

Machine Learning & Statistical Data Mining

- Stock Trend Prediction Analysis - The JPMorgan Chase & Co (JPM) case

Table of Contents

Table of Contents.....	2
Abstract.....	4
Definition of Terms.....	4
Software Dependencies.....	4
Hardware.....	4
How to Run the Code.....	4
Introduction.....	5
Data Description	6
Data Pre-processing	7
Simple Moving Average (Sma) 20/50/100/200.....	7
Exponential Moving Average (Ema) 20/50/100/200	7
Relative Strength Index (RSI)	8
Average True Range (ATR)	8
Stochastic Momentum Index (SMI)	8
Moving Average Convergence/Divergence (MACD).....	9
Bollinger Bands	9
Money Flow Index (MFI).....	9
Parabolic Stop And Reverse (SAR).....	9
Volatility.....	10
The price lags.....	10
The Volume	10
Dummification	11
Skewness Reduction	12
Missing Data	13
Data Visualisation.....	14
The Sliding Time Windows	17
Feature Extraction.....	18
Normal Distribution Test.....	19
Models Training & Hyper Parameters Tuning	22
The Ridge	23
The Lasso.....	26
Linear Discriminant Analysis (LDA)	29
Quadratic Discriminant Analysis (QDA)	32
Decision Tree.....	35
RandomForest.....	41
Boosting.....	46
Support Vector Machine (SVM)	49
Evaluation.....	53

Challenges & Potential Improvements.....	54
Conclusion	55
References.....	56
Appendices.....	57
Appendix A – The Kolmogorov Smirnov Test Details	57
Appendix B – Code for Graphics Generation	59

Abstract

This aim of this study is to establish whether the stock trend is influenced by past prices only, or whether external factors also play a role. In this report, the analysis only focus one of stock; namely JPMorgan Chase & Co (JPM). The trend at date t is generated from the log returns between $t+1$ and t . When the log return is positive, then the stock trend is classified as *Up*, when it is negative, it is classified as *Down*. Otherwise, it is considered as *Neutral*. For this, several supervised classification models, such as LDA, QDA, SVM, etc. have been selected to assess the predictive power of the past prices, the volume as well as numerous technical indicators to forecast the trend. The result of the experiment indicates that, after feature selection and model optimizations, the best model test performance is obtained by the Quadratic Discriminant Analysis (QDA), the Ridge and the Linear Discriminant Analysis (LDA) models, with an average test accuracy of approximately 80% (and standard deviation of 1.4%) across these models. Consequently, it can be concluded that most of the price variation is mainly attributable to past prices fluctuations (endogenous factors). However, there is still a large proportion (20.00%) that does not seem to be explained by price movements. The potential impact of exogenous factors (such as fundamental or sentiment indicators) is not part of this experiment.

Definition of Terms

- Basis Point (bp): This represents a percent of one percent, i.e. 0.0001.
- Technical Analysis: This a trading tool employed to forecast securities future movement by analyzing statistics gathered from trading activity [1].
- Technical Indicator: A statistical or otherwise fabricated metric (usually using the underlying price/volume information) to help predict the underlying future price move [2]. Technical indicators used in this study are detailed in the *Data Pre-processing* section.
- API: The application program interface (API) defines how software components should interact. [3]
- XLF Index: this is an exchange traded fund based on the financial sector industry. It contains approximately 60 stocks. The XLF index price moves is an indicator of the sector performance.
- Fundamental analysis: it is a method of evaluating a security net asset value by measuring its intrinsic value, examining related economic, financial and other qualitative and quantitative factors [4].
- Sentiment analysis: this is the process of opinion mining to gauge positive, neutral or negative 'feeling' the market holds, at given point time, relating to a stock performance.

Software Dependencies

R - Version 3.3.3 & R Studio - Version 0.99.903

Hardware

Windows 10 64bits platform, supported by a Intel Core i7-6700HQ processor / RAM: DDR4 8GB.

How to Run the Code

The code can be run in its entirety (or section by section) in R studio. For this, the following files need to be stored under a unique directory:

- Stock_Trend_Following_Analysis_Assignment2.rmd
- JPM.csv
- XLF.csv

These extra files and folder are not required for the correct running of the R code

- StockPriceCropper.ipynb
- The result folder contains the Training/Optimisation/Testing accuracy rate for each model, as well the QDA confusion matrix which lists the sensitivity, precision, F1, etc. measures of the tested model.

