



UNIVERSITY OF CAPE TOWN
IYUNIVESITHI YASEKAPA • UNIVERSITEIT VAN KAAPSTAD

Data Analysis for High-Frequency Trading

Assignment 3

Hawkes Process and Order Book Simulation.

Julian Albert
ALBJUL005



Department of Statistical Sciences
University of Cape Town
South Africa
October 5, 2019

Contents

| | |
|--|-----------|
| Introduction | 2 |
| Classification | 2 |
| Estimate a Hawkes Process for Order-Book Events | 3 |
| Four Dimensional Hawkes Process | 8 |
| Simulate the Order-Book | 12 |
| Conclusions | 13 |

Introduction

The Hawkes process is a class of point processes whose future depends on their own history. This class of processes is particularly useful for estimating and simulating order book events as they occur asynchronously in event-time and there exists relationships between the occurrence of order-book events. As seen in previous Assignments trade events exhibit auto-correlation in the event arrival times which leads to the observation that event arrival causes the conditional intensity function to rise. This event arrival process can thus be thought of as an example of a “self-exciting” process which can be modelled using the Hawkes process — a counting process that models a sequence of “arrivals” of some type over time. Each arrival excites the process in that the likelihood of a subsequent arrival is increased for some time period after the initial arrival. We will estimate the parameters of a multivariate Hawkes process using a single day of trading for a particular stock. The resulting non-homogeneous parameter estimates will be used to simulate the arrival of order-book events (classified into eight subcategories). Using these arrivals times, we then simulate the order-book by means of plotting the best-ask price, best-bid price and mid-price over a day of trading.

Classification

The first part of this assignment is the classification of order-book events into eight subcategories, this is done using Table (1). First we classify trades as either buyer or seller initiated and then we determine whether or not they have moved the order book. To determine whether the trades are buyer (+1) or seller (-1) initiated market orders we use the quote data. We can further subdivide these into buyer and seller initiated market orders that do or do not move the best bid/ask price. For buyer (seller) initiated trades, if the best sell (buy) price after the trades are higher (lower) than the best sell (buy) price before the trade, then we classify the trade as “Aggressive buy (sell) moves the ask (bid)” otherwise the trade is classified as “Aggressive buy (sell) doesn’t move the ask (bid)”. The LO classification are self explanatory, where the bid (ask) is at or lower (higher) then the previous bid (ask) the classification is an order than does not move the best bid (ask). Similarly if the bid (ask) lies in the spread it is classified as between quotes.

| No. | Buy/Sell | Order Type | Price Move? | Classification |
|-----|----------|------------|-------------|------------------------------|
| 1 | Sell | (MO) | ✓ | Market Sell Moves Bid |
| 2 | Sell | (MO) | ✗ | Market Sell Doesn’t Move Bid |
| 3 | Buy | (MO) | ✓ | Market Buy Moves Ask |
| 4 | Buy | (MO) | ✗ | Market Buy Doesn’t Move Ask |
| 5 | Buy | (LO) | ✓ | Bid Between Quotes |
| 6 | Buy | (LO) | ✗ | Bid \leq Best-Bid |
| 7 | Sell | (LO) | ✓ | Ask Between Quotes |
| 8 | Sell | (LO) | ✗ | Ask \geq Best-Ask |

Table 1: Order-Event Classification as Per [3], MO refers to Market-Order whilst LO refers to Limit-Order.

Estimate a Hawkes Process for Order-Book Events

Before we dive into Estimating the process for different order-book events we can introduce the process. A Hawkes process is a point-processes $N(t)$ given the number of events on the time interval $[0, t]$ for some $s \leq t$ with $N(s) \leq N(t)$

$$\begin{aligned}\mathbb{P}[N(t+h) - N(t) = 1 | N(s)] &= \lambda(t)h + o(h) \\ \mathbb{P}[N(t+h) - N(t) > 1 | N(s)] &= o(h)\end{aligned}\tag{1}$$

where the intensity $\lambda(t)$ is self-exciting with event arrival-times t_m and kernel $\phi = \alpha\beta e^{-\beta(t)}$ we get

$$\lambda(t) = \mu + \int_{-\infty}^t \phi(t-s) dN(s) = \mu + \sum_{t_m < t} \phi(t-t_m)\tag{2}$$

a point-process with an intensity function that is linear in past events, where the baseline intensity μ are events not stimulated by prior events. Defining

$$d\lambda(t) = -\beta\lambda(t)dt + \alpha\beta dN(t)\tag{3}$$

and multiplying by the kernel $e^{-\beta(t)}$

$$g(\lambda, t) = \lambda e^{\beta t}\tag{4}$$

Such that

$$\frac{\partial g}{\partial t} = \lambda\beta e^{\beta t} \quad \frac{\partial g}{\partial \lambda} = e^{\beta t} \quad \frac{\partial^2 g}{\partial \lambda^2} = 0\tag{5}$$

Applying *Ito's lemma* we get

$$dg = \lambda\beta e^{\beta t} dt + e^{\beta t} d\lambda\tag{6}$$

$$d\lambda e^{\beta t} = \lambda\beta e^{\beta t} dt + e^{\beta t} (-\beta\lambda dt + \alpha\beta dN)\tag{7}$$

$$= (\lambda\beta e^{\beta t} - \lambda\beta e^{\beta t}) dt + \alpha\beta e^{\beta t} dN\tag{8}$$

$$= \alpha\beta e^{\beta t} dN\tag{9}$$

we therefore obtain

$$\lambda(t)e^{\beta t} = \lambda(0) + \alpha\beta \int_0^t e^{-\beta(t)} dN(t)\tag{10}$$

$$\lambda(t) = \mu + \alpha\beta \int_0^t e^{-\beta(t-s)} dN(s)\tag{11}$$

where $\mu = e^{-\beta t} \lambda(0)$. In [2] we get an exact approximation of this integral which will be used for calculating the intensities and plotting them as in Figure (3), (4) and (5) expressed by Equation (12).

$$\lambda(t) = \mu + \sum_{t_i < t} \alpha\beta e^{-\beta(t-t_i)}\tag{12}$$

In a multivariate setting we have that the intensity function is

$$\lambda_m(t) = \mu_m + \alpha_{mn}\beta_m \int_0^t e^{-\beta_m(t-s)} dN^m(s) \quad (13)$$

where α , $\vec{\beta}$, $\vec{\mu}$ represent a matrix and vectors of $k \times k$ and $k \times 1$ respectively with k being the number of dimensions (events).

One challenge when modelling using self-exciting point processes is estimating parameters from observed data, to do this we can use a maximum likelihood estimate for our parameters where the likelihood function is give by [1] in Equation (14)

$$L(\Theta|t_n) = \prod_{i=1}^N \lambda(t_i; \Theta) \exp\left(-\int_{t_{i-1}}^{t_i} \lambda(s; \Theta) ds\right) \quad (14)$$

To obtain the MLE we take the natural logarithms to get the log-likelihood in Equation (15)

$$\ell(\Theta|t_n) = -\int_0^{t_n} \lambda(s; \Theta) ds + \sum_{i=1}^n \ln(\lambda(t_i; \Theta)) \quad (15)$$

minimising the negative log-likelihood is equivalent to maximising the likelihood function, however, the optimisation is less expensive and more computationally stable optimisations. Note here that Θ represent the parameter set containing α , $\vec{\beta}$, and $\vec{\mu}$.

To perform this optimisation we use `optim()` in R, the nragtive log-likelihood function to be minimised corresponds to the `likelihoodHawkes()` function from the `Hawkes()` package [6]. The function does require the following constraints;

$$\begin{aligned} \alpha, \vec{\mu} &\geq 0 \\ \text{eigenvalue}(\text{diag}(\vec{\beta}) - \alpha) &> 0 \end{aligned}$$

which are implemented in the code by creating a new function which either calls `likelihoodHawkes()` if the constraints are met or calls `likelihoodHawkes()` plus some very large penalty — making the solution infeasible as we are minimising the objective.

The first issue with trying to estimate Θ for all 8 event types presented previously is that of local maxima. As the dimensions increase the shape of the negative log-likelihood function can become complex and may not be globally convex. The lack of convexity in the objective function would result in the MLE being the local maximum rather than the global maximum [5]. A usual approach used in trying to identify the global maximum involves using several sets of different initial values for the maximum likelihood estimation, however, this proved unsuccessful.

