddings for LOB data
 - ▶ Price and volume data
 - ▶ Return and volume data
 - ▶ Factor models
 - Semiparametric Factor Model (Hautsch, Härdle, Mihoci, 2012)
 - ▶ Measure for buying/selling pressure



Measure for buying/selling pressure

- Absolute difference from mid price

$$\delta_{t,d}^b = |P_{t,d}^b - P_t| \quad (1)$$

$$\delta_{t,d}^a = |P_{t,d}^a - P_t| \quad (2)$$

- Mid price $P_t \in [P_{t,1}^b, P_{t,1}^a]$

$$P_t = \sqrt{P_{t,1}^b P_{t,1}^a} \quad (3)$$

$$P_t = \frac{1}{2} (P_{t,1}^b + P_{t,1}^a) \quad (4)$$

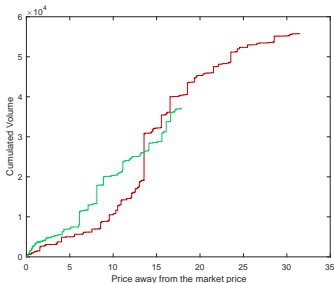


Figure: Cumulated volume over difference from mid price P_t in t



Theta θ_t

- Theta $\theta_t \in \mathbb{R}$

$$\theta_t = f(\bar{V}_{t,d}, \delta_{t,d}) \quad (5)$$

- Theta - difference

$$\bar{\theta}_{t,d} = \bar{V}_{t,d}^a - \bar{V}_{t,d}^b \quad (6)$$

- $d \in \{1, \dots, D\}$ is the specific depth chosen

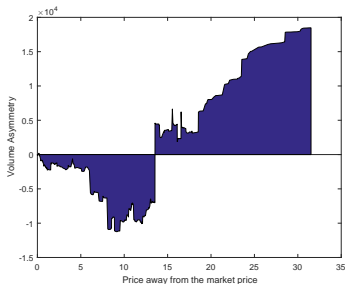


Figure: Cumulated volume difference $\bar{\theta}_{t,d}$



Theta θ_t

□ Theta - ratio

$$\tilde{\theta}_{t,d} = \begin{cases} \tilde{\theta}_{t,d}^+ = \overline{V}_{t,d}^b / \overline{V}_{t,d}^a, & \overline{V}_{t,d}^b > \overline{V}_{t,d}^a \\ \tilde{\theta}_{t,d}^- = -\overline{V}_{t,d}^a / \overline{V}_{t,d}^b, & \overline{V}_{t,d}^b \leq \overline{V}_{t,d}^a \end{cases} \quad (7)$$

□ Interpretation

- ▶ $\tilde{\theta}_{t,d}^+$ as multiple of bid over ask cum volume (buying pressure)
- ▶ $\tilde{\theta}_{t,d}^-$ as multiple of ask over bid cum volume (selling pressure)

□ Depth weights

- ▶ Uniform
- ▶ Exponential



Return forecasting

- Return forecast up to 600 seconds

$$r_{t+i} = f(X_t, X_{t-1}, \dots) + \varepsilon_{t+i}, \quad (8)$$

$$\varepsilon_{t+i} \sim F(), \quad i = 1, \dots, 600 \quad (9)$$

- Dependent variable: Log return

$$r_t = \log P_t - \log P_{t-1} = p_t - p_{t-1} \quad (10)$$

- Independent variable(s) X_t, X_{t-1}, \dots

$$X_t = \begin{cases} 1) [P_{t,1}^b, V_{t,1}^b, P_{t,1}^a, V_{t,1}^a, \dots, P_{t,D}^b, V_{t,D}^b, P_{t,D}^a, V_{t,D}^a] \\ 2) [r_{t,1}^b, V_{t,1}^b, r_{t,1}^a, V_{t,1}^a, \dots, r_{t,D}^b, V_{t,D}^b, r_{t,D}^a, V_{t,D}^a] \\ 3) \theta_t = f(\bar{V}_{t,d}, \delta_{t,d}) \end{cases} \quad (11)$$



Modeling Approaches

Model	Embedding	1 Prices and volumes	2 Log returns and volumes	3 Asymmetry estimator θ
1 Linear Regression		11	12	13
2 RNN (e.g. LSTM)		21	22	23
3 TCN		31	32	33