The file configuration is shown in *Code Snippet 0*. This is for information only; the code should not be altered.

```
54 - #####
55 - #####
56 - ##### Setup / User Configuration
57 - #####
58 - #####
59 - ```{r set_up}
60 index_name = "XLF"
61 constituents_full_path = paste(getwd(), "/", index_name, ".csv", sep="")
62 stock_mk_price_dir = paste(getwd(), "/", sep="")
63 out_dir = paste(getwd(), "/", sep="")
64
```

Code Snippet 0 – Configuration

Introduction

The aim of this study is to evaluate the prediction power of lagged prices, the volume and numerous technical indicators for forecasting stock price movements (*Up/Down* and *Neutral*). This experiment only focuses on one stock, namely JPM, but the methodology can be extended to other stocks, indices, and asset classes (such as FX, Commodities, etc). To achieve this, the stock price data is downloaded from the Yahoo Finance web site [5] using a hand-crafted python web scraping program (c.f. section *Data Description*). The daily log stock returns are transformed into *Up/Down* and *Neutral* categories, indicating respectively whether the stock increased/decreased or remained unchanged between a date t and $t+1$. These categories are then stored under a new explained variable called *Direction*. Numerous technical indicators are generated (c.f. section *Data Pre-processing*) and used as attributes, amongst other explanatory variables (such as lagged prices and the volume), to attempt to predict the trend *Direction*.

Missing data produce a prediction bias. Therefore, records containing missing data are removed (c.f. section *Missing Data*). Following this, the analysis starts with a feature selection. The aim is to reduce the model complexity and focus on attributes concentrating the most prediction power (i.e. the most variance). The next step focuses on training and optimising several supervised classification models (such as LDA, QDA, SVM, etc.,) against the JPM times series *Direction* response variable. The training and optimisation phase are performed on a sliding time window, as the usual cross validation methodology should be avoided in this instance. Indeed, time series contain patterns that need to be chronologically preserved (c.f. section *The Sliding Time Windows*). The last phase involves choosing the best fitting model on training data and optimising the model hyper-parameters. Once the best training model is selected, the model is ready for testing against an unseen test data set.

Data Description

The JPM stock historical price data (a.k.a. the time series) is at the heart of the analysis. It was downloaded from the *Yahoo IChart API*, via a custom-made *Python* program. The following piece of code describes the common functions used for downloading data from the Yahoo API (c.f. 'In [6]' the *pull_historical_data()* function). The main calling function lives in section 'In [8]' below. The code is fully commented in the *Code Snippet 1* section below. A fully running version of the code is available in *StockPriceCropper.ipynb*. The result of running this code is the production of the *JPM.csv* file that contains the historical price/volume information.

The time series for the JPM stock starts on 30-Dec-1983 and ends on 14-Dec-2016. It contains the *Open/High/Low/Close/Volume* and *Adj Close* for each date. They represent successively: i) the daily opening price, ii) the highest/lowest daily price, iii) the daily price close level, iv) the daily closing price amendment (i.e. adjustment), i.e. the close stock price including any distributions and corporate actions that occurred at any time prior to the next day's opening, as well as iv) the volume. The volume represents the number of shares that changed hands during the trading hours.

Stock Prices Cropper Functions

```
In [6]: print ("Start Stock Prices Cropper Functions")

base_url = "http://ichart.finance.yahoo.com/table.csv?s="
input_path = output_path = os.getcwd()

#Create the directory, if it does not exist
def mkdir(directory_full_path):
    if not os.path.exists(directory_full_path):
        os.makedirs(directory_full_path)

#Generate the full URL based on the base_url and the ticker name (e.g. JPM)
def make_url(ticker_symbol):
    return base_url + ticker_symbol

#Returns the directory output name as well as the full file path
def make_output_filename(ticker_symbol, directory="Unknown"):
    return output_path + "/" + directory + "/", output_path + "/" + directory + "/" + ticker_symbol + ".csv"

def pull_historical_data(ticker_symbol, directory="Unknown"):
    try:
        #Generate the output directory and full path for the stock file name (e.g. ./JPM.csv)
        directory_full_path, file_full_path = make_output_filename(ticker_symbol, directory)
        #Generate the full URL
        ticker_url = make_url(ticker_symbol)
        #Make the directory if it does not exist
        mkdir(directory_full_path)
        #Get the data from the url and store it in the defined file path
        urlretrieve(ticker_url, file_full_path)
    except Exception as e: # catch *all* exceptions
        directory_full_path, file_full_path = make_output_filename(ticker_symbol, directory)
        outfile = open(file_full_path, "w")
        print(e)
        outfile.write(str(e))
        outfile.close()

print ("End Stock Prices Cropper Functions")
```

```
Start Stock Prices Cropper Functions
End Stock Prices Cropper Functions
```

Run Stock Prices Cropper...