Alternatively, I look at a bivariate Hawkes process for two sets of cases, the first is aggressive buy (sell) market orders where the price is greater (less) than the best ask (bud). The second is passive limit orders where the new bid (ask) is higher (lower) than the previous bid (ask) thus updating the quote in the LOB. Initial parameters are chosen arbitrarily with the assumption that the optimisation algorithm will converge, the starting values for α , $\vec{\beta}$, and $\vec{\mu}$ are chosen to be a small diagonal matrix, a vector of values between 0.5 and 1, and the mean number of events divided by seconds in the day respectively. The starting parameters for aggressive and passive orders are

$$\vec{\lambda}_{agg} = [0.014 \quad 0.010] \quad \vec{\lambda}_{pass} = [0.133 \quad 0.090]$$

with kernel starting parameters

$$\begin{aligned}\boldsymbol{\alpha}_{agg} &= \begin{bmatrix} 0.01 & 0 \\ 0 & 0.01 \end{bmatrix} & \vec{\beta}_{agg} &= [0.05 \quad 0.05] \\ \boldsymbol{\alpha}_{pass} &= \begin{bmatrix} 0.01 & 0 \\ 0 & 0.01 \end{bmatrix} & \vec{\beta}_{pass} &= [0.05 \quad 0.05]\end{aligned}$$

And optimal solutions

$$\begin{aligned}\vec{\lambda}_{agg}^* &= [0.006 \quad 0.006] & \boldsymbol{\alpha}_{agg}^* &= \begin{bmatrix} 0.023 & 0.003 \\ 0.000 & 0.017 \end{bmatrix} & \vec{\beta}_{agg}^* &= [0.044 \quad 0.051] \\ \vec{\lambda}_{pass}^* &= [0.033 \quad 0.023] & \boldsymbol{\alpha}_{pass}^* &= \begin{bmatrix} 0.066 & 0.002 \\ 0.000 & 0.075 \end{bmatrix} & \vec{\beta}_{pass}^* &= [0.088 \quad 0.100]\end{aligned}$$

We can see that the optimal parameters have moved off of the initial parameters, in both cases optimisation converged (likely to local optima). It is worth noting that changes in α impact the number of events that arise from simulating the Hawkes process with the optimal parameters. Table (2) shows the four events and how the number of events observed in the data compare to the number of events simulated using a bivariate Hawkes process. Although these results are not necessarily outstanding it was decided to leave them unaltered as any further tuning would simply introduce bias and likely lead to an extreme case of overfitting. We can see however that the 2-dimensional Hawkes process approximates relatively well.

| No. | Classification | n_{true} | n_{sim} |
|-----|----------------|------------|-----------|
| 1 | Agg. Sell | 375 | 449 |
| 3 | Agg. Buy | 262 | 223 |
| 5 | Pass. Buy | 3658 | 3825 |
| 7 | Pass. Sell | 2477 | 2339 |

Table 2: Number of Event Arrivals True vs Simulated

To evaluate the parameters influence on the likelihood function I simulate 10000 (100×100) combinations of $\vec{\beta}$, and $\vec{\mu}$ as well as 10000 $\boldsymbol{\alpha}$ diagonal matrices. Where the ranges are

$$\begin{aligned}0.0010 &\leq \mu \leq 0.1500 \\ 0.0500 &\leq \beta \leq 5.0000 \\ 0.0001 &\leq \alpha_{ii} \leq 0.0400\end{aligned}$$

For aggressive and passive the range was tuned. In the univariate case we would vary the parameters as opposed to varying combinations of parameters making insightful inference much easier. We plot the contours for the various combinations of parameters in the multi-dimensional likelihoods. This is achieved by fixing one of the three parameters at their MLE and plotting the remaining two parameters against one another to create a profile for visualising the relationship between these parameters, the resulting objectives for the bivariate Hawkes process are illustrated in Figure (1) and (2).

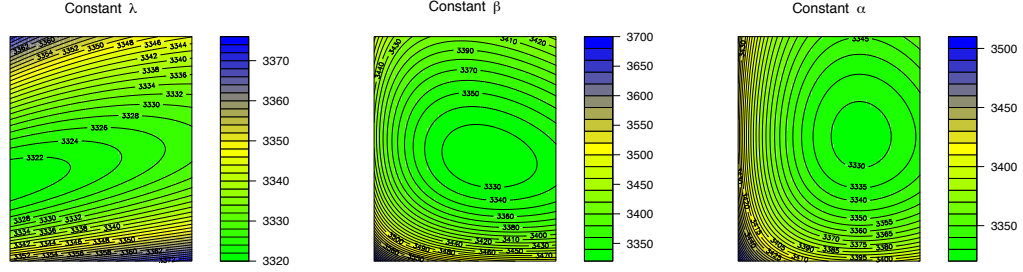


Figure 1: Objective Surfaces for Aggressive Order Types - Bivariate Hawkes

In both case we observe relatively “well-behaved” objective functions. Where we have smaller numbers of events it appears better results for simulation are achieved using smaller parameters for α and β , thus our initial starting point may somewhat be “directly” related to the number of event arrivals. However, the results for starting parameters on a mix of number of events is relatively robust (i.e. choosing different starting parameters for each event breaks down, maybe due to some iid. assumptions).

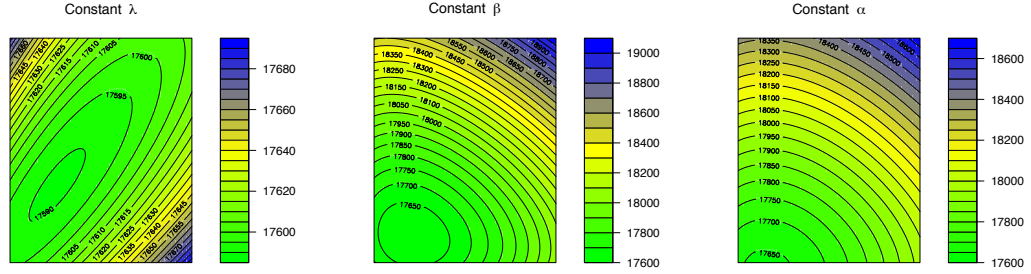


Figure 2: Objective Surfaces for Passive Order Types - Bivariate Hawkes

We can see that the objective surface is relatively stable in two-dimensions as the likelihood function (note $\lambda \equiv \mu$) is relatively simple. However, it is worth noting that the range for α is very small. Also unclear here is exactly how the dependencies and auto-correlations between parameters effect the process simulations. Recall that $N(t)$ is the number of “event arrivals” of the process by time t and that the sequence of event times is assumed to be observed within $[0, T]$ for finite T [5]. A property of the Hawkes process is that the events usually arrive clustered in time. Empirically what may occur is that the process might have started sometime in the past, prior to the moment when we start observing it $t = 0$. Hence, there may be unobserved event times which occurred before time zero, which could have generated “offspring” events during the interval $[0, T]$ [5]. It is possible that the unobserved event times could have had an impact during the observation period, i.e. some time $t > 0$, but because we are not aware of them, their contribution to the event intensity is not recorded. Such phenomenon are referred to as edge effects [5].

Of particular interest to us may be the intensities which can be recovered from our simulated Hawkes processes using Equation (10). We can plot the intensities in Figures (3) and (4), we see that in the aggressive case the events number of simulated events are similar. Also worth noting is the sporadic nature of both aggressive buys and sells appear similar with occasionally higher peaks being observe in the aggressive sells. With the passive order types we observe a higher baseline intensity and extreme jumps (LOB stuffing), in a South African context this may be due to the lack of market makers (liquidity providers). We also observe less aggressive buy and passive sell events relative to passive sell and aggressive buy events.

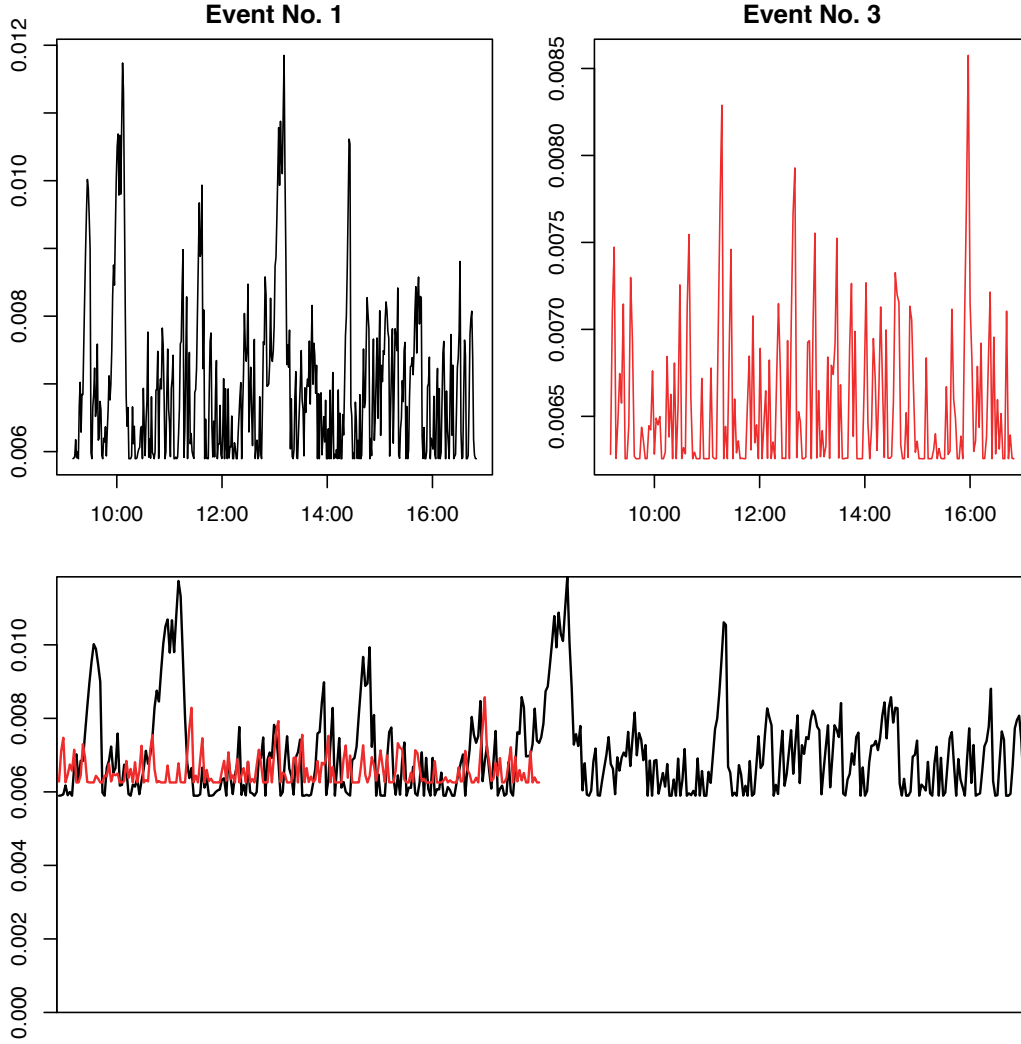


Figure 3: Intensities for Aggressive Order Types $k = 2$.

