



**UNIVERSITY OF CAPE TOWN**  
IYUNIVESITHI YASEKAPA • UNIVERSITEIT VAN KAAPSTAD

# Data Analysis for High-Frequency Trading

Assignment 1

---

TAQ Data workflow and EDA

---

Julian Albert  
ALBJUL005



Department of Statistical Sciences  
University of Cape Town  
South Africa  
August 7, 2019

# Contents

<b>Introduction</b>	<b>2</b>
<b>Data</b>	<b>2</b>
<b>Initial EDA</b>	<b>3</b>
<b>Order-Flow Auto-Correlation</b>	<b>4</b>
<b>Inter-arrival Times</b>	<b>5</b>
Frequency . . . . .	5
QQ Plot . . . . .	7
ACF . . . . .	8

## Introduction

In this assignment I am tasked with creating a trade and quote (TAQ) data workflow and performing some basic exploratory data analysis (EDA). TAQ data contains intraday transactions data (trades and quotes) providing insight into a securities supply and demand that drives price, an essential part of high-frequency data analysis. The data is pulled from Bloomberg [1] such that the trades are in one .csv and the quotes are split into bid and ask with each month being a separate .csv, as such the data workflow will pull in the months of quotes and aggregate them with the trades after which the data will be processed into a user friendly form. The data workflow has been optimised using `data.tables` in R as they are more memory efficient and ultimately faster than `data.frames`. The runtime to pull in the 5 months of TAQ data, clean, classify and compact for both AGL and NPN is approximately 2 minutes.

## Data

The TAQ data pulled from [1] consists of trades and quotes for the period 2019-01-17 until 2019-06-19, each set (trades, bids and asks) contains the date, price of event and volume of the event. The bid and ask data are pulled separately and merged to keep all the data, where there is no bid (or ask) data on a specified day the data is filled with the previous non-missing values. The trades consist of larger trades executed against collections of smaller trades, this sequence is still a single trade of a given volume and as such must be *compacted* to reflect it as such. In the data we see two event types (quote and trade), where the data displays same time quotes - sequences of quotes with the same date-time stamps - the quote most recent in sequence quote is retained and the remainder removed. For same time trades - sequences of trades with the same date-time stamps - as well as the same order types (buyer or seller initiated) the sequence of trades is aggregated into a single trade with a single date-time stamp but aggregated volume and the trade price is set to be the volume weighted average price (VWAP). The assumption here is that orders of different sign +1(-1) are not aggregated. To identify order type we must *classify* the trades using the methodology, as originally developed by [3], involving a two-step approach defined by a quote test followed by a tick test

1. In the Quote Test, if the price of a trade is lower than the midpoint of the matched bid and ask, then the trade is classified as seller-initiated. If the price is higher, the trade is considered to be buyer-initiated.
2. When they are equal, the tick test is utilized. In the Tick Test, if the price is lower than the previous price, it is classified as seller-initiated and if it is higher, then the trade is considered to be buyer-initiated.

If the tick test results are not conclusive the trades will be classified as undetermined.

In order to carry out basic EDA I must first prepare the dataset to include some useful measures. The complete dataset must be time ordered and include: date, time, mid-price, mid-price change, volume, normalised volume, bid, bid volume, ask, ask volume, event type, transaction price, transaction volume and volume-weighted average prices (VWAP). To calculate some of these measures we can use the definitions in Table 1. The mid-price represents the mid-point of the current bid-ask spread quoted for a security and acts as a proxy for the trade price and the mid-price change is the change in quote due to a transaction event (take to be log differences to preserve ergodicity). Where the mid-price may be unrealistic when the volume of limit orders

at the best bid and ask prices differ significantly the micro-price is considered as it represents a volume weighted mid-price.

Measure	Formulae
Mid-price ( $m_{ij}$ )	$m_{ij} = \frac{1}{2} (b_{ij} + a_{ij})$
Mid-price Change ( $\Delta m_{ij}$ )	$\Delta m_{ij} = \log \left( \frac{m_{i+1,j}}{m_{i,j}} \right)$
Micro-price ( $S_t$ )	$S_t = \left( \frac{v_{ij}^a}{v_{ij}^a + v_{ij}^b} \right) a_{ij} + \left( \frac{v_{ij}^b}{v_{ij}^a + v_{ij}^b} \right) b_{ij}$
Inter-arrival Times ( $t_{ij}$ )	$t_{ij} = t_{i+1,j} - t_{ij}$
Normalised Trade Volume ( $\omega_{ij}$ )	$\omega_{ij} = \frac{v_{ij}}{\sum_{k=1}^{T_j} v_{kj}} \left( \frac{\sum_{j=1}^N T_j}{N} \right)$
VWAP ( $vwap_{t+1,t}$ )	$vwap_{t+1,t} = \frac{\sum_{i=t}^{t+1} v_i p_i}{\sum_{i=t}^{t+1} v_i}$

Table 1: Useful Definitions for Dataset Measures.

## Initial EDA

With the prepared dataset this section will visualise an hour of trading on the two different stocks between 10h00 and 11h00, and 16h00 and 17h00. This visualisation is achieved by plotting the micro-price (in green) in calendar time, the bid and (offer) (best-bid and best-offer) order book events as blue (and red) bubbles that are proportional to the volume on bid (and on offer) and transactions in yellow at the transaction event-times as bubbles proportional to the transaction volumes. For the visualisation a random trading day is selected (a day where a trade certainly occurred) and the a data subset is taken for a specified trading interval, the resulting plots can be seen in Figure 1.

Starting with AGL we can visualise the two time intervals, starting at 10:00am we see the hourly high proceeded by big bid volumes driving the mid-price (and ultimately the transaction and micro-price) downwards, this is met with large transaction volumes between 10:15 and 10:20 from which point on the data seems more stable. Starting at 16:00pm a week later we see the price is higher and rallying up, the price increase is not evidently due to an excess of buyers or sellers driving the price. Curiously at approximately 16:31 we see a large transaction volume (possibly at the best-ask at the time but given the transaction bubble clearly covers any of the best ask bubble) we may determine this to be an error in the data-feed. This data point obscures the visual analysis of this graph.

Moving on to NPN at 10:00am we see a fairly consistent rise and fall over the hour with high asks at around 10:37 and a big transaction at approximately 10:56am. More interesting is NPN a month later, there is a big increase in the price. During the interval from 16:00 until 16:49:59 we see a large bid-ask spread with the ask volumes being distinctly larger than those of the bids

(there are more sellers than buyers). The bid-ask spread may be much wider due to increased volatility prompting market makers to profit from the spread. When securities are increasing in value, investors are willing to pay more giving market makers the opportunity to charge higher premiums. When volatility is low, and uncertainty and risk are at a minimum, the bid-ask spread is narrow.