The data is sourced for free from Yahoo, via the Yahoo API.

```
In [8]: print ("Start Cropping...")
#Get the JPM historical data from beg of time...
pull_historical_data ('JPM', "")
print ("End Cropping...")
```

```
Start Cropping...
End Cropping...
```

Code Snippet 1 – JPM Time Series Download

Data Pre-processing

The process of analysis trends in times series usually involves the fabrication of numerous indicators directly derived from the data under analysis. In this case, the historical stock Close price. Trend filtering methods uses econometric/statistical estimators to extract trends from time series [6]. In this study, we generate and study the impact of a selection of the indicators presented in [7]. Unless stated otherwise, the Close price is used to generate technical indicators. Some metrics, such as the *Average True Range (ATR)* may use the HLC (high/Low/Close) price. The HLC price aggregate is built as follows: $HLC = (High + Low + Close) / 3$. The following sections details all the technical indicators and other derived data formulas that are used as part of the model feature selection. All technical indicators have been implemented by calling the relevant R functions from the *quantmod* library [7bis] (c.f. Code Snippet 2). In the following sections [-1] means today minus one business day. This is also known as $t-1$.

```
192 #Add a the indicators column (the average value of the High,Low and Close)
193 sma_df["Hlc"] = (sma_df["High"] + sma_df["Low"] + sma_df["close"])/3
194 sma_df["Hl"] = (sma_df["High"] + sma_df["Low"]) /2
195 sma_df["Sma20"] = SMA(sma_df$close, n=20)
196 sma_df["Sma50"] = SMA(sma_df$close, n=50)
197 sma_df["Sma100"] = SMA(sma_df$close, n=100)
198 sma_df["Sma200"] = SMA(sma_df$close, n=200)
199 sma_df["Ema20"] = EMA(sma_df$close, n=20)
200 sma_df["Ema50"] = EMA(sma_df$close, n=50)
201 sma_df["Ema100"] = EMA(sma_df$close, n=100)
202 sma_df["Ema200"] = EMA(sma_df$close, n=200)
203 sma_df["Rsi"] = RSI(sma_df$close, n=14)
204 sma_df["Atr_tr"] = ATR(HLC(sma_df[,c("High", "Low", "close")]))[,1]
205 sma_df["Atr_atr"] = ATR(HLC(sma_df[,c("High", "Low", "close")]))[,2]
206 sma_df["Atr_trueHigh"] = ATR(HLC(sma_df[,c("High", "Low", "close")]))[,3]
207 sma_df["Atr_trueLow"] = ATR(HLC(sma_df[,c("High", "Low", "close")]))[,4]
208 sma_df["Smi_smi"] = SMI(HLC(sma_df[,c("High", "Low", "close")]))[,1]
209 sma_df["Smi_signal"] = SMI(HLC(sma_df[,c("High", "Low", "close")]))[,2]
210 sma_df["Macd_macd"] = MACD(sma_df[,c("close")]))[,1]
211 sma_df["Macd_signal"] = MACD(sma_df[,c("close")]))[,2]
212 sma_df["Bb_dn"] = BBands(HLC(sma_df[,c("High", "Low", "close")]))[,1]
213 sma_df["Bb_mavg"] = BBands(HLC(sma_df[,c("High", "Low", "close")]))[,2]
214 sma_df["Bb_up"] = BBands(HLC(sma_df[,c("High", "Low", "close")]))[,3]
215 sma_df["Bb_pctB"] = BBands(HLC(sma_df[,c("High", "Low", "close")]))[,4]
216 sma_df["Mfi"] = MFI(HLC(sma_df[,c("High", "Low", "close")]), sma_df[,c("Volume")])
217 sma_df["Sar"] = SAR(sma_df[,c("High", "close")])
218 sma_df["Volatility"] = volatility(OHLC(sma_df[,c("open", "High", "Low", "close")]), calc = "garman")
219
```