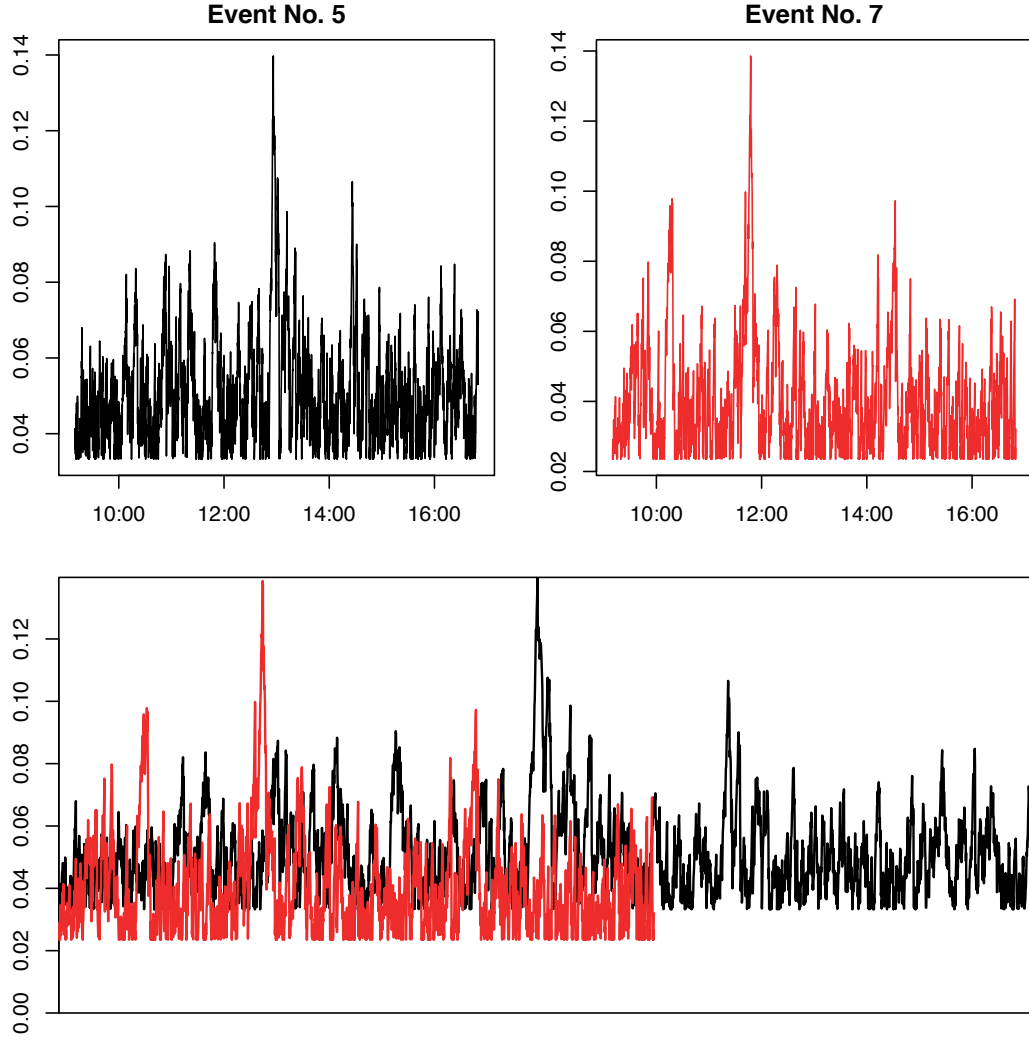


Figure 4: Intensities for Passive Order Types $k = 2$.

Four Dimensional Hawkes Process

Now we consider the case where we combine the two bivariate cases into a four-dimensional Hawkes process. Here the numerical instability is much higher, relatively many more iterations are used for convergence using the Nelder-Mead Simplex algorithm.

$$\vec{\lambda} = [0.014 \quad 0.009 \quad 0.133 \quad 0.090] \quad \alpha = \begin{bmatrix} 0.01 & 0 & 0 & 0 \\ 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0.01 & 0 \\ 0 & 0 & 0 & 0.01 \end{bmatrix} \quad \vec{\beta} = [0.05 \quad 0.05 \quad 0.05 \quad 0.05]$$

which are optimised to

$$\vec{\lambda}^* = [0.020 \quad 0.008 \quad 0.153 \quad 0.111] \quad \vec{\beta}^* = [0.005 \quad 0.000 \quad 0.163 \quad 0.092]$$

$$\boldsymbol{\alpha}^* = \begin{bmatrix} 0.000 & 0.011 & 0.001 & 0.000 \\ 0.005 & 0.003 & 0.000 & 0.000 \\ 0.011 & 0.014 & 0.005 & 0.001 \\ 0.018 & 0.013 & 0.000 & 0.000 \end{bmatrix}$$

Using the new parameter set Θ^* we obtain the following estimates for number of event arrivals

| No. | Classification | n_{true} | n_{sim} |
|-----|----------------|------------|-----------|
| 1 | Agg. Sell | 375 | 188 |
| 3 | Agg. Buy | 262 | 82 |
| 5 | Pass. Buy | 3658 | 4895 |
| 7 | Pass. Sell | 2477 | 3478 |

Table 3: Number of Event Arrivals True vs Simulated - 4D Hawkes

Again, this proves less stable and the difference between empirical and simulated value is greater than in the bivariate case — intuitive as the number of events is highly variate the parameters should be less efficient. The process has difficulty finding the solution that can account for low MO activity and high LO activity and the `Hawkes()` package implicitly defines constraints that makes it difficult to initialise the parameters for each dimension in Θ^* at different points. Higher dimensional solutions would require advanced usage of non-linear constrained optimisations. Again we can plot the intensities, we see that the passive orders are much more sporadic with the jumps. The combined graph is very skewed as we observe many more passive orders and higher intensities hinder any visual inference from relative intensities, we can plot relative to time scale.

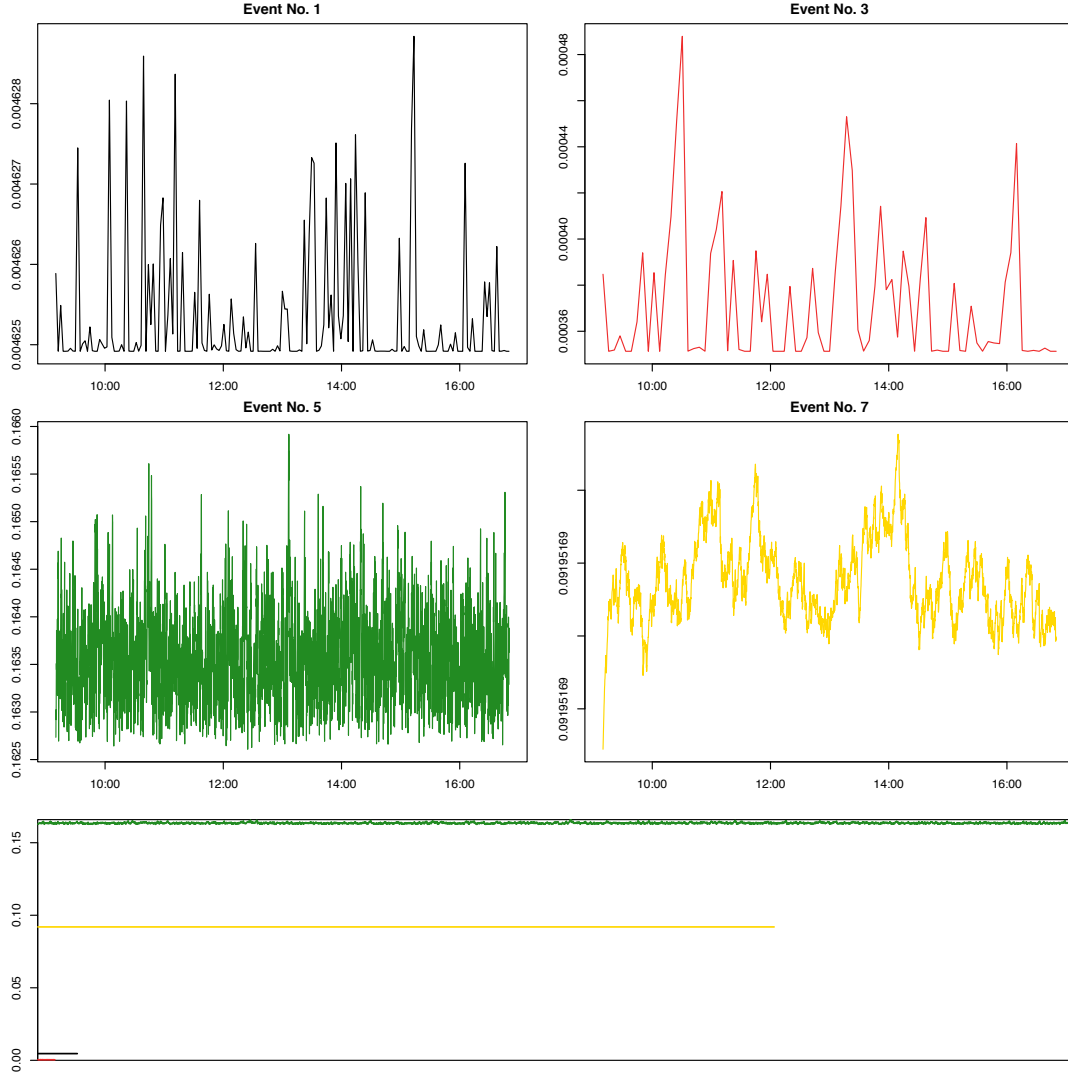


Figure 5: Comparison of Intensities $\lambda(t)$ for Different Order Types and Order 4 Hawkes Process.