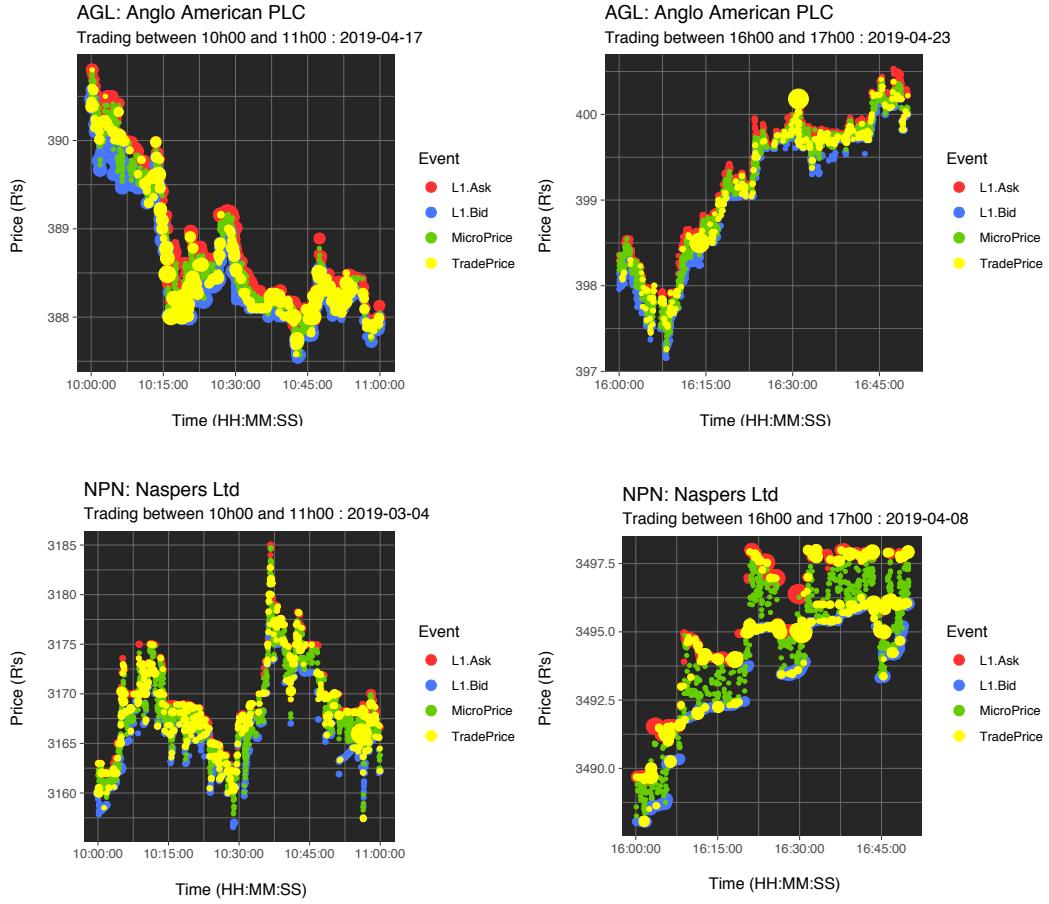


Figure 1: EDA of 1-hour TAQ data.

## Order-Flow Auto-Correlation

Order flow can be defined as the sequence/process which assumes value +1 for buyer initiated trades and -1 for seller initiated trades. Order flow is persistent in the sense that orders to buy tend to be followed by more orders to buy and orders to sell tend to be followed by more orders to sell [4]. Figure 2 illustrates that the autocorrelation function for market order signs decays very slowly and that the autocorrelation function is roughly linear in a double logarithmic scale suggesting that a power law relation might be a reasonable description

for the sample autocorrelation function [2].

Using the `ACF()` function in R, the auto-correlations for both the trade signs and inter-arrival times are calculated and plotted in Figure 2 and 6 respectively. The auto-correlations are plotted on the log-log scale, as a result the transformation sometimes sets data to null (hence the discontinuities).

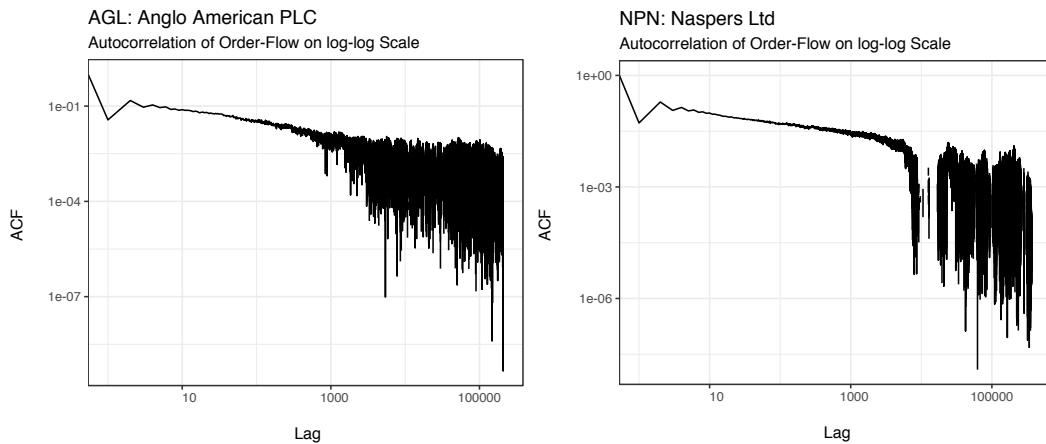


Figure 2: Auto-Correlation of Order Flow for Two Equities.

Autocorrelation is an important feature to understand, if a time series is autocorrelated we can anticipate future states of the series, thus if there were auto-correlations for the trade-signs, high-frequency traders could predict whether subsequent traders will be seller or buyer initiated based on the most recent trade signs.

## Inter-arrival Times

Trade inter-arrival times  $\Delta t_{ij}$  represent the time between trades and provide useful insight into the order flow and liquidity of a security. In this section I will consider the frequency of the inter-arrival times, a quantile-quantile analysis and the auto-correlations.

### Frequency

We may be interested to know the distribution of inter-arrival times  $\Delta t_{ij}$  as this will give us a sense of how “fast” the security is trading. In Figure 3 we plot the histogram of inter-arrival times over the entire set of trading events for the 135 trading days. We can see that the amount of trading events are higher for Naspers Ltd compared to AGL PLC. In both cases we observe a right skew distribution indicating that the trades occur more frequently together (a sign of higher liquidity), the left column is on the absolute scale where the skew is heavy thus we apply a log scale to get a better sense of the behaviour of inter-arrival times. We can see a subtle difference in that Naspers Ltd is more right skewed indicating more liquidity.