Code Snippet 2 – Technical Indicators implementations

Simple Moving Average (Sma) 20/50/100/200

Each output value is the average of the previous n day values (e.g. 20/50/100/200 days). As each value in the period carries equal weight, it makes it less responsive to recent changes [8].

$$SMA = \frac{\sum_{i=1}^n price}{n}$$

With n the number of lag day for the moving average (e.g. 20/50/100/200)

Exponential Moving Average (Ema) 20/50/100/200

Exponential Moving Average is a cumulative calculation, including all data (even values outside of the period). More recent values have a greater contribution to the average, past values have a diminishing contribution [8][9].

Multiplier: $(2 / (\text{time periods} + 1))$

EMA: $\{\text{Close} - \text{EMA}(\text{previous day})\} \times \text{multiplier} + \text{EMA}(\text{previous day})$.

With *time periods* being the number of lag day for the moving average (e.g. 20/50/100/200)

Relative Strength Index (RSI)

It produces a ratio of the recent upward price movements to the absolute price movement. The index ranges from 0 to 100. A RSI strictly greater 70 is interpreted as an overbought indicator. When the RSI is strictly below 30, it is interpreted as an oversold indicator [8].

```
If  $close > close_{[-1]}$  then
   $up = close - close_{[-1]}$ 
   $dn = 0$ 
else
   $up = 0$ 
   $dn = close_{[-1]} - close$ 
 $upavg = \frac{upavg \times (n-1) + up}{n}$ 
 $dnavg = \frac{dnavg \times (n-1) + dn}{n}$ 
 $RMI = 100 \times \frac{upavg}{upavg + dnavg}$ 
```

Average True Range (ATR)

The ATR relates to the True Range Welles Wilder moving average. The ATR is a measure of volatility. High ATR values indicate high volatility, and low values indicate low volatility. The True range is used to determine the usual trading range of a stock [8].

```
TrueHigh = Highest of  $high[0]$  or  $close[-1]$ 
TrueLow = Lowest of  $low[0]$  or  $close[-1]$ 
 $TR = TrueHigh - TrueLow$ 
```

The formula is sometimes stated as:

$TR =$ The greatest of the following:

```
|  $high[0] - low[0]$  |
|  $high[0] - close[-1]$  |
|  $low[0] - close[-1]$  |
```

```
TrueHigh = Highest of  $high[0]$  or  $close[-1]$ 
TrueLow = Lowest of  $low[0]$  or  $close[-1]$ 
 $TR = TrueHigh - TrueLow$ 
 $ATR = \frac{TR_{[-1]} \times (n-1) + TR}{n}$ 
```

Stochastic Momentum Index (SMI)

The SMI is built on the Stochastic Oscillator. The Stochastic Oscillator calculates the position of the close price in relation to the high/low range. The SMI calculates the position of the Close price in relation to high/low range midpoint. SMI values range between +100 and -100. The SMI is greater than zero, when a close is greater than the midpoint (and reverse).

Extreme high/low SMI values indicate overbought/oversold conditions. A buy signal is issued when the SMI rises above -50, or when it crosses above the signal line (and reverse).

```
 $cm = close - \frac{highesthigh - lowestlow}{2}$ 
 $hl = highesthigh - lowestlow$ 
 $cm = EMA(EMA(cm))$ 
 $hl = EMA(EMA(hl))$ 
 $SMI = 100 \times \left( \frac{cm}{hl/2} \right)$ 
 $Signal = EMA(SMI)$ 
```


Moving Average Convergence/Divergence (MACD)

The MACD is the difference between two Exponential Moving Averages; a short and long one. The Signal line is an Exponential Moving Average of the MACD [8].

High/low values indicate an overbought/oversold asset. When the MACD line crosses above/below the signal line a buy/sell signal is generated. The signal is confirmed when the MACD is above/below zero buy/sell.

$$\begin{aligned} \text{shortema} &= 0.15 \times \text{price} + 0.85 * \text{shortema}_{[-1]} \\ \text{longema} &= 0.075 \times \text{price} + 0.925 * \text{longema}_{[-1]} \\ \text{MACD} &= \text{shortema} - \text{longema} \end{aligned}$$

Bollinger Bands

Bollinger Bands consist of three lines: the simple moving average and the upper/lower bands. The latter are usually two standard deviations above/below [8].