Having applied an exponential kernel we can check whether the resulting inter-arrival times are exponentially distributed. To do this we assume no cross-excitation between different order types (i.e. set the cross diagonal values of the α matrix to 0). We can define a recursive duration [4] (time-change) as

$$\Lambda(t_{j-1}, t_j) = [t_j - t_{j-1}] \mu + \alpha \left[1 - e^{-\beta(t_j - t_{j-1})} \right] E(j-1) \quad (16)$$

Where,

$$E(j-1) = 1 + e^{-\beta(t_{j-1} - t_{j-2})} A(j-2)$$

From Equation (16) we can generate empirical quantile plots for the distribution of the durations of the simulated Hawkes process with a particular choice of estimated parameters — optimal

parameter set Θ^* . The results of checking exponential results are illustrated in Figure (6) which shows the durations of the simulated results as a QQ - plot compared to an exponential distribution with rate parameter $1/\bar{\Lambda}$. The plots show that the exponential assumption holds well for lower quantiles and for the aggressive orders seems relatively well-fitted, however, deviate about the exponential quantile line for larger quantiles especially for passive orders.

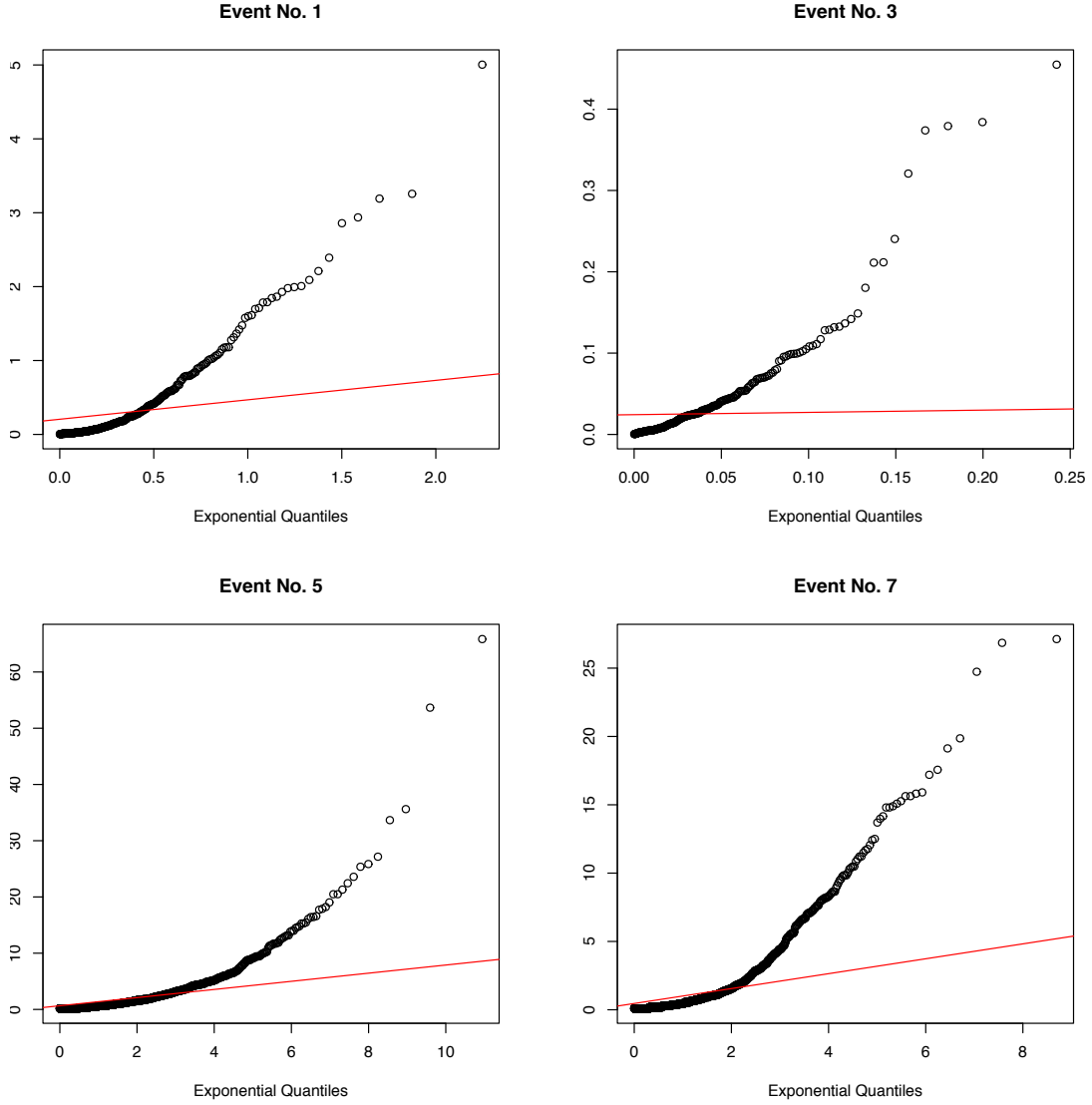


Figure 6: Exponential QQ-Plots for Different Events.

Simulate the Order-Book

Simulating the Order book is a unique problem, after having simulated the event arrivals using a calibrated Hawkes process we know when the events arrive. What is unknown is the value and quantity of the events. To generate an accurate reflection of the order-book is to make many assumptions about the order type dynamics, in a real-world attempt we would simulate order arrival volumes and use the volumes to determine a relative price impact. For the purposes of this Assignment we simulate prices only, noting that in order for the simulation to conform to empirical results we are almost surely overfitting. In both the simulated and empirical LOB the ask is represented in blue whilst the bid is in black. Figure (7) shows the results of 3 runs for a simulated order-book. We know the event arrivals from calibration in the previous steps, what is unknown is the dynamics of the orders which arrive. I assume that for an aggressive buy (sell) order the price impact follows an exponential distribution causing an uptick (downtick) on the best ask (bid) of $\exp(\lambda)$ where λ is a rate of one. For a passive buy (sell) I assume a truncated normal distribution with mean zero and standard deviation one hundred, the order will move the bid (ask) up (down) with this distribution. Note here the truncated distribution is used to ensure the spread is not crossed.

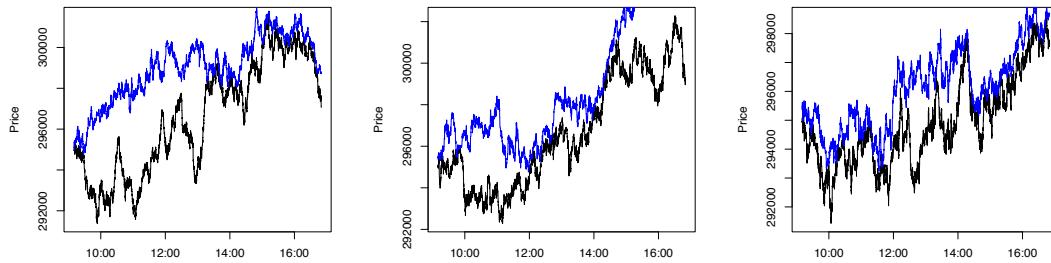


Figure 7: Simulated Order-Book for Four Event Types using Hawkes Process

Figure (7) is clearly not very robust in capturing the market dynamics, this is likely due to a rather “guess-like” justification for the distributions of prices. As mentioned previously a better representation would be to simulate volumes and place assumptions on the price impact. It is worth mentioning that if we simulate the process many times and take the mean of the paths with the above mentioned distributions I actually find that bids decrease whilst asks increase (the spread increases causing a big divergence). Thus restrictions for simulating the order book would likely need to capture the spread contracting and relaxing based on market pressures. When the simulated LOB is compared to Figure (8) we can loosely create a justification for the simulation being a success, but it is rather chosen to be like this and the power of the simulation is particularly weak.

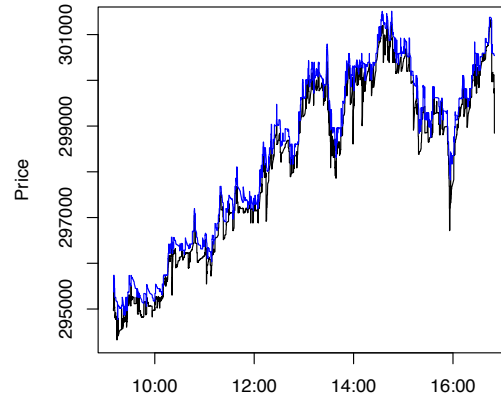


Figure 8: True LOB.

Conclusions

This assignment was particularly challenging, the whole premise of simulating order-book dynamics is full of biases due to lack of formal definition. Simply classifying trades — the first step in the process — is subjective and in reality there are many different order types that would make calibrating a Hawkes process infeasible. The process provides a good baseline for event arrival simulation, however, knowing when an event occurs is not enough to model market dynamics. The problem is synonymous with knowing I have an exam on the 22nd. I know when the exam (event) is occurring but I can not precisely model the content (volume and price impact) in the exam. Sure, I have an outline based on my knowledge of the examiner and content load (history and domain knowledge) which allows me to make assumptions about the event dynamics. But there is almost surely no way of providing exact simulation in the presence of noise, this begs the question as to whether we are better off making simplistic model assumptions or whether complexity trumps all.