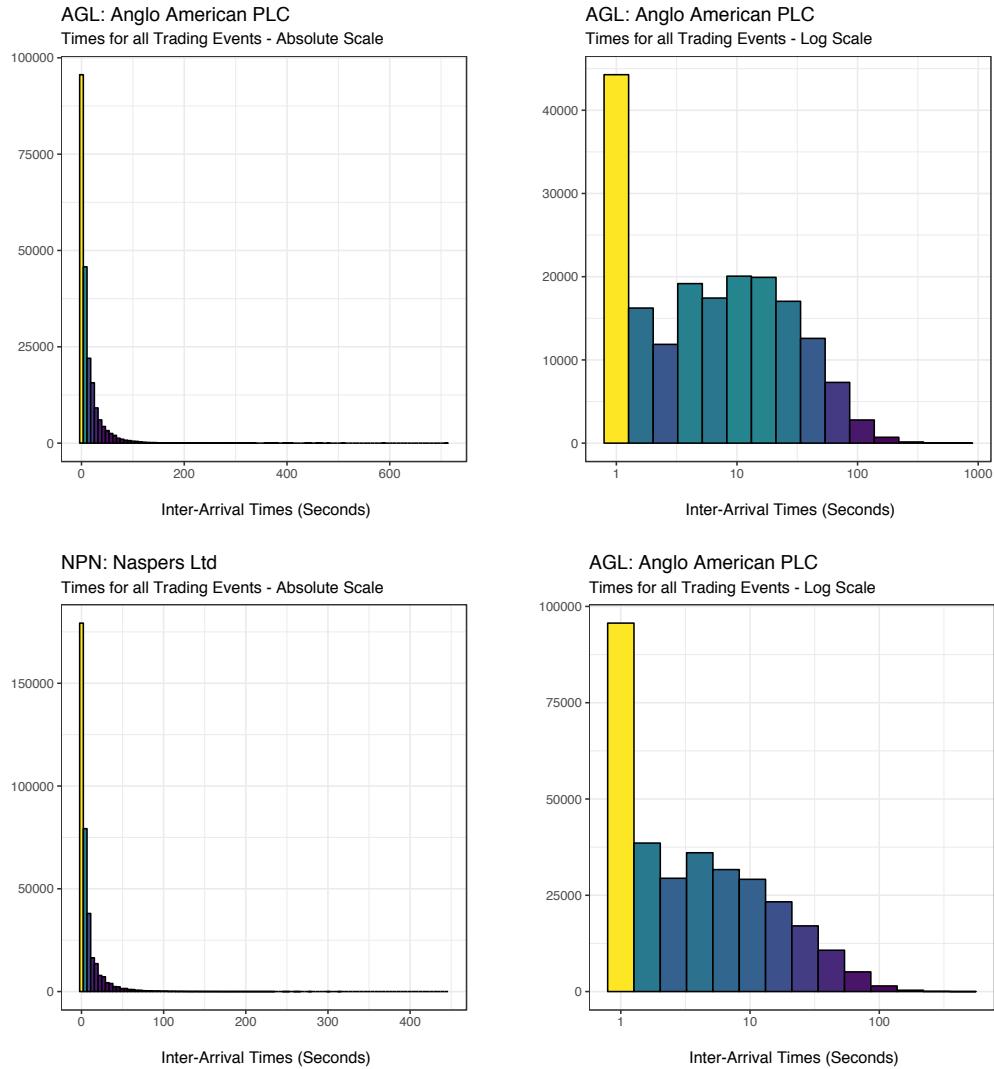


Figure 3: Frequency of Inter-Arrival Times in Absolute and Log Scales

To further analyse the trade frequency we can plot the frequency of inter-arrival times explicitly less than 10 seconds, 1 minute and 10 minutes. The resulting plots are seen in Figure 4 where clearly the bulk of the mass is captured in 1 second whilst we can determine that the inter-arrival times are majority less than 1 minute for both securities indicating that they are both relatively liquid. Interestingly the high frequency of inter-arrival times around 2 seconds fits well with the exponential distribution up to 2 seconds (see the next section and Figure 5).

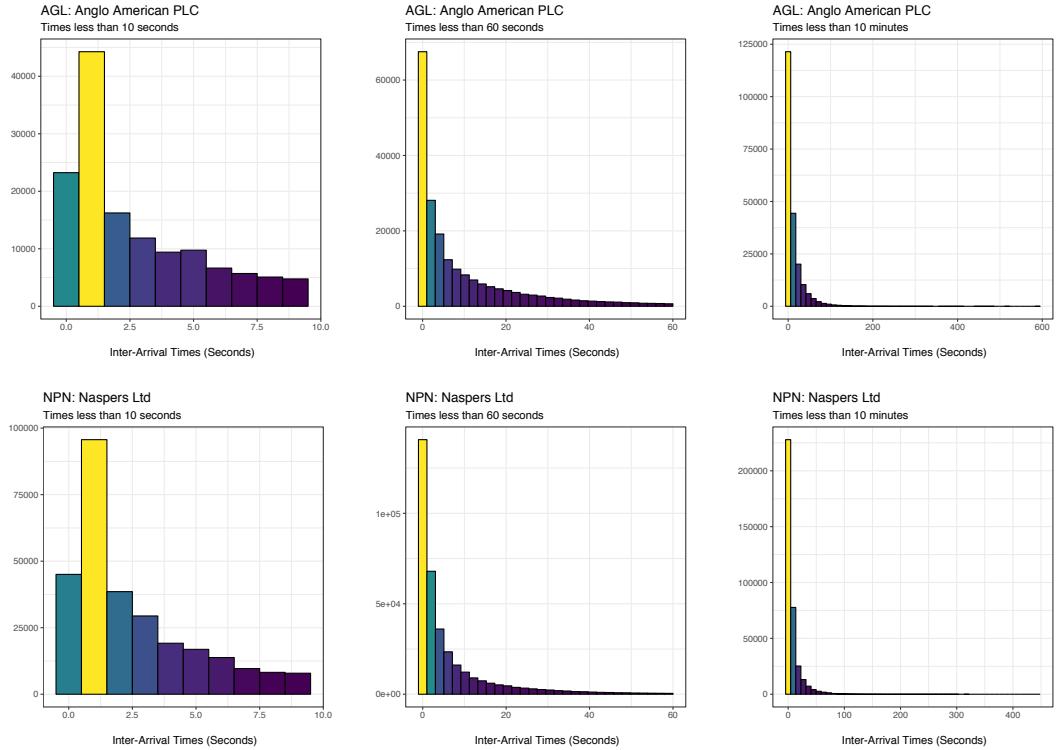


Figure 4: Frequency of Inter-Arrival Times less than 10 seconds, 1 minute and 10 minutes.

## QQ Plot

The shape of the frequency for inter-arrival times in Figure 3 and 4 appear Exponential in form. To check this we can consider the quantile-quantile plot (QQ Plot) of the empirical quantiles obtained from the inter-arrival times against the theoretical quantiles of the Exponential distribution. The QQ plot is illustrated in Figure 5 from which we can see that the inter-arrival times of trades are not Exponentially distributed (as a good fit is one in which the points lie along the 45 degree reference line through the origin). For inter-arrival times less than approximately 2.5 seconds the Exponential fits relatively well, however, for trade inter-arrivals greater than 2.5 seconds the empirical quantiles are larger than the empirical quantiles indicating that the data has a fatter right tail than the Exponential distribution. This phenomenon of leptokurtosis is a known real world problem of financial data, the fat tail suggests that inter-arrival times of trades decay slower than the Exponential distribution.

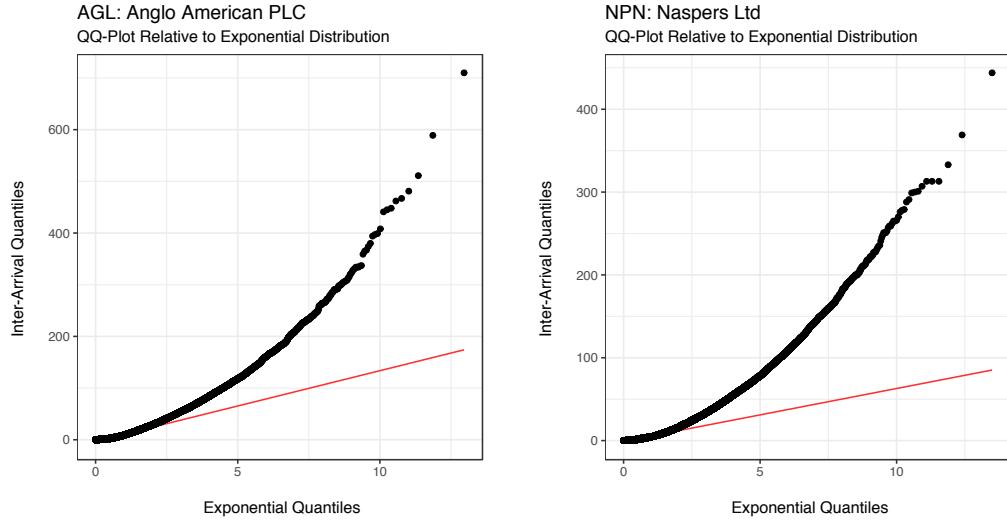


Figure 5: QQ Plots of Inter-Arrival Times Relative to Exponential Distribution.