Bollinger Bands are not, in themselves, buy or sell signals. They indicate overbought or oversold conditions. A price near to the upper or lower band indicates a potential imminent reversal. The moving average becomes a support or resistance level.

$$\begin{aligned} TP &= \frac{\text{high} + \text{low} + \text{close}}{3} \\ \text{MidBand} &= \text{SimpleMovingAverage}(TP) \\ \text{UpperBand} &= \text{MidBand} + F \times \sigma(TP) \\ \text{LowerBand} &= \text{MidBand} - F \times \sigma(TP) \end{aligned}$$

With F usually set to 2.

Money Flow Index (MFI)

The MFI calculates the ratio of money flowing into and out of a security [8]. Money Flow Index ranges between 0 and 100. Values above 80/below 20 indicate market tops/bottoms.

$$\begin{aligned} \text{TypicalPrice} &= \frac{\text{high} + \text{low} + \text{close}}{3} \\ \text{MoneyFlow} &= \text{TypicalPrice} \times \text{volume} \\ \text{If } \text{TypicalPrice} > \text{TypicalPrice}_{[-1]} & \\ \quad \text{PositiveMoneyFlow} &= \text{PositiveMoneyFlow}_{[-1]} + \text{MoneyFlow} \\ \text{else} & \\ \quad \text{NegativeMoneyFlow} &= \text{NegativeMoneyFlow}_{[-1]} + \text{MoneyFlow} \\ \text{MoneyRatio}_i &= \frac{\sum_{i=1}^n \text{PositiveMoneyFlow}}{\sum_{i=1}^n \text{NegativeMoneyFlow}} \\ \text{MoneyFlowIndex} &= 100 - \left(\frac{100}{1 + \text{MoneyRatio}} \right) \end{aligned}$$

Parabolic Stop And Reverse (SAR)

SAR calculates a trailing stop [8]. SAR advises to exit when the price crosses the SAR.

$$\begin{aligned} \text{If long and } \text{high} > \text{xp} \text{ then} & \\ \quad \text{xp} &= \text{high} \\ \quad \text{af} &= \text{af} + \text{step} \\ \text{If short and } \text{low} < \text{xp} \text{ then} & \\ \quad \text{xp} &= \text{low} \\ \quad \text{af} &= \text{af} + \text{step} \\ \text{SAR} &= (\text{xp} - \text{SAR}_{[-1]}) \times \text{af} + \text{SAR}_{[-1]} \\ \text{xp is extreme point} & \\ \text{af is acceleration factor} & \end{aligned}$$

Volatility

This indicator generates the Close price degree of variation over a business 260 days' period.

$$\sigma_{cl} = \sqrt{\frac{Z}{n-2} \sum_{i=1}^{n-1} (r_i - \bar{r})^2}$$

where $r_i = \log\left(\frac{C_i}{C_{i-1}}\right)$

and $\bar{r} = \frac{r_1 + r_2 + \dots + r_{n-1}}{n-1}$

The price lags

The $t-1$, $t-2$, $t-3$, $t-4$ and $t-5$ are the price lags, i.e. the $t-n$ prices used at date t . These price lags are used as attributes to establish whether they have an explanatory power on the trend at date t .

```
223 #Add day lags
224 sma_df = add_close_price_day_lag(sma_df, 1)
225 sma_df = add_close_price_day_lag(sma_df, 2)
226 sma_df = add_close_price_day_lag(sma_df, 3)
227 sma_df = add_close_price_day_lag(sma_df, 4)
228 sma_df = add_close_price_day_lag(sma_df, 5)
229

168 #Add the close price day lag to the analysis
169 add_close_price_day_lag = function (stock_df, nb_days){
170   #Ensure the dataset is ordered in data descending order,
171   #as the calculations assume the ordering for the calculation
172
173   #Generate the column dynamically based on the nb_days
174   close_price_lag = paste("close_price_", nb_days, "day_lag", sep="")
175   for (k in 1:nrow(stock_df)) {
176     stock_df[k,close_price_lag] = stock_df[k+nb_days,"close"]
177   }
178   return (stock_df)
179 }
180
```

Code Snippet 3 – The Price lag main function

The Volume

The volume corresponds to the amount of stock traded during each day. It is provided as part of the original data. No data transformation is performed on this quantity.