- END -

References

- [1] DALEY, D., AND VERE-JONES, D. Basic properties of the poisson process. *An Introduction to the Theory of Point Processes: Volume I: Elementary Theory and Methods* (2003), 19–40.
- [2] DASSIOS, A., ZHAO, H., ET AL. Exact simulation of hawkes process with exponentially decaying intensity. *Electronic Communications in Probability* 18 (2013).
- [3] LARGE, J. Measuring the resiliency of an electronic limit order book. *Journal of Financial Markets* 10, 1 (2007), 1–25.
- [4] OZAKI, T. Maximum likelihood estimation of hawkes’ self-exciting point processes. *Annals of the Institute of Statistical Mathematics* 31, 1 (1979), 145–155.
- [5] RIZOIU, M.-A., LEE, Y., MISHRA, S., AND XIE, L. A tutorial on hawkes processes for events in social media. *arXiv preprint arXiv:1708.06401* (2017).
- [6] ZAATOUR, R. *hawkes: Hawkes process simulation and calibration toolkit*, 2014. R package version 0.0-4.

Appendix

R Code

```
1 ##
2 #
3 # Author: Julian Albert
4 # Date: 22 September 2019
5 #
6 # Description:
7 # HFT Assignment 3 - Basically want to work with TAQ data and try simulate the
8 # order book activity using hawkes processes
9 #
10 #-----#
11
12 # 0. Clean Workspace, Directory and Library ----
13
14 ## Clean Workspace
15 rm(list=ls())
16 dev.off() # Close Plots
17 setwd("~/") # Clear Path to User
18
19 ## Locations
20 project_folder <- "/Documents/UCT/Coursework/HFT"
21 loc_script <- "/Assignment_3/UCT_Assignment/Code/HFT_Ass3"
22 loc_figs <- "/Assignment_3/UCT_Assignment/Figs"
23 loc_data <- "/Data"
24
25 ## Directories
26 dir_script <- paste("~/", project_folder, loc_script, sep = '')
27 dir_figs <- paste("~/", project_folder, loc_figs, sep = '')
28 dir_data <- paste("~/", project_folder, loc_data, sep = '')
29
30 ## Filenames
31 file_dat <- "/dat_TAQ.rds"
32
33 ## Set Working Directory to Script Location
34 setwd(dir_script)
35
36 ## Libraries - Lazyload
37 if (!require("pacman")) install.packages("pacman")
38 p_load(tidyverse, lubridate, zoo, data.table, Cairo, hawkes, truncdist)
```

```

39
40 ## Import Clean Data
41 dat_TAQ <- readRDS(file = paste(dir_data, file_n_dat, sep = ""))
42 #setwd(dir_data)
43 #load(file = "dat_SimHawkes.RData")
44
45 ## Options
46 options(digits.secs = 3)
47
48 # 0.1 Simulate Accept/Reject for Inhomogenous Poisson Process ----
49
50 ## Plot an Inhomogeneous Poisson process by thinning
51 set.seed(123)
52 func.lambda <- function(t){2 + sin(t)}
53
54 M = 4; Time = 4*pi
55 p = numeric(); py = numeric()
56 r = numeric(); ry = numeric()
57 t = 0
58
59 while(t < Time)
60 {
61   t = t + rexp(1, M)
62   if(t < Time){u = runif(1)}
63   if(u <= func.lambda(t)/M){
64     p = c(p, t)
65     py = c(py, M*u)}
66   else{
67     r = c(r, t)
68     ry = c(ry, M*u)}
69 }
70
71 t <- seq(0, Time, by = 0.01)
72 plot(c(0, Time), c(0, M), type = 'n', xlab = "Time", ylab = "",
73      xaxs = 'i', yaxs = 'i')
74 lines(t, func.lambda(t), type = 'l')
75 lines(c(0, Time), c(M, M), lwd = 4)
76 points(p, py, pch = 16, col = 'darkgreen')
77 points(r, ry, pch = 3, col = 'red')
78 segments(p, rep(0, length(p)), p, py, lty = 'dashed', col = 'darkgreen')
79 dev.off()
80
81 # 1. Data Classification ----
82
83 ## Let's work with Naspers
84 dat_NPN <- dat_TAQ$NPN
85
86 # Select one day of trading
87 random_date <- sample(unique(as.Date(dat_NPN$DateTimeL)), 1)
88 dat_NPN_oneday <- dat_NPN %>%
89   filter(as.Date(DateTimeL) == random_date)
90
91 # Classify Market Orders... >> There is a bug here, if sell follows buy?
92 dat_NPN_oneday_classified <- dat_NPN_oneday %>%
93   mutate(Test.Sign = Trade.Sign, After_Bid = lead(L1.Bid),
94          Before_Bid = lag(L1.Bid), After_Ask = lead(L1.Ask),
95          Before_Ask = lag(L1.Ask)) %>%
96   mutate(MO_Class =
97     ifelse(Trade.Sign == -1 & After_Bid < Before_Bid, "Sell Moves Bid",
98            ifelse(Trade.Sign == -1, "Sell Doesn't Move Bid",
99                   ifelse(Trade.Sign == 1 & After_Ask > Before_Ask, "Buy Moves Ask",
100                          ifelse(Trade.Sign == 1, "Buy Doesn't Move Ask", NA))))))
101
102 # Classify Limit Orders... >> Bids and Asks need to be separate
103
104 ## Bid Classification
105 tmp_bids_class <- dat_NPN_oneday %>%
106   select(DateTimeL, L1.Bid, L1.Ask) %>%
107   mutate(Before_Bid = lag(L1.Bid),
108          Before_Ask = lag(L1.Ask)) %>%
109   mutate(LO_Bid_Class =
110     ifelse(L1.Bid > Before_Bid & L1.Bid < Before_Ask, "Bid Btwn Quotes",
111            ifelse(L1.Bid <= Before_Bid, "Bid At/Below Best Bid", NA)))

```



```

112
113 ## Ask Classification
114 tmp_asks_class <- dat_NPN_oneday %>%
115   select(DateTimeL, L1.Bid, L1.Ask) %>%
116   mutate(Before_Bid = lag(L1.Bid),
117          Before_Ask = lag(L1.Ask)) %>%
118   mutate(LO_Ask_Class =
119     ifelse(L1.Ask > Before_Bid & L1.Ask < Before_Ask, "Ask Btwn Quotes",
120            ifelse(L1.Ask >= Before_Ask, "Ask At/Above Best Ask", NA)))
121
122 ## Final DataFrame with Classified Events
123 dat_NPN_oneday_classified <- dat_NPN_oneday_classified %>%
124   mutate(LO_Bid_Class = tmp_bids_class$LO_Bid_Class,
125          LO_Ask_Class = tmp_asks_class$LO_Ask_Class)
126
127 ## We need event arrival-times... subset data into events, take arrival times
128
129 ## All of the events
130 Event_vector <- c("Sell Moves Bid", "Sell Doesn't Move Bid",
131                  "Buy Moves Ask", "Buy Doesn't Move Ask",
132                  "Bid Btwn Quotes", "Bid At/Below Best Bid",
133                  "Ask Btwn Quotes", "Ask At/Above Best Ask")
134 ## Column names of data
135 Column_name_vector <- c("MO_Class", "LO_Bid_Class", "LO_Ask_Class")
136
137 ## function to subset data to get arrival times
138 func.get_events <- function(data, Event, Column_name, Type)
139 {
140   ## Subset data into event
141   dat_filtered <- data %>%
142     filter(get(Column_name) == Event)
143
144   if(Type == "Quote"){
145     ## Get Quote relevant columns
146     dat_filtered <- dat_filtered %>%
147       select(DateTimeL, Type, L1.Ask, L1.Bid, MidPrice, Volume.Ask, Volume.Bid,
148              MicroPrice)
149   } else{
150     ## Get Trade relevant columns
151     dat_filtered <- dat_filtered %>%
152       select(DateTimeL, Type, Price, Trade.Sign, Volume.Trade)
153   }
154   return(dat_filtered)
155 }
156
157 ## Initialise lists
158 MO_Classes = LO_Bid_Classes = LO_Ask_Classes = list()
159
160 ## Market Order Events are the first 4
161 for(i in 1:4){
162   MO_Classes[[i]] <- func.get_events(dat_NPN_oneday_classified,
163                                     Event_vector[i], Column_name_vector[1],
164                                     "Trade")
165 }
166 names(MO_Classes) <- Event_vector[1:4]
167
168 ## Bid Limit Order Events are 5th and 6th
169 for(i in 1:2){
170   LO_Bid_Classes[[i]] <- func.get_events(dat_NPN_oneday_classified,
171                                     Event_vector[4+i], Column_name_vector[2],
172                                     "Quote")
173 }
174 names(LO_Bid_Classes) <- Event_vector[5:6]
175
176 ## Ask Limit Order Events are 7th and 8th
177 for(i in 1:2){
178   LO_Ask_Classes[[i]] <- func.get_events(dat_NPN_oneday_classified,
179                                     Event_vector[6+i], Column_name_vector[3],
180                                     "Quote")
181 }
182 names(LO_Ask_Classes) <- Event_vector[7:8]
183
184 ## Event Arrivals Relative to Day Start