## ACF

If the inter-arrival times were autocorrelated, then we could anticipate the time of the next trade using recent historical inter-arrival times. We can see that the auto-correlations in the inter-arrival times decay slowly. This result indicates that if one observes a market order now, based on this information alone there is some nonvanishing predictability of the market order arrival times [2].

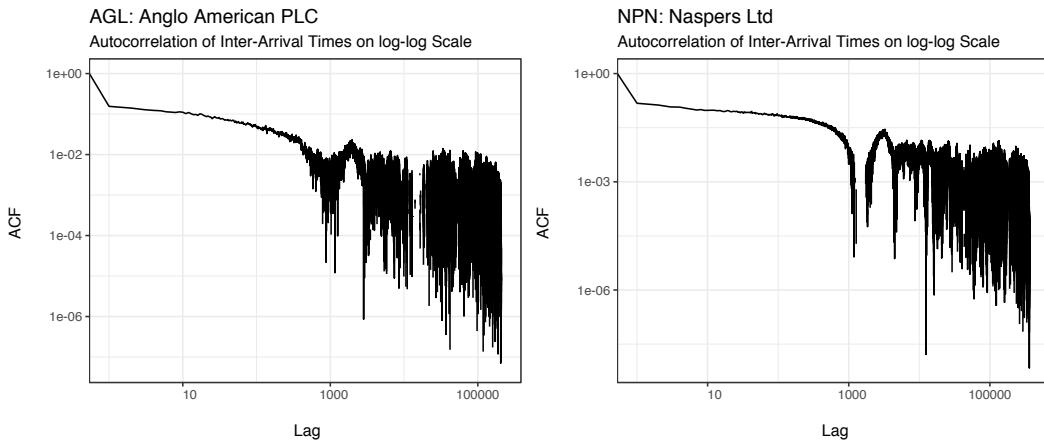


Figure 6: Auto-Correlation of Inter-Arrival Times.

- END -

## References

- [1] BLOOMBERG DATABASE. *Bloomberg Professional*. Bloomberg Subscription Service, Online, 2019.
- [2] BOUCHAUD, J.-P., FARMER, J. D., AND LILLO, F. How markets slowly digest changes in supply and demand. In *Handbook of financial markets: dynamics and evolution*. Elsevier, 2009, pp. 57–160.
- [3] LEE, C. M., AND READY, M. J. Inferring trade direction from intraday data. *The Journal of Finance* 46, 2 (1991), 733–746.
- [4] TOTH, B., PALIT, I., LILLO, F., AND FARMER, J. D. Why is equity order flow so persistent? *Journal of Economic Dynamics and Control* 51 (2015), 218–239.

# Appendix

## R Code

Listing 1: Data Pull From Bloomberg

```
1 ##  
2 #  
3 # Author: Julian Albert  
4 # Date: 01 August 2019  
5 #  
6 # Description:  
7 # HFT Data Pull from Bloomberg  
8 #  
9 #-----#  
10  
11  
12 # 0. Clean Workspace, Directory and Library ----  
13  
14 ## Clean Workspace  
15 rm(list=ls())  
16 dev.off() # Close Plots  
17  
18 ## Packages for Bloomberg PC  
19 # Some might/might not install  
20 install.packages("Rblpapi")  
21 install.packages("lubridate")  
22 library(Rblpapi)  
23 library(lubridate)  
24  
25 ## Bloomberg Connection  
26 con <- blpConnect()  
27  
28 ## Directory  
29 dir <- "G:/HFT/Data/"  
30  
31 # 1. Variables for Pulling ----  
32  
33 # security <- c("APN SJ EQUITY", "CPI SJ EQUITY", "BIL SJ EQUITY", "INP SJ EQUITY",  
34 #           "AGL SJ EQUITY", "TFG SJ EQUITY", "BTI SJ EQUITY", "MTN SJ EQUITY",  
35 #           "NPN SJ EQUITY", "SHP SJ EQUITY")  
36  
37 ## Securities to pull  
38 security <- c("AGL SJ EQUITY", "NPN SJ EQUITY")  
39 n.securities <- length(security)  
40 ## Timezones  
41 timezone <- Sys.timezone()  
42 StartTime <- as.POSIXct("2019-01-01 08:00:00", tz = timezone)  
43 CurrentTime <- Sys.time()  
44 n.months <- interval(StartTime, test1) %/% months(1)  
45  
46 # 2. Pulling data ----  
47  
48 ## Daily Sampled OHLCV  
49 barInt <- 540  
50  
51 daily_sampled <- list()  
52 for (i in 1:n.months){  
53   daily_sampled[[i]] <- getBars(security = security[i], barInterval = 540,  
54                                startTime = StartTime, tz = timezone, con = con)  
55   filen <- paste(dir, "Daily/", substr(security[i], 1, 3), ".csv", sep = "")  
56   write.csv(daily_sampled[[i]], filen, row.names = F)  
57 }  
58  
59 getBars(security = security[1], barInterval = 540,  
60          startTime = StartTime, tz = timezone, con = con)  
61  
62 ## Intraday  
63 StartTime <- as.POSIXct("2019-01-01 08:00:00", tz = timezone)  
64  
65 ### Trade data
```

```

66 for (i in 1:n.months){
67   filen_trades = paste(dir, "Intraday/trade_",
68                         substr(security[i], 1, 3), ".csv", sep = "")
69   temp = getTicks(security = security[i], eventType = c("TRADE"),
70                   startTime = StartTime, tz = timezone, con = con)
71   write.csv(temp, filen_trades, row.names = FALSE)
72 }
73
74 ### Ask data
75 for (i in 1:n.securities){
76   StartTime = StartTime
77   EndTime = StartTime
78   month(EndTime) <- month(StartTime) + 1
79   for (k in 1:6){
80     filen_ask <- paste(dir,"Intraday/asks/ask_",
81                           substr(security[i], 1, 3), "00", k, ".csv", sep = "")
82     tmp_filen <- getTicks(security = security[i], eventType = c("BEST_ASK"),
83                           startTime = StartTime, endTime = EndTime,
84                           tz = timezone, con = con)
85     write.csv(tmp_filen, filen_ask, row.names = FALSE)
86     StartTime <- EndTime
87     month(EndTime) <- month(EndTime) + 1
88   }
89 }
90
91 ### Bid data
92 for (i in 1:n.securities){
93   StartTime <- StartTime
94   EndTime <- StartTime
95   month(EndTime) <- month(StartTime) + 1
96   for (k in 1:n.months){
97     filen_bid <- paste(dir, "Intraday/bids/bid_",
98                           substr(security[i], 1, 3), "00", k, ".csv", sep = "")
99     tmp_filen = getTicks(security = security[i], eventType = c("BEST_BID"),
100                           startTime = StartTime, endTime = EndTime,
101                           tz = timezone, con = con)
102     write.csv(tmp_filen, filen_bid, row.names = FALSE)
103     StartTime <- EndTime
104     month(EndTime) <- month(EndTime) + 1
105   }
106 }
107
108 #-----#