Dummification

As the study focuses on predicting stock market trends, the price is not used directly as the regressor. The main reason is that the price is not log normally distributed. However, the log return is. Prices have therefore been transformed into log returns for each date t , as follows:

$$\text{Return}_t = \log(\text{Price}_{t+1}/\text{Price}_t).$$

A new attribute named *Direction* is generated from Return_t . It is defined as per the formula below. The ε (*Epsilon*) is a user defined constant, currently set to 0.0025 (i.e. 25 basis points). This was the level chosen for the JPM stock to ensure a balanced *Direction* class. More details about class balancing are provided in the section 'Skew Reduction'.

A date t :

$\left\{ \begin{array}{l} \text{Return @}t > \varepsilon \text{ then } \textit{Direction} \text{ is set to Up} \\ \text{Return @}t < -\varepsilon \text{ then } \textit{Direction} \text{ is set to Down} \\ \text{Abs}(\text{Return @}t) \leq \varepsilon \text{ then } \textit{Direction} \text{ is set to Neutral} \end{array} \right.$

The *Direction* attribute is the predicted variable that is used in all models. The code implementation can be found in *Code Snippet 4* below.

```
129 #This function the log_returns. It indicates whether the market when up, down or sideways (neutral) from one day to the next.
130 #Please ensure the dataset is ordered in data descending order, as the calculations assume the ordering for the calculation
131 generate_log_returns_and_direction = function(stock_df, epsilon){
132   stock_df["Direction"] = NA
133   stock_df["Log_returns"] = NA
134
135   for (k in 1:nrow(stock_df)) {
136     #Generate the log return value
137     stock_df[k,"Log_returns"] = log10(stock_df$close[k+1] / stock_df$close[k])
138     #Indicate whether it is going up, down or neutral
139     if (k < nrow(stock_df)) {
140       if (!is.na( stock_df[k,"Log_returns"])){
141         #This is the default values of the 'Direction' and 'Direction_Flag'
142         stock_df[k,"Direction"] = 'NA'
143         stock_df[k,"Direction_Flag"] = -1
144         #Do the Epsilon check
145         if (abs(stock_df[k,"Log_returns"]) <= epsilon){
146           stock_df[k,"Direction"] = 'Neutral'
147           stock_df[k,"Direction_Flag"] = 1
148         }
149         else {
150           #when the return is greater today compared to the day before,
151           #then set the Direction to Up and Direction_Flag to 2
152           if (stock_df[k,"Log_returns"] > epsilon){
153             stock_df[k,"Direction"] = 'Up'
154             stock_df[k,"Direction_Flag"] = 2
155           }
156           #Else set the Direction to Up and Direction_Flag to 2
157           else {
158             stock_df[k,"Direction"] = 'Down'
159             stock_df[k,"Direction_Flag"] = 0
160           }
161         }
162       }
163     }
164   }
165   return(stock_df)
```

Code Snippet 4 – The predicted variable Dummification

Skewness Reduction

As mentioned in [10], skew or imbalance in data affect most performance metrics (e.g. Kappa, F-1 scores, etc.). The below tables show the share of each trend category in the *Direction* class. The first table shows that once rows containing missing data are removed, and ϵ is set to 0, the *Neutral* category is greatly underrepresented in the *Direction* class (only 4.17%). To rebalance the class distribution, several values for ϵ were tested. The one offering the best-balanced ratio between the classes was retained, where ϵ equals to 0.0025.

Row Labels	Count of Epsilon 0.0000	% Class Allocation
DOWN	3853	47.53%
NEUTRAL	338	4.17%
UP	3916	48.30%
Grand Total	8107	100.00%

The following tables shows the *Direction* class balance for different level of ϵ .