```

```

185 DayStartTime <- as.POSIXct(paste(random_date, "09:10:00", sep = " "),
186                             timezone = "Africa/Johannesburg")
187 DayEndTime <- dat_NPN_oneday[NROW(dat_NPN_oneday), ]$DateTimeL
188
189 ## Get Arrival times for each Event Type relative to day starting time
190 func.calc_arrival <- function(data)
191 {
192   dat_new <- data %>%
193     mutate(arrival_times = difftime(data$DateTimeL, DayStartTime, units = "secs"))
194
195   return(dat_new)
196 }
197
198 MO_Classes <- lapply(MO_Classes, func.calc_arrival)
199 LO_Bid_Classes <- lapply(LO_Bid_Classes, func.calc_arrival)
200 LO_Ask_Classes <- lapply(LO_Ask_Classes, func.calc_arrival)
201
202 ## Now We Want just arrival times
203 func.select_arrival <- function(data)
204 {
205   dat_new <- data %>% select(arrival_times)
206   return(dat_new)
207 }
208
209 MO_Arrivals <- lapply(MO_Classes, func.select_arrival)
210 LO_Bid_Classes <- lapply(LO_Bid_Classes, func.select_arrival)
211 LO_Ask_Classes <- lapply(LO_Ask_Classes, func.select_arrival)
212
213 # All arrival times in a single list
214 All_Arrivals <- list(MO_Arrivals[[1]], MO_Arrivals[[2]],
215                     MO_Arrivals[[3]], MO_Arrivals[[4]],
216                     LO_Bid_Classes[[1]], LO_Bid_Classes[[2]],
217                     LO_Ask_Classes[[1]], LO_Ask_Classes[[2]])
218
219 Event_No_Names_vec <- c()
220 for(i in 1:8) {Event_No_Names_vec[i] = (paste("Event_No", i, sep = ""))}
221 names(All_Arrivals) <- Event_No_Names_vec
222 # lapply(All_Arrivals, function(x) any(is.na(x))) # No NA values
223
224 ## Pull removes tibble
225 All_Arrivals_list_of_vectors <- lapply(All_Arrivals, function(x) as.vector(pull(x)))
226
227 # 2. Estimate Hawkes Process ----
228
229 k = 2
230 ## Likelihood function with constraints
231 func.likelihoodHawkes = function(params, k, columns)
232 {
233   lambda = matrix(params[1:k], nrow = k)
234   alpha = matrix(params[(k+1):(k^2 + k)], nrow = k, byrow = TRUE)
235   beta = matrix(params[(k^2 + k + 1):(k^2 + 2*k)], nrow = k)
236   eigenval = Re(eigen(diag(c(beta)) - alpha)$values)
237   if(min(eigenval) <= 0 | min(alpha) < 0 | min(lambda) < 0){
238     pen <- 1000000
239     ll <- likelihoodHawkes(lambda, alpha, beta,
240                           All_Arrivals_list_of_vectors[columns])
241     return(ll + pen)
242   }else{
243     pen <- 0
244     ll <- likelihoodHawkes(lambda, alpha, beta,
245                           All_Arrivals_list_of_vectors[columns])
246     return(ll + pen)
247   }
248 }
249
250 ## convert vector of pars into matrix
251 func.get_pars <- function(params, k)
252 {
253   lambda0 = matrix(params[1:k], nrow = k)
254   alpha = matrix(params[(k+1):(k^2 + k)], nrow = k, byrow = TRUE)
255   beta = matrix(params[(k^2 + k + 1):(k^2 + 2*k)], nrow = k)
256   return(list(lambda0, alpha, beta))
257 }

```

```

258
259 ## Horizon and true number of events
260 secs_in_day <- as.numeric(difftime(DayEndTime, DayStartTime, units = "secs"))
261 n_true <- lapply(All_Arrivals, function(x) NROW(x)) %>% unlist()
262
263 # lambda0 will be 8x1 | alpha will be 8x8 | beta will be 8x1
264 ## Market Sell and Buy (aggressive)
265 lambda_init <- lapply(All_Arrivals_list_of_vectors[c(1, 3)],
266   function(x) NROW(x)/secs_in_day) %>% unlist()
267 lambda_init_3 <- lapply(All_Arrivals_list_of_vectors[c(5, 7)],
268   function(x) NROW(x)/secs_in_day) %>% unlist()
269 lambda_init_4 <- lapply(All_Arrivals_list_of_vectors[c(1, 3, 5, 7)],
270   function(x) NROW(x)/secs_in_day) %>% unlist()
271 ## Random starting parameters >> satisfy constraints
272 set.seed(123)
273 alpha_init <- diag(0.01, k) %>% as.vector()
274 beta_init <- rep(0.05, k)
275 params_init <- as.numeric(c(lambda_init, alpha_init, beta_init))
276
277 alpha_init_3 <- diag(0.01, k) %>% as.vector()
278 beta_init_3 <- rep(0.05, k)
279 params_init_3 <- as.numeric(c(lambda_init_3, alpha_init_3, beta_init_3))
280
281 alpha_init_4 <- diag(c(0.01, 0.01, 0.01, 0.01), 4) %>% as.vector()
282 beta_init_4 <- c(0.05, 0.05, 0.05, 0.05)
283 params_init_4 <- as.numeric(c(lambda_init_4, alpha_init_4, beta_init_4))
284 # func_likelihoodHawkes(params_init, k = 2)
285 ## Market Sell and Buy (aggressive)
286 params_hawkes_optim_NM <- optim(params_init, func_likelihoodHawkes, k = 2,
287   columns = c(1, 3), method = "Nelder-Mead",
288   control = list(maxit = 5000))
289 ## Market Sell and Buy (passive)
290 params_hawkes_optim_NM_3 <- optim(params_init_3, func_likelihoodHawkes, k = 2,
291   columns = c(5, 7), method = "Nelder-Mead",
292   control = list(maxit = 5000))
293 ## Market Sell and Buy (passive) and aggressive
294 params_hawkes_optim_NM_4 <- optim(params_init_4, func_likelihoodHawkes, k = 4,
295   columns = c(1, 3, 5, 7), method = "Nelder-Mead",
296   control = list(maxit = 10000))
297
298 ## Converges
299 params_hawkes_optim_NM$convergence
300 params_hawkes_optim_NM_3$convergence
301 params_hawkes_optim_NM_4$convergence
302 ## Optimal Pars
303 optim_pars <- params_hawkes_optim_NM$par
304 optim_pars_3 <- params_hawkes_optim_NM_3$par
305 optim_pars_4 <- params_hawkes_optim_NM_4$par
306 ## Get parameter matrices
307 params_mat_hawkes <- func.get_pars(optim_pars, k)
308 params_mat_hawkes_3 <- func.get_pars(optim_pars_3, k)
309 params_mat_hawkes_4 <- func.get_pars(optim_pars_4, 4)
310 ## Clean names
311 names(params_mat_hawkes) <- c("lambda0_opt", "alpha_opt", "beta_opt")
312 names(params_mat_hawkes_3) <- c("lambda0_opt", "alpha_opt", "beta_opt")
313 names(params_mat_hawkes_4) <- c("lambda0_opt", "alpha_opt", "beta_opt")
314
315 ## Simulate(estimate) hawkes processes
316 res.sim_hawkes <- simulateHawkes(lambda0 = params_mat_hawkes$lambda0_opt,
317   alpha = params_mat_hawkes$alpha_opt,
318   beta = params_mat_hawkes$beta_opt,
319   horizon = secs_in_day)
320
321 res.sim_hawkes_3 <- simulateHawkes(lambda0 = params_mat_hawkes_3$lambda0_opt,
322   alpha = params_mat_hawkes_3$alpha_opt,
323   beta = params_mat_hawkes_3$beta_opt,
324   horizon = secs_in_day)
325
326 res.sim_hawkes_4 <- simulateHawkes(lambda0 = params_mat_hawkes_4$lambda0_opt,
327   alpha = params_mat_hawkes_4$alpha_opt,
328   beta = params_mat_hawkes_4$beta_opt,
329   horizon = secs_in_day)
330 ## Clean names

```

```

331 names(res.sim_hawkes) <- Event_No_Names_vec[c(1, 3)]
332 names(res.sim_hawkes_3) <- Event_No_Names_vec[c(5, 7)]
333 names(res.sim_hawkes_4) <- Event_No_Names_vec[c(1, 3, 5, 7)]
334
335 ## Check number of simulated vs true
336 n_sim <- lapply(res.sim_hawkes, function(x) NROW(x)) %>% unlist()
337 n_sim_3 <- lapply(res.sim_hawkes_3, function(x) NROW(x)) %>% unlist()
338 n_sim_4 <- lapply(res.sim_hawkes_4, function(x) NROW(x)) %>% unlist()
339
340 n_true[c(1, 3)]; n_sim
341 n_true[c(5, 7)]; n_sim_3
342 n_true[c(1, 3, 5, 7)]; n_sim_4
343
344 ## Calculate Intensities of events
345 func.intensity <- function(lambda_opt, alpha_opt, beta_opt, res, tend)
346 {
347
348   mu <- lambda_opt
349   mult_constant <- alpha_opt*beta_opt
350   integral <- sum(exp(-beta_opt*(res[tend] - res[c(1:(tend-1))]))
351   lambda_t <- mu + mult_constant*integral
352
353   return(lambda_t)
354 }
355
356 ## Objectives checking... Lots of local minimum in optim()
357 func.obj_calc <- function(par_mat, k, n_tests, alpha_range, beta_range,
358                           lambda_range, columns)
359 {
360
361   alpha_diags <- lapply(seq(alpha_range[1], alpha_range[2],
362                             length.out = n_tests^2), diag, nrow = 2)
363   beta_comb <- expand.grid(seq(beta_range[1], beta_range[2], length.out = n_tests),
364                             seq(beta_range[1], beta_range[2], length.out = n_tests))
365   lambda_comb <- expand.grid(seq(lambda_range[1], lambda_range[2], length.out = n_tests),
366                               seq(lambda_range[1], lambda_range[2], length.out = n_tests))
367
368   z_lambda <- numeric(dim(lambda_comb)[1])
369   z_alpha <- numeric(dim(lambda_comb)[1])
370   z_beta <- numeric(dim(lambda_comb)[1])
371
372   for(i in 1:dim(lambda_comb)[1])
373   {
374     ## Constant Lambda
375     tmp.pars <- c(as.numeric(par_mat$lambda0_opt),
376                   as.numeric(alpha_diags[[i]]),
377                   as.numeric(beta_comb[i, ]))
378     z_lambda[i] <- func_likelihoodHawkes(tmp.pars, k, columns = columns)
379     ## Constant Alpha
380     tmp.pars <- c(as.numeric(lambda_comb[i, ]),
381                   as.numeric(par_mat$alpha_opt),
382                   as.numeric(beta_comb[i, ]))
383     z_alpha[i] <- func_likelihoodHawkes(tmp.pars, k, columns = columns)
384     ## Constant Beta
385     tmp.pars <- c(as.numeric(lambda_comb[i, ]),
386                   as.numeric(alpha_diags[[i]]),
387                   as.numeric(par_mat$beta_opt))
388     z_beta[i] <- func_likelihoodHawkes(tmp.pars, k, columns = columns)
389   }
390
391   Z_lambdas <- matrix(z_lambda, nrow = n_tests)
392   Z_betas <- matrix(z_beta, nrow = n_tests)
393   Z_alphas <- matrix(z_alpha, nrow = n_tests)
394
395   return(list(Z_lambdas = Z_lambdas, Z_betas = Z_betas, Z_alphas = Z_alphas))
396 }
397
398
399 alpha_range_1 <- c(0.01, 0.04)
400 beta_range_1 <- c(0.041, 0.06)
401 lambda_range_1 <- c(0.001, 0.01)
402 test1 <- func.obj_calc(params_mat_hawkes, 2, 100,
403                       alpha_range_1, beta_range_1, lambda_range_1, c(1, 3))