```

Listing 2: Data Workflow - Clean, Classify and Compact

```

1 ##
2 #
3 # Author: Julian Albert
4 # Date: 01 August 2019
5 #
6 # Description:
7 # HFT Assignment 1 - Data Load and Clean
8 #
9 #-----#
10
11 # O. Clean Workspace, Directory and Library ----
12
13 ## Clean Workspace
14 rm(list=ls())
15 dev.off() # Close Plots
16 setwd("~/") # Clear Path to User
17
18 ## Locations
19 project_folder <- "/Documents/UCT/Coursework/HFT"
20 loc_script <- "/Assignment_1/UCT_Assignment/Code"
21 loc_data <- "/Data"
22
23 ## Directories
24 dir_script <- paste("~/", project_folder, loc_script, sep = ' ')
25 dir_data <- paste("~/", project_folder, loc_data, sep = ' ')
26
27 ## Set Working Directory to Script Location

```

```

28| setwd(dir_script)
29|
30## Libraries - Lazyload
31if (!require("pacman")) install.packages("pacman")
32p_load(tidyverse, openxlsx, lubridate, zoo, data.table)
33|
34## Options
35options(digits.secs = 3)
36|
37# 1. Load in data ----
38|
39## Variables
40securities <- c("AGL", "NPN")
41n.securities <- length(securities)
42n.months <- 5
43timezone <- "Africa/Johannesburg"
44dat_TAQ_final <- list()
45|
46# 2. Cleaning Everything ----
47|
48## Start Time
49ptm <- proc.time()
50|
51## For each security
52for(security in 1:length(securities)){
53|
54    # Quotes are in multiple .csv by month >> read into list
55    dat.quotes_CLEAN <- list()
56|
57    # Per month into one list
58    for(month in 1:n.months){
59|
60        ## Filenames for Bid and Ask
61        filen_ask <- paste(dir_data, "/ask_",
62                            substr(securities[security], 1, 3), "00", month, ".csv",
63                            sep = " ")
64|
65        filen_bid <- paste(dir_data, "/bid_",
66                            substr(securities[security], 1, 3), "00", month, ".csv",
67                            sep = " ")
68|
69        ## Read Data for Bid and Ask
70        tmp.dat_ask <- fread(filen_ask, drop = "concode")
71        tmp.dat_bid <- fread(filen_bid, drop = "concode")
72|
73        ## Join the Data and put into Desired Format >> X = Ask, Y = Bid
74        tmp.dat_quotes <- full_join(tmp.dat_ask, tmp.dat_bid, by = "times") %>%
75            rename(DateTimeL = times,
76                  L1.Ask = value.x, Volume.Ask = size.x,
77                  L1.Bid = value.y, Volume.Bid = size.y) %>%
78            arrange(DateTimeL) %>%
79            mutate(Type = "Quote", type.x = NULL, type.y = NULL) %>%
80            setDT()
81|
82        ## Clean the quotes >> Only take Last Quote >> Do now to reduce size ASAP
83        dat.quotes_CLEAN[[month]] <- tmp.dat_quotes[, .SD[.N], DateTimeL]
84|
85    }
86|
87    ## Remove the Auction >> Set DateTime Format >> Fill NA Asks/Bids
88    df_Quals <- dat.quotes_CLEAN %>%
89        bind_rows() %>%
90        [, DateTimeL := parse_date_time(DateTimeL, "Ymd HMS", tz = timezone)] %>%
91        setDF() %>%
92        mutate(hms = format(as.POSIXct(DateTimeL), "%H:%M:%S")) %>%
93        [.hms >= "09:10:00" & .hms <= "16:49:59", ] %>%
94        mutate(hms = NULL, Price = NA, Volume.Trade = NA) %>%
95        fill(L1.Ask, Volume.Ask, L1.Bid, Volume.Bid) %>%
96        as_tibble()
97|
98    ## Read in the trades (only one .csv with all trades)
99    filen_trades <- paste(dir_data, "/trade_",
100

```

```

101             securities[security], ".csv",
102             sep = "")
103
104 ## Remove the Auction
105 tmp.dat.trades <- fread(filen_trades, drop = c("condcode", "type")) %>%
106   rename(DateTimeL = times, Price.Trade = value, Volume.Trade = size) %>%
107   .[, DateTimeL := parse_date_time(DateTimeL, "Ymd_HMS", tz = timezone)] %>%
108   mutate(Type = "Trade", hms = format(as.POSIXct(DateTimeL), "%H:%M:%S")) %>%
109   .[$hms >= "09:10:00" & .$hms <= "16:49:59", ] %>%
110   mutate(hms = NULL) %>%
111   setDT()
112
113 ## CLEAN TRADES
114 df_Trades <- tibble(DateTimeL = tmp.dat.trades$DateTimeL,
115                       L1.Ask = NA, Volume.Ask = NA,
116                       L1.Bid = NA, Volume.Bid = NA,
117                       Type = tmp.dat.trades$type,
118                       Price = tmp.dat.trades$Price,
119                       Volume.Trade = tmp.dat.trades$Volume.Trade)
120
121 ## TAQ data is Trades and Quotes
122 df_TAQ <- bind_rows(df_Quotes, df_Trades) %>%
123   arrange(DateTimeL) %>%
124   setDT() %>%
125   .[, MidPrice := 0.5*(L1.Ask + L1.Bid)] %>%
126   fill(MidPrice) %>%
127   mutate(QuoteRule = ifelse(Price < MidPrice, -1,
128                           ifelse(Price > MidPrice, 1, 0)),
129         tmpdiff = c(NA, diff(Price)),
130         TickRule = ifelse(tmpdiff > 0, 1,
131                           ifelse(tmpdiff < 0, -1, 0)),
132         Trade.Sign = ifelse(!is.na(QuoteRule) & QuoteRule != 0, QuoteRule,
133                           ifelse(QuoteRule == 0, TickRule, NA)))
134
135 df_TAQ_CLEAN <- df_TAQ %>%
136   mutate(QuoteRule = NULL, tmpdiff = NULL, TickRule = NULL)
137
138 ## Calculate measures on the Trades
139 dat_transactions <- df_TAQ_CLEAN %>% filter(Type == "Trade") %>% setDT()
140
141 dat_transactions <- dat_transactions %>%
142   .[, keyby = .(DateTimeL, Trade.Sign),
143     c("VWAP", "TotVolume") :=
144       .(weighted.mean(Price, Volume.Trade),
145         sum(Volume.Trade))] %>% # Calc VWAP and sum(Volume)/Trade Type/day
146   .[, keyby = .(DateTimeL, Trade.Sign), .SD[.N]] %>%
147   .[, by = floor_date(DateTimeL, "day"),
148     c("DailyPropVol", "No.DailyTrades", "day", "Inter.Arrival.times") :=
149       .(TotVolume/sum(TotVolume), .N,
150         floor_date(DateTimeL, "day"),
151         c(NA, diff(DateTimeL)))] %>% # Calc. Total Volume on Day, no Trades, t
152   .[, c("Price", "Volume.Trade", "VWAP", "TotVolume",
153         "TotalTradingDays") :=
154     .(VWAP, TotVolume, NULL, NULL, uniqueN(day))] %>% # Clean and get N
155   .[, NormalisedVolume :=
156     DailyPropVol*(No.DailyTrades/TotalTradingDays)] %>% # Normalised Volume
157   .[, c("day", "MidPrice", "Mid.Price.Change", "MicroPrice") :=
158     .(NULL, NA, NA, NA)] # Clean
159
160 ## Calculate measures on the Quotes
161 dat_quotes <- df_TAQ_CLEAN %>% filter(Type == "Quote")%>% setDT()
162
163 dat_quotes <- dat_quotes %>%
164   .[, .SD[.N], by = DateTimeL] %>%
165   .[, c("DailyPropVol", "No.DailyTrades", "TotalTradingDays",
166         "NormalisedVolume", "Inter.Arrival.times", "MidPrice",
167         "Mid.Price.Change", "MicroPrice") :=
168     .(NA, NA, NA, NA, NA, 0.5*(L1.Bid + L1.Ask), c(NA, diff(log(MidPrice))),
169       (Volume.Ask/(Volume.Ask+Volume.Bid))*L1.Ask +
170       (Volume.Bid/(Volume.Ask+Volume.Bid))*L1.Bid)]
171
172 ## Reorder into a Nice Format
173 tmp_order_names <- c("DateTimeL", "Type", "L1.Ask", "L1.Bid", "MidPrice",