Row Labels	Count of Epsilon 0.0015	% Class Allocation
DOWN	3214	39.65%
NEUTRAL	1579	19.48%
UP	3314	40.87%
Grand Total	8107	100.00%

Row Labels	Count of Epsilon 0.0020	% Class Allocation
DOWN	2979	36.75%
NEUTRAL	2065	25.46%
UP	3063	37.79%
Grand Total	8107	100.00%

Row Labels	Count of Epsilon 0.0025	% Class Allocation
DOWN	2763	34.08%
NEUTRAL	2500	30.84%
UP	2844	35.08%
Grand Total	8107	100.00%

Row Labels	Count of Epsilon 0.0030	% Class Allocation
DOWN	2538	31.31%
NEUTRAL	2937	36.22%
UP	2632	32.47%
Grand Total	8107	100.00%

Missing Data

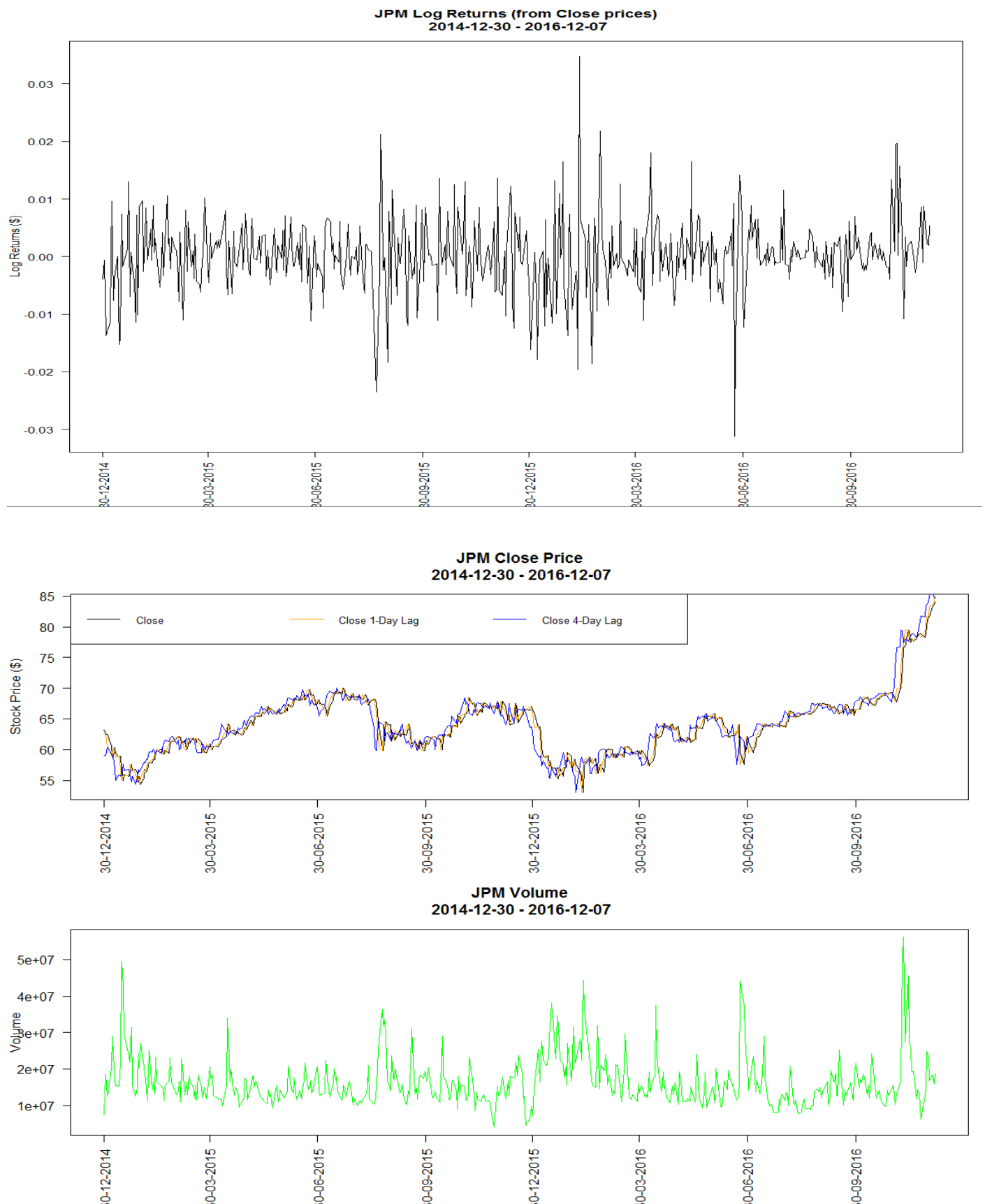
The initial raw market data, downloaded from the 'Yahoo API', do not contain any missing data. However, the data pre-processing step generates missing data, for many dates. For example, the Moving Average indicators need numerous historical dates to produce a value for a given date. Consequently, lagging indicators cannot produce data for the initial number of days corresponding to their day range lag. Data imputation, using median or K-Nearest Neighbours methods were considered. However, this would have created 'synthetic' patterns in the time series data, that do not correspond to reality. Indeed, the literature [11] suggests that financial returns do not follow a Brownian motion, i.e. a random motion. On the contrary, in case of financial crisis, volatility clustering can be observed. Volatility can also display some multifractal behaviour, which intrinsically implies regularities and volatilities patterns. Consequently, the preferred solution was simply to delete the missing data rows. The *Code Snippet 5* shows the code relating to the deletion of missing data. Approximately 12% of the rows contained in the original dataset were removed.

```
683 #####
684 #####
685 ##### Missing Data #####
686 #####
687 #####
688
689 #Remove all rows containing NA values
690 #count the number of na in the entire dataset
691 count_na = sum(is.na(sma_df))
692 #remove all the rows containing na values
693 sma_df = na.omit(sma_df)
694 #count the perc of removed row and perform some sanity check
695 removed_missing_values_perc = 100 * count_na / nrow(sma_df)
696 cat("Removed missing values represent ", removed_missing_values_perc, " % of the dataset\n", sep ="" )
697 count_na = sum(is.na(sma_df))
698 if (count_na > 0) {
699     cat("**** WARNING... **** \n")
700     cat("The NA count in the dataframe should be zero. The actual result is: ", count_na, "\n", sep ="" )
701 }
702
703 #####
```