```

```

404
405 alpha_range_3 <- c(0.06, 0.08)
406 beta_range_3 <- c(0.081, 0.11)
407 lambda_range_3 <- c(0.02, 0.1)
408 test3 <- func.obj_calc(params_mat_hawkes_3, 2, 100,
409                       alpha_range_3, beta_range_3, lambda_range_3, c(5, 7))
410
411 func.plot_obj <- function(n_tests, z, levels, title)
412 {
413   cols = c('green', 'yellow', 'blue')
414
415   filled.contour(seq(0, 1, length.out = n_tests), seq(0, 1, length.out = n_tests), z,
416                 plot.axes = {contour(seq(0, 1, length.out = n_tests),
417                                     seq(0, 1, length.out = n_tests),
418                                     z, add = T, nlevels = 30)}, main = title,
419                 color.palette = colorRampPalette(cols))
420 }
421
422 setwd(dir_figs)
423 cairo_pdf("HFT_Ass3_fig_objs_agg_E13_lam.pdf", height = 5, width = 5)
424 func.plot_obj(100, test1$Z_lambdas, 30, expression("Constant ~lambda"))
425 dev.off()
426
427 cairo_pdf("HFT_Ass3_fig_objs_agg_E13_beta.pdf", height = 5, width = 5)
428 func.plot_obj(100, test1$Z_betas, 30, expression("Constant ~beta"))
429 dev.off()
430
431 cairo_pdf("HFT_Ass3_fig_objs_agg_E13_alpha.pdf", height = 5, width = 5)
432 func.plot_obj(100, test1$Z_alphas, 30, expression("Constant ~alpha"))
433 dev.off()
434
435 cairo_pdf("HFT_Ass3_fig_objs_agg_E57_lam.pdf", height = 5, width = 5)
436 func.plot_obj(100, test3$Z_lambdas, 30, expression("Constant ~lambda"))
437 dev.off()
438
439 cairo_pdf("HFT_Ass3_fig_objs_agg_E57_beta.pdf", height = 5, width = 5)
440 func.plot_obj(100, test3$Z_betas, 30, expression("Constant ~beta"))
441 dev.off()
442
443 cairo_pdf("HFT_Ass3_fig_objs_agg_E57_alpha.pdf", height = 5, width = 5)
444 func.plot_obj(100, test3$Z_alphas, 30, expression("Constant ~alpha"))
445 dev.off()
446
447 func.res_sim_hawkes <- function(res, k_events, param_mat, event_no)
448 {
449
450   colour_vec <- c("black", "firebrick2", "forestgreen", "gold",
451                 "darkorange", "darkorchid3", "violetred1", "dodgerblue3")
452
453   ## Optimal Parameters per Event
454   tmp.mu <- as.numeric(param_mat$lambda0_opt)
455   tmp.alpha <- diag(param_mat$alpha_opt)
456   tmp.beta <- as.numeric(param_mat$beta_opt)
457   lambda_t <- numeric()
458   lambda_t_list <- list()
459
460   ## Calculate for each event at each arrival
461   for(event in 1:k_events){
462     for(i in 2:NROW(res[[event]])){
463       lambda_t[i] <- func.intensity(tmp.mu[event], tmp.alpha[event],
464                                   tmp.beta[event], res[[event]], i)
465     }
466     lambda_t_list[[event]] <- na.omit(lambda_t)
467     lambda_t <- numeric()
468   }
469
470   ## find max/min for plots
471   tmp_max_intensity <- max(unlist(lapply(lambda_t_list, max)))
472   tmp_min_lambdas <- min(unlist(lapply(lambda_t_list, NROW)))
473   tmp_max_lambdas <- max(unlist(lapply(lambda_t_list, NROW)))
474
475   if(k_events > 2){

```

```

477     m <- matrix(c(1,2,3,4,5,5), nrow = 3, ncol = 2, byrow = TRUE)
478     layout(mat = m, heights = c(2,2,1.5))
479   }else{
480     m <- matrix(c(1,2,3,3), nrow = 2, ncol = 2, byrow = TRUE)
481     layout(mat = m, heights = c(0.4,0.4,0.2))}
482
483   par(mar = c(2, 2, 2, 2))
484   #par(mfrow = c(1, 2))
485   for(i in 1:k_events){
486     plot(y = lambda_t_list[[i]],
487          x = seq(DayStartTime, DayEndTime, length.out = NROW(lambda_t_list[[i]])),
488          type = 'l', col = colour_vec[i],
489          ylab = expression(lambda~"(t)"), xlab = "Time",
490          main = paste("Event No.", event_no[i], sep = " "))
491   }
492   #par(mfrow = c(1, 1))
493   plot(y = c(0, tmp_max_intensity), x = c(0, tmp_max_lambdas),
494        ylab = expression(lambda~"(t)"), type = 'n',
495        xlab = '', yaxis = 'i', xaxis = 'i', xaxt="n")
496   for(i in 1:k_events){
497     lines(lambda_t_list[[i]], col = colour_vec[i], lwd = 1.5) }
498 }
499
500 setwd(dir_figs)
501 cairo_pdf("HFT_Ass3_fig_intensities_agg.pdf", height = 7, width = 7)
502 func.res_sim_hawkes(res.sim_hawkes, 2, params_mat_hawkes, c(1, 3))
503 dev.off()
504
505 cairo_pdf("HFT_Ass3_fig_intensities_pass.pdf", height = 7, width = 7)
506 func.res_sim_hawkes(res.sim_hawkes_3, 2, params_mat_hawkes_3, c(5, 7))
507 dev.off()
508
509 cairo_pdf("HFT_Ass3_fig_intensities_agg_and_pass.pdf", height = 10, width = 10)
510 func.res_sim_hawkes(res.sim_hawkes_4, 4, params_mat_hawkes_4, c(1, 3, 5, 7))
511 dev.off()
512
513 colour_vec <- c("black", "firebrick2", "forestgreen", "gold",
514                 "darkorange", "darkorchid3", "violetred1", "dodgerblue3")
515
516 ## Optimal Parameters per Event
517 tmp.mu <- as.numeric(params_mat_hawkes_4$lambda0_opt)
518 tmp.alpha <- diag(params_mat_hawkes_4$alpha_opt)
519 tmp.beta <- as.numeric(params_mat_hawkes_4$beta_opt)
520 lambda_t <- numeric()
521
522 lambda_t_list <- list()
523 ## Calculate for each event at each arrival
524 for(event in 1:length(res.sim_hawkes_4)){
525   for(i in 2:NROW(res.sim_hawkes_4[[event]])){
526     lambda_t[i] <- func.intensity(tmp.mu[event], tmp.alpha[event],
527                                   tmp.beta[event], res.sim_hawkes_4[[event]], i)
528   }
529   lambda_t_list[[event]] <- na.omit(lambda_t)
530   lambda_t <- numeric()
531 }
532
533 list_intensity <- list()
534 types_vec <- c('AggSell', 'AggBuy', 'PassBuy', 'PassSell')
535
536 for(i in 1:length(res.sim_hawkes_4)){
537   df_intensity <- data.frame(Intensity = lambda_t_list[[i]],
538                              Time = DayStartTime + res.sim_hawkes_4[[i]][-1],
539                              Count = 1:(NROW(res.sim_hawkes_4[[i]])-1),
540                              Event = types_vec[i])
541   list_intensity[[i]] <- df_intensity
542 }
543
544
545
546 setwd(dir_figs)
547 cairo_pdf("HFT_Ass3_fig_intensities_agg_wcount.pdf", height = 7, width = 7)
548 dat_plot_intensities <- do.call(rbind, list_intensity[c(1,2)]) %>% arrange(Time)
549 scaleFactor <- max(dat_plot_intensities$Intensity) / max(dat_plot_intensities$Count)