```

```

174     "Price", "Trade.Sign", "Volume.Ask", "Volume.Bid",
175     "Volume.Trade", "NormalisedVolume", "MicroPrice",
176     "No.DailyTrades", "TotalTradingDays", "DailyPropVol",
177     "Inter.Arrival.times", "Mid.Price.Change")
178
179 setcolorder(dat_transactions, tmp_order_names)
180 setcolorder(dat_quotes, tmp_order_names)
181
182 ## Some Reason dates are mismatched, start from min(date) end at min(lastdate)
183 first_date <- min(first(dat_quotes$DateTimeL),
184                     first(dat_transactions$DateTimeL))
185
186 last_date <- min(last(dat_quotes$DateTimeL),
187                     last(dat_transactions$DateTimeL))
188
189 ## Clean and Compact Data >> Store Per Ticker
190 dat_TAQ_final[[security]] <- bind_rows(dat_transactions, dat_quotes) %>%
191   arrange(DateTimeL) %>%
192   setDT() %>% # data with everything needed
193   .[DateTimeL %between% c(first_date, last_date)] %>%
194   as_tibble()
195
196 }
197
198 names(dat_TAQ_final) <- securities
199 proc.time() - ptm
200
201
202 # 3. Save Cleaned, Classified, Compacted Dataset ----
203 saveRDS(dat_TAQ_final, file = "dat_TAQ.rds")
204
205 #-----#

```

Listing 3: EDA with Processed Data

```

1 ##
2 #
3 # Author: Julian Albert
4 # Date: 02 August 2019
5 #
6 # Description:
7 # HFT Assignment 1 - Basically want to work with TAQ data and perform EDA to
8 # better understand the data and the trading environment.
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_1/UCT_Assignment/Code"
22 loc_figs <- "/Assignment_1/UCT_Assignment/Figs"
23
24 ## Directories
25 dir_script <- paste("~/", project_folder, loc_script, sep = ' ')
26 dir_figs <- paste("~/", project_folder, loc_figs, sep = ' ')
27
28 ## Set Working Directory to Script Location
29 setwd(dir_script)
30
31 ## Libraries - Lazyload
32 if (!require("pacman")) install.packages("pacman")
33 p_load(tidyverse, lubridate, zoo, data.table, Cairo, viridis)
34
35 ## Pull Clean Data
36 dat_TAQ <- readRDS(file = "dat_TAQ.rds")
37
38 ## Options

```

```

39 | options(digits.secs = 3)
40 |
41 | # 1. Create Plots >> Trades/Quotes/MicroPrice ----
42 |
43 | StartTime1 <- "10:00:00"
44 | EndTime1 <- "11:00:00"
45 | StartTime2 <- "16:00:00"
46 | EndTime2 <- "17:00:00"
47 | timezone <- "Africa/Johannesburg"
48 |
49 | Stock1 <- "AGL: Anglo American PLC"
50 | Stock2 <- "NPN: Naspers Ltd"
51 |
52 | func.plot_EDA_1hour <- function(Data, StartTime, EndTime, Stock, Legend = FALSE)
53 | {
54 |
55 |   ## Randomly Select a Trading Day
56 |   tmp_random_date <- Data %>%
57 |     filter(Type == 'Trade') %>%
58 |     slice(runif(1, 1, NROW(.))) %>%
59 |     dplyr::select(DateTimeL)
60 |
61 |   ## Define the Day
62 |   tmp_ymd <- paste(year(tmp_random_date$DateTimeL),
63 |                      month(tmp_random_date$DateTimeL),
64 |                      day(tmp_random_date$DateTimeL), sep = "-")
65 |
66 |   ## Need to subset for the Trading Hours
67 |   date_start <- as.POSIXct(paste(tmp_ymd, StartTime, sep = " "), tz = timezone)
68 |   date_end <- as.POSIXct(paste(tmp_ymd, EndTime, sep = " "), tz = timezone)
69 |   plot_interval <- interval(date_start, date_end)
70 |
71 |   dat_hourly <- Data[Data$DateTimeL %within% plot_interval, ]
72 |
73 |   ## Data in nice plotting format >> Long
74 |   testprice <- dat_hourly %>%
75 |     rename(TradePrice = Price) %>%
76 |     gather(key, price, c("L1.Ask", "L1.Bid", "MicroPrice", "TradePrice")) %>%
77 |     select(DateTimeL, key, price)
78 |
79 |   testvols_quotes <- dat_hourly %>%
80 |     gather(key, volume, Volume.Ask:Volume.Bid) %>%
81 |     select(DateTimeL, key, volume)
82 |
83 |   testvols_mp <- data.frame(DateTimeL = dat_hourly$DateTimeL,
84 |                                key = "Volume.MicroPrice",
85 |                                volume = ifelse(!is.na(dat_hourly$L1.Ask),
86 |                                               0.5*min(testvols_quotes$volume,
87 |                                                       na.rm = TRUE), NA)) %>%
88 |     mutate(key = as.character(key))
89 |
90 |   testvols_trades <- dat_hourly %>%
91 |     gather(key, volume, Volume.Trade) %>%
92 |     select(DateTimeL, key, volume)
93 |
94 |   testvols <- bind_rows(testvols_quotes, testvols_mp, testvols_trades)
95 |
96 |   test_comb <- data.frame(DateTimeL = testprice$DateTimeL,
97 |                             Event = as.factor(testprice$key),
98 |                             Value = testprice$price,
99 |                             Volume = testvols$volume)
100 |
101 |   ## Generic Titles
102 |   time.title <- paste("Trading between ",
103 |                        unlist(strsplit(StartTime, ":"))[1],
104 |                        "h", unlist(strsplit(StartTime, ":"))[2], " and ",
105 |                        unlist(strsplit(EndTime, ":"))[1],
106 |                        "h", unlist(strsplit(EndTime, ":"))[2], sep = "")
107 |
108 |   date.title <- date(date_start)
109 |
110 |   ## Plot EDA
111 |