Table: Modeling matrix



13 Linear Regression + θ_t

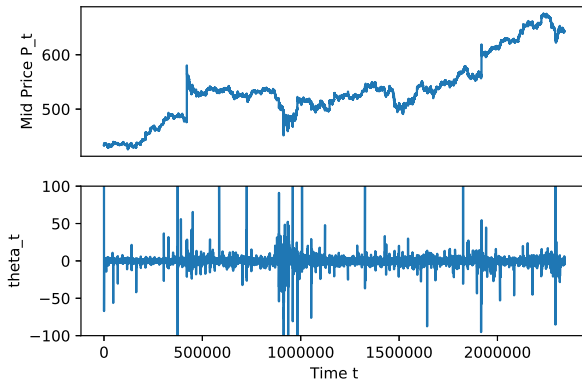


Figure: (upper) Amazon prices series (lower) θ_t over time



13 Linear Regression + θ_t

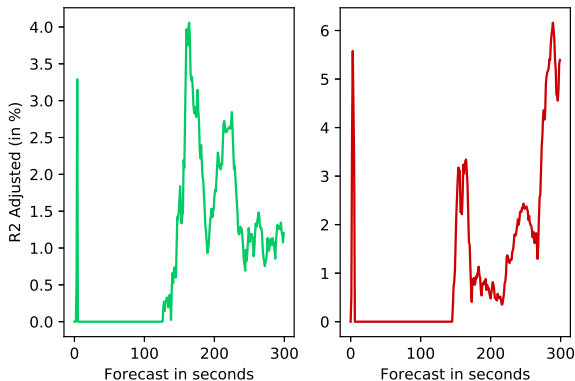


Figure: (left) R^2 for $X_t = [\theta_t^+]$ (right) R^2 for $X_t = [\theta_t^-]$



23 Linear Regression $+\theta_t + \theta_t^2$

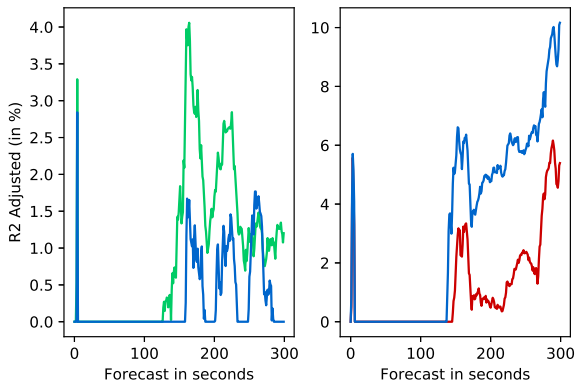


Figure: (left) R^2 for $X_t = [\theta_t^+, (\theta_t^+)^2]$ (right) R^2 for $X_t = [\theta_t^-, (\theta_t^-)^2]$



TCN inspired model architecture

□ Layers

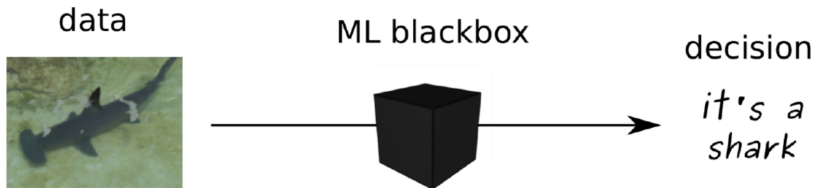
- ▶ Feature generating layer
- ▶ Inception layers
- ▶ Dimension reduction layers

□ Training parameters

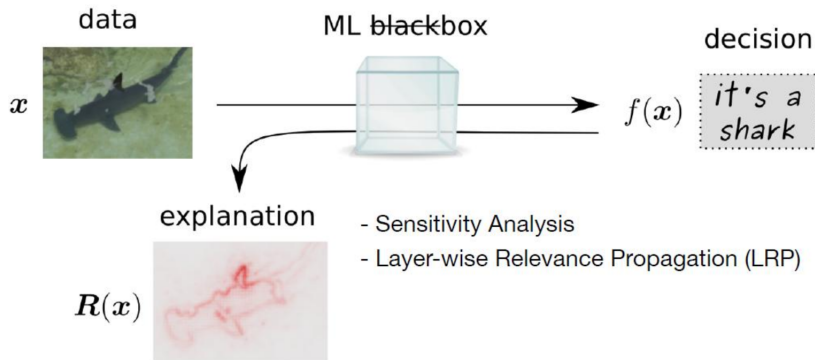
- ▶ Depth: 10
- ▶ Lag: 64 seconds
- ▶ Prediction horizon: 15 seconds
- ▶ Loss: MSE
- ▶ Optimizer: Adam
- ▶ Batchsize: 128



NNs as a blackbox



LRP



Input relevance for 15 second forecast

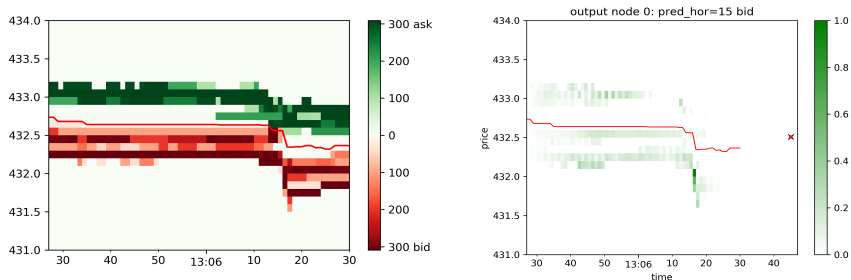


Figure: (left) LOB input and (right) scaled forecast relevance (LRP) for input of 15 sec ahead bid price prediction, — mid price, • price, × predicted price



Input relevance for 15 second forecast

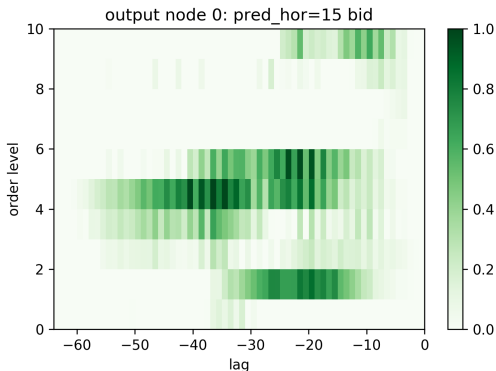


Figure: Dept and Lag relevance



Research goals

- Supervised methods
 - ▶ Stock markets returns forecasts for high frequencies
- Embeddings
 - ▶ Will unstructured information still lead to superior NNs?
 - ▶ Compare parametric depth structure (uniform, exponential) to the forecast relevance of NNs (LRP)
- Out-of-sample Test
 - ▶ Performance holds under transactions costs
 - ▶ Out-of-sample and out-of-stock