```

```

550 ggplot(dat_plot_intensities, aes(x=Time, color = Event)) +
551   geom_line(aes(y=Intensity)) +
552   geom_line(aes(y=Count * scaleFactor)) +
553   scale_y_continuous(name="Intensity",
554                      sec.axis=sec_axis(~.*1/scaleFactor, name="Count")) +
555   theme_bw()
556 dev.off()
557
558 cairo_pdf("HFT_Ass3_fig_intensities_pass_wcount.pdf", height = 7, width = 7)
559 dat_plot_intensities <- do.call(rbind, list_intensity[c(3,4)]) %>% arrange(Time)
560 scaleFactor <- max(dat_plot_intensities$Intensity) / max(dat_plot_intensities$Count)
561 ggplot(dat_plot_intensities, aes(x=Time, color = Event)) +
562   geom_line(aes(y=Intensity)) +
563   geom_line(aes(y=Count * scaleFactor)) +
564   scale_y_continuous(name="Intensity",
565                      sec.axis=sec_axis(~.*1/scaleFactor, name="Count")) +
566   theme_bw()
567 dev.off()
568
569 cairo_pdf("HFT_Ass3_fig_intensities_agg_and_pass_wcount.pdf", height = 7, width = 7)
570 dat_plot_intensities <- do.call(rbind, list_intensity) %>% arrange(Time)
571 scaleFactor <- max(dat_plot_intensities$Intensity) / max(dat_plot_intensities$Count)
572 ggplot(dat_plot_intensities, aes(x=Time, color = Event)) +
573   geom_line(aes(y=Intensity)) +
574   geom_line(aes(y=Count * scaleFactor)) +
575   scale_y_continuous(name="Intensity",
576                      sec.axis=sec_axis(~.*1/scaleFactor, name="Count")) +
577   theme_bw()
578 dev.off()
579
580 # 3. Check Exponential results of data ----
581
582 # make list of dfs..
583 test_interarrivals <- lapply(All_Arrivals_list_of_vectors[c(1, 3, 5, 7)],
584                             diff, differences = 1)
585
586 df <- list(); k_events = 4
587 for(j in 1:k_events){
588   df[[j]] <- list(inter_arrivals = test_interarrivals[[j]],
589                  alpha_opt = diag(params_mat_hawkes_4$alpha_opt)[j],
590                  beta_opt = params_mat_hawkes_4$beta_opt[j],
591                  lambda_opt = params_mat_hawkes_4$lambda0_opt[j])
592 }
593
594 names(df) <- Event_No_Names_vec[c(1, 3, 5, 7)]
595
596 func.exponential_check <- function(df)
597 {
598
599   A = c()
600   E = c()
601   A0 = 1
602   E0 = 1
603   A[1] = E0*exp(-df$beta_opt * df$inter_arrivals[1])
604   E[1] = exp (- df$beta_opt*df$inter_arrivals[1])*A[1]
605
606   for( i in 2:(length(df$inter_arrivals)) ){
607     E[i] = 1 + exp(-df$beta_opt*df$inter_arrivals[i-1])*A[i -1]
608     A[i] = E[i - 1]*exp(-df$beta_opt*df$inter_arrivals[i])
609   }
610
611   tmp = c()
612   tmp[1] = df$inter_arrivals[1]*df$lambda_opt +
613     df$alpha_opt*(1- exp(- df$beta_opt*df$inter_arrivals[1]))*E0
614
615   for( i in 2:(length(df$inter_arrivals)) ){
616     tmp[i] = df$inter_arrivals[i]*df$lambda_opt +
617     df$alpha_opt*(1- exp(- df$beta_opt*df$inter_arrivals[i-1]))*E[i -1]
618   }
619
620   Total_duration = tmp
621   return(Total_duration)
622

```

```

623 }
624
625 check_exps <- lapply(df, func.exponential_check)
626 tmp_exp_names <- c(1, 3, 5, 7)
627
628 setwd(dir_figs)
629 cairo_pdf("HFT_Ass3_fig_check_exps.pdf", height = 10, width = 10)
630 par(mfrow = c(2, 2))
631 for(i in 1:k_events){
632   main = paste("Event No.", tmp_exp_names[i], sep = " ")
633   qqplot(qexp(ppoints(length(check_exps[[i]])),
634             rate = 1/mean(check_exps[[i]]), check_exps[[i]],
635             xlab = "Exponential Quantiles", ylab = '', main = main)
636   qqline(check_exps[[i]], col = 'red')
637 }
638 dev.off()
639
640 # 4. Simulate the Order-Book ----
641
642 ## Save workspace so we can use same data...
643 setwd(dir_data)
644 save.image(file = "dat_SimHawkes.RData")
645 setwd(dir_script)
646
647 ## Initialise Midprice and spread
648 MP_init <- 100
649 spread_init <- 4
650 ## Sample initial prices and volumes using distributions???
651 sample_init_prices <- rnorm(1, 100, 20)
652 sample_init_vols <- rtrunc(1, "norm", mean = 200, sd = 300)
653
654 ##
655 simulated_arrivals <- res.sim_hawkes_4
656
657 ## Update the Order-Book
658 types_vec <- c('AggSell', 'AggBuy', 'PassBuy', 'PassSell')
659 class(simulated_arrivals$Event_No1)
660 sim_dfs <- list()
661
662 for(i in 1:length(simulated_arrivals)){
663   sim_dfs[[i]] <- data.frame(time = simulated_arrivals[[i]],
664                               type = types_vec[i])
665 }
666
667 dat_ALL_sim_arrivals <- do.call(rbind, sim_dfs) %>%
668   arrange(time)
669
670 DayLength <- NROW(dat_ALL_sim_arrivals)
671
672 sim_bid_prices <- c()
673 sim_ask_prices <- c()
674 spread_t <- c()
675 dat_true_init <- dat_NPN_oneday_classified[1, ]
676 sim_bid_prices[1] <- dat_true_init$L1.Bid
677 sim_ask_prices[1] <- dat_true_init$L1.Ask
678
679 func.sim_LOB <- function(daylength)
680 {
681   for(t in 2:daylength){
682     order <- dat_ALL_sim_arrivals[(t-1), ]
683     spread <- max(sim_ask_prices[t-1] - sim_bid_prices[t-1], 0.00001)
684     spread_t[t] <- spread
685     if(order$type == 'AggSell'){
686       tmp_neg_move = rexp(1, rate = 1) # some distribution?
687       sim_bid_prices[t] = sim_bid_prices[t-1] - tmp_neg_move
688       sim_ask_prices[t] = sim_ask_prices[t-1]
689     }else if(order$type == 'AggBuy'){
690       tmp_pos_move = rexp(1, rate = 1) # some distribution?
691       sim_bid_prices[t] = sim_bid_prices[t-1]
692       sim_ask_prices[t] = sim_ask_prices[t-1] + tmp_pos_move
693     }else if(order$type == 'PassBuy'){
694       tmp_pos_move_pass_buy = rtrunc(1, "norm", b = spread, mean = 0, sd = 100) # some
695         distribution?

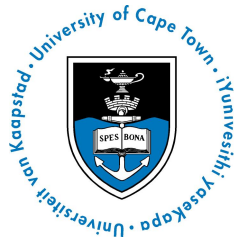
```



```

695     sim_bid_prices[t] = sim_bid_prices[t-1] + tmp_pos_move_pass_buy
696     sim_ask_prices[t] = sim_ask_prices[t-1]
697   } else if (order$type == 'PassSell') {
698     tmp_neg_move_pass_sell = rtrunc(1, "norm", b = spread, mean = 0, sd = 100) # some
699     distribution?
700     sim_bid_prices[t] = sim_bid_prices[t-1]
701     sim_ask_prices[t] = sim_ask_prices[t-1] - tmp_neg_move_pass_sell
702   }
703   return(data.frame(Sim_L1.Bid = sim_bid_prices, Sim_L1.Ask = sim_ask_prices,
704                     spread = spread_t))
705 }
706
707 test_OB <- func.sim_LOB(DayLength)
708
709 plot(y=test_OB$Sim_L1.Bid, x = (DayStartTime + dat_ALL_sim_arrivals$time),
710      type = 'l', xlab = "Time", ylab = 'Price')
711 lines(y=test_OB$Sim_L1.Ask, x = (DayStartTime + dat_ALL_sim_arrivals$time), col = 'blue')
712
713 setwd(dir_figs)
714 cairo_pdf('HFT_Ass3_fig_simulatedLOB3.pdf', height = 5, width = 5)
715 plot(y=test_OB$Sim_L1.Bid, x = (DayStartTime + dat_ALL_sim_arrivals$time),
716      type = 'l', xlab = "Time", ylab = 'Price')
717 lines(y=test_OB$Sim_L1.Ask, x = (DayStartTime + dat_ALL_sim_arrivals$time), col = 'blue')
718 dev.off()
719
720 cairo_pdf('HFT_Ass3_fig_trueLOB.pdf', height = 5, width = 5)
721 plot(na.omit(dat_NPN_oneday_classified$L1.Bid), type = 'l',
722      xlab = "Time", ylab = 'Price', x = seq(DayStartTime, DayEndTime,
723        length.out = length(na.omit(dat_NPN_oneday_classified$L1.Bid))))
724 lines(na.omit(dat_NPN_oneday_classified$L1.Ask), col = 'blue',
725       x = seq(DayStartTime, DayEndTime,
726         length.out = length(na.omit(dat_NPN_oneday_classified$L1.Ask))))
727 dev.off()

```



Plagiarism Declaration Form

A copy of this form, completed and signed, to be attached to all coursework submissions to the Statistical Sciences Department.

COURSE CODE: **STA5091Z**
COURSE NAME: **Data Analysis for High-Frequency Trading**
STUDENT NAME: **Julian Albert**
STUDENT NUMBER: **ALBJUL005**
GROUP NUMBER: **1**

PLAGIARISM DECLARATION

- I know that plagiarism is wrong. Plagiarism is to use another's work and pretend that it is one's own.
- I have used a generally accepted citation and referencing style. Each contribution to, and quotation in, this tutorial/report/project from the work(s) of other people has been attributed, and has been cited and referenced.
- This tutorial/report/project is my own work.
- I have not allowed, and will not allow, anyone to copy my work with the intention of passing it off as his or her own work.
- I acknowledge that copying someone else's assignment or essay, or part of it, is wrong, and declare that this is my own work.
- Agreement to this statement does not exonerate me from the University's plagiarism rules.

Signature:

A handwritten signature in black ink, appearing to read 'Julian Albert', written over a horizontal line.

Date: October 5, 2019