```

```

112 plot <- test_comb %>% na.omit() %>%
113   ggplot(aes(x = DateTimeL, y = Value/100, size = Volume, color = Event)) +
114   geom_point() +
115   theme_dark() +
116   scale_color_manual(breaks = c("L1.Ask", "L1.Bid",
117                               "MicroPrice", "TradePrice"),
118                       values=c("firebrick1", "royalblue1",
119                               "chartreuse3", "yellow")) +
120   labs(x = "\n Time (HH:MM:SS)", y = "Price (R's) \n", title = Stock,
121         subtitle = paste(time.title, date.title, sep = " : ")) +
122   scale_x_datetime(date_labels = "%H:%M:%S",
123                     timezone = "Africa/Johannesburg") +
124   theme(panel.background = element_rect(fill = 'grey15')) +
125   theme(legend.key=element_blank()) +
126   scale_size_continuous(guide = FALSE) +
127   guides(color = guide_legend(override.aes = list(size = 3)))
128
129 if(Legend == TRUE) {plot <- plot + theme(aspect.ratio=1)}
130 else{plot <- plot + theme(legend.position ="none") + theme(aspect.ratio=1)}
131
132 return(plot)
133
134 }
135
136 setwd(dir_figs)
137
138 # cairo_pdf("HTF_Ass1_fig_EDA_AGL_1011.pdf", height = 5, width = 5)
139 func.plot_EDA_1hour(dat_TAQ$AGL , StartTime1, EndTime1, Stock1, Legend = TRUE)
140 # dev.off()
141
142 # cairo_pdf("HTF_Ass1_fig_EDA_AGL_1617.pdf", height = 5, width = 5)
143 func.plot_EDA_1hour(dat_TAQ$AGL , StartTime2, EndTime2, Stock1, Legend = TRUE)
144 # dev.off()
145
146 # cairo_pdf("HTF_Ass1_fig_EDA_NPN_1011.pdf", height = 5, width = 5)
147 func.plot_EDA_1hour(dat_TAQ$NPN , StartTime1, EndTime1, Stock2, Legend = TRUE)
148 # dev.off()
149
150 # cairo_pdf("HTF_Ass1_fig_EDA_NPN_1617.pdf", height = 5, width = 5)
151 func.plot_EDA_1hour(dat_TAQ$NPN , StartTime2, EndTime2, Stock2, Legend = TRUE)
152 # dev.off()
153
154 setwd(dir_script)
155
156 # 2. Create Plots >> ACF ----
157
158 ## Get TAQ that has +1(-1) Trade Signs only
159 dat_tradesigns <- lapply(dat_TAQ, function(x) x %>%
160                           setDT() %>% .[!is.na(Trade.Sign) & Trade.Sign != 0])
161
162 ## Total Length for N days
163 counts <- lapply(dat_tradesigns, function(x) x %>%
164                   setDT() %>% .[, length(Trade.Sign)])
165
166 ## Get ACFs
167 acfs <- lapply( dat_tradesigns, function(x) x %>%
168                  .[, acf(Trade.Sign, lag.max = length(Trade.Sign))])
169
170 ## Nice dataframes
171 dat_AGL_ACF <- tibble(ACF = as.numeric(acfs$AGL$acf),
172                        Lag = as.numeric(acfs$AGL$lag))
173
174 dat_NPN_ACF <- tibble(ACF = as.numeric(acfs$NPN$acf),
175                        Lag = as.numeric(acfs$NPN$lag))
176
177 ## Plot the things
178 setwd(dir_figs)
179
180 # cairo_pdf("HTF_Ass1_fig_ACF_Orderflow_AGL.pdf", height = 5, width = 5)
181 ggplot(dat_AGL_ACF, aes(x = Lag, y = ACF)) +
182   geom_line() +
183   scale_x_log10(labels = function(x) format(x, scientific = FALSE)) +
184   scale_y_log10(labels = function(x) format(x, scientific = TRUE)) +

```

```

185  labs(y = "ACF \n", x = "\n Lag", title = Stock1,
186      subtitle = "Autocorrelation of Order-Flow on log-log Scale") +
187  theme_bw() +
188  theme(aspect.ratio = 0.75)
189 # dev.off()
190
191 # cairo_pdf("HFT_Ass1_fig_ACF_Orderflow_NPN.pdf", height = 5, width = 5)
192 ggplot(dat_NPN_ACF, aes(x = Lag, y = ACF)) +
193   geom_line() +
194   scale_x_log10(labels = function(x) format(x, scientific = FALSE)) +
195   scale_y_log10(labels = function(x) format(x, scientific = TRUE)) +
196   labs(y = "ACF \n", x = "\n Lag", title = Stock2,
197        subtitle = "Autocorrelation of Order-Flow on log-log Scale") +
198   theme_bw() +
199   theme(aspect.ratio = 0.75)
200 # dev.off()
201
202 setwd(dir_script)
203
204 # 3. Create Plots >> Inter-Arrival Times ----
205
206 ## Frequency of Inter-Arrival Times
207
208 setwd(dir_figs)
209
210 ### Anglo
211 dat_AGL <- dat_TAQ$AGL
212
213 # cairo_pdf("HFT_Ass1_fig_Freq10s_InterArrival_AGL.pdf", height = 5, width = 5)
214 ggplot(dat_AGL[Inter.Arrival.times < 10], aes(x = Inter.Arrival.times)) +
215   geom_histogram(aes(fill = ..count..), colour="black", bins = 10) +
216   scale_fill_viridis() +
217   labs(title = Stock1,
218        subtitle = "Times less than 10 seconds",
219        y = "", x = "\n Inter-Arrival Times (Seconds)") +
220   theme_bw() +
221   theme(aspect.ratio = 1, legend.position = "none")
222 # dev.off()
223
224 # cairo_pdf("HFT_Ass1_fig_Freq60s_InterArrival_AGL.pdf", height = 5, width = 5)
225 ggplot(dat_AGL[Inter.Arrival.times < 60], aes(x = Inter.Arrival.times)) +
226   geom_histogram(aes(fill = ..count..), colour="black", bins = 30) +
227   scale_fill_viridis() +
228   labs(title = Stock1,
229        subtitle = "Times less than 60 seconds",
230        y = "", x = "\n Inter-Arrival Times (Seconds)") +
231   theme_bw() +
232   theme(aspect.ratio = 1, legend.position = "none")
233 # dev.off()
234
235 # cairo_pdf("HFT_Ass1_fig_Freq10m_InterArrival_AGL.pdf", height = 5, width = 5)
236 ggplot(dat_AGL[Inter.Arrival.times < 600], aes(x = Inter.Arrival.times)) +
237   geom_histogram(aes(fill = ..count..), colour="black", bins = 50) +
238   scale_fill_viridis() +
239   labs(title = Stock1,
240        subtitle = "Times less than 10 minutes",
241        y = "", x = "\n Inter-Arrival Times (Seconds)") +
242   theme_bw() +
243   theme(aspect.ratio = 1, legend.position = "none")
244 # dev.off()
245
246 # cairo_pdf("HFT_Ass1_fig_FreqAll_InterArrival_AGL.pdf", height = 5, width = 5)
247 ggplot(dat_AGL, aes(x = Inter.Arrival.times)) +
248   geom_histogram(aes(fill = ..count..), colour="black", bins = 100) +
249   scale_fill_viridis() +
250   labs(title = Stock1,
251        subtitle = "Times for all Trading Events - Absolute Scale",
252        y = "", x = "\n Inter-Arrival Times (Seconds)") +
253   theme_bw() +
254   theme(aspect.ratio = 1, legend.position = "none")
255 # dev.off()
256
257 # cairo_pdf("HFT_Ass1_fig_FreqAllLog_InterArrival_AGL.pdf", height = 5, width = 5)

```

```

258| ggplot(dat_AGL, aes(x = Inter.Arrival.times)) +
259|   geom_histogram(aes(fill = ..count..), colour="black", bins = 15) +
260|   scale_fill_viridis() +
261|   scale_x_log10() +
262|   labs(title = Stock1,
263|     subtitle = "Times for all Trading Events - Log Scale",
264|     y = "", x = "\n Inter-Arrival Times (Seconds)") +
265|   theme_bw() +
266|   theme(aspect.ratio = 1, legend.position = "none")
267| # dev.off()
268|
269| ### Naspers
270| dat_NPN <- dat_TAQ$NPN
271|
272| # cairo_pdf("HFT_Ass1_fig_Freq10s_InterArrival_NPN.pdf", height = 5, width = 5)
273| ggplot(dat_NPN[Inter.Arrival.times < 10], aes(x = Inter.Arrival.times)) +
274|   geom_histogram(aes(fill = ..count..), colour="black", bins = 10) +
275|   scale_fill_viridis() +
276|   labs(title = Stock2,
277|     subtitle = "Times less than 10 seconds",
278|     y = "", x = "\n Inter-Arrival Times (Seconds)") +
279|   theme_bw() +
280|   theme(aspect.ratio = 1, legend.position = "none")
281| # dev.off()
282|
283| # cairo_pdf("HFT_Ass1_fig_Freq60s_InterArrival_NPN.pdf", height = 5, width = 5)
284| ggplot(dat_NPN[Inter.Arrival.times < 60], aes(x = Inter.Arrival.times)) +
285|   geom_histogram(aes(fill = ..count..), colour="black", bins = 30) +
286|   scale_fill_viridis() +
287|   labs(title = Stock2,
288|     subtitle = "Times less than 60 seconds",
289|     y = "", x = "\n Inter-Arrival Times (Seconds)") +
290|   theme_bw() +
291|   theme(aspect.ratio = 1, legend.position = "none")
292| # dev.off()
293|
294| # cairo_pdf("HFT_Ass1_fig_Freq10m_InterArrival_NPN.pdf", height = 5, width = 5)
295| ggplot(dat_NPN[Inter.Arrival.times < 600], aes(x = Inter.Arrival.times)) +
296|   geom_histogram(aes(fill = ..count..), colour="black", bins = 50) +
297|   scale_fill_viridis() +
298|   labs(title = Stock2,
299|     subtitle = "Times less than 10 minutes",
300|     y = "", x = "\n Inter-Arrival Times (Seconds)") +
301|   theme_bw() +
302|   theme(aspect.ratio = 1, legend.position = "none")
303| # dev.off()
304|
305| # cairo_pdf("HFT_Ass1_fig_FreqAll_InterArrival_NPN.pdf", height = 5, width = 5)
306| ggplot(dat_NPN, aes(x = Inter.Arrival.times)) +
307|   geom_histogram(aes(fill = ..count..), colour="black", bins = 100) +
308|   scale_fill_viridis() +
309|   labs(title = Stock2,
310|     subtitle = "Times for all Trading Events - Absolute Scale",
311|     y = "", x = "\n Inter-Arrival Times (Seconds)") +
312|   theme_bw() +
313|   theme(aspect.ratio = 1, legend.position = "none")
314| # dev.off()
315|
316| # cairo_pdf("HFT_Ass1_fig_FreqAllLog_InterArrival_NPN.pdf", height = 5, width = 5)
317| ggplot(dat_NPN, aes(x = Inter.Arrival.times)) +
318|   geom_histogram(aes(fill = ..count..), colour="black", bins = 14) +
319|   scale_fill_viridis() +
320|   scale_x_log10() +
321|   labs(title = Stock1,
322|     subtitle = "Times for all Trading Events - Log Scale",
323|     y = "", x = "\n Inter-Arrival Times (Seconds)") +
324|   theme_bw() +
325|   theme(aspect.ratio = 1, legend.position = "none")
326| # dev.off()
327|
328| ## QQ Plots
329|
330| # cairo_pdf("HFT_Ass1_fig_QQExp_InterArrival_AGL.pdf", height = 5, width = 5)

```

```

331 | ggplot(dat_AGL, aes(sample = Inter.Arrival.times)) +
332 |   stat_qq_line(distribution = stats::qexp, colour = "firebrick1") +
333 |   stat_qq(distribution = stats::qexp, geom = "point") +
334 |   labs(x = "\n Exponential Quantiles", y = "Inter-Arrival Quantiles \n",
335 |         title = Stock1,
336 |         subtitle = "QQ-Plot Relative to Exponential Distribution") +
337 |   theme_bw() +
338 |   theme(aspect.ratio = 1)
339 | # dev.off()
340 |
341 | # cairo_pdf("HFT_Ass1_fig_QQExp_InterArrival_NPN.pdf", height = 5, width = 5)
342 | ggplot(dat_NPN, aes(sample = Inter.Arrival.times)) +
343 |   stat_qq_line(distribution = stats::qexp, colour = "firebrick1") +
344 |   stat_qq(distribution = stats::qexp, geom = "point") +
345 |   labs(x = "\n Exponential Quantiles", y = "Inter-Arrival Quantiles \n",
346 |         title = Stock2,
347 |         subtitle = "QQ-Plot Relative to Exponential Distribution") +
348 |   theme_bw() +
349 |   theme(aspect.ratio = 1)
350 | # dev.off()
351 |
352 ## ACF Plots of Inter-Arrival Times
353 |
354 acfs2 <- lapply( dat_tradesigns, function(x) x %>%
355                   .[, acf(na.omit(Inter.Arrival.times),
356                           lag.max = length(na.omit(Inter.Arrival.times)))] )
357 |
358 ## Nice dataframes
359 dat_AGL_ACF_times <- tibble(ACF = as.numeric(acfs2$AGL$acf),
360                               Lag = as.numeric(acfs2$AGL$lag))
361 |
362 dat_NPN_ACF_times <- tibble(ACF = as.numeric(acfs2$NPN$acf),
363                               Lag = as.numeric(acfs2$NPN$lag))
364 |
365 ## Plot the things
366 setwd(dir_figs)
367 |
368 # cairo_pdf("HFT_Ass1_fig_ACF_InterArrival_AGL.pdf", height = 5, width = 5)
369 ggplot(dat_AGL_ACF_times, aes(x = Lag, y = ACF)) +
370   geom_line() +
371   scale_x_log10(labels = function(x) format(x, scientific = FALSE)) +
372   scale_y_log10(labels = function(x) format(x, scientific = TRUE)) +
373   labs(y = "ACF \n", x = "\n Lag", title = Stock1,
374         subtitle = "Autocorrelation of Inter-Arrival Times on log-log Scale") +
375   theme_bw() +
376   theme(aspect.ratio = 0.75)
377 # dev.off()
378 |
379 # cairo_pdf("HFT_Ass1_fig_ACF_InterArrival_NPN.pdf", height = 5, width = 5)
380 ggplot(dat_NPN_ACF_times, aes(x = Lag, y = ACF)) +
381   geom_line() +
382   scale_x_log10(labels = function(x) format(x, scientific = FALSE)) +
383   scale_y_log10(labels = function(x) format(x, scientific = TRUE)) +
384   labs(y = "ACF \n", x = "\n Lag", title = Stock2,
385         subtitle = "Autocorrelation of Inter-Arrival Times on log-log Scale") +
386   theme_bw() +
387   theme(aspect.ratio = 0.75)
388 # dev.off()
389 |
390 setwd(dir_script)
391 |
392 #-----#

```



### 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".

Date: August 7, 2019