Code Snippet 5 – Missing Data Removal

Data Visualisation

The purpose of this section is to provide several graphics to visualise the Log Returns pattern (i.e. the implicit direction) over a two-year period, alongside several of the main attributes that are part of the feature selection process. The R code for the graphics generation can be found in Appendix B.



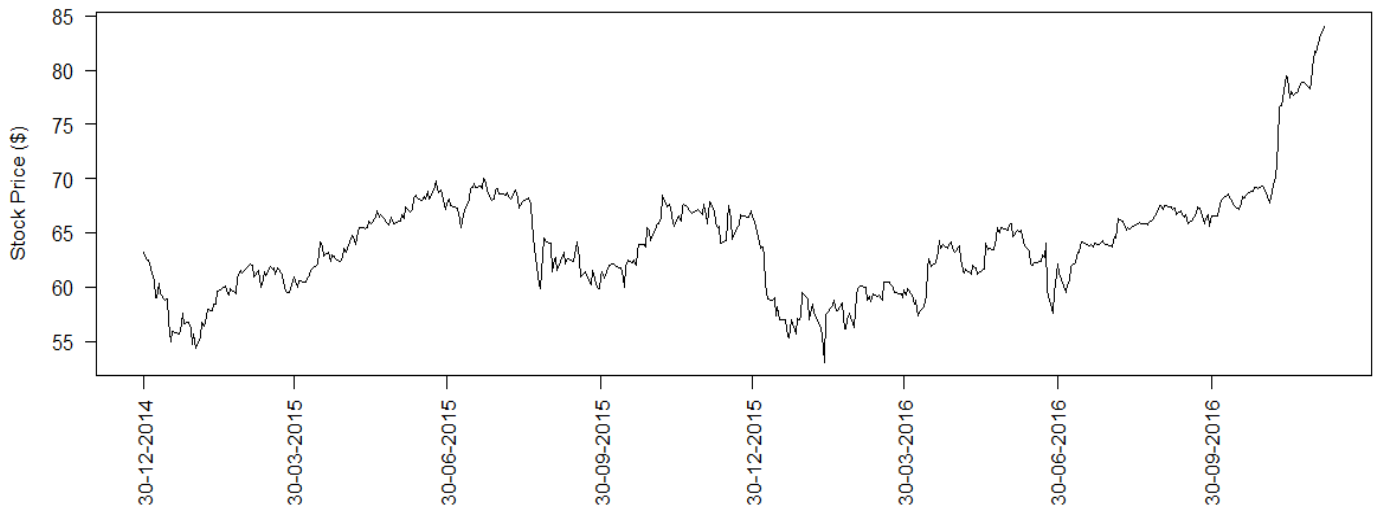
JPM Close Price Vs Simple Moving Average (SMA) Crossovers
2014-12-30 - 2016-12-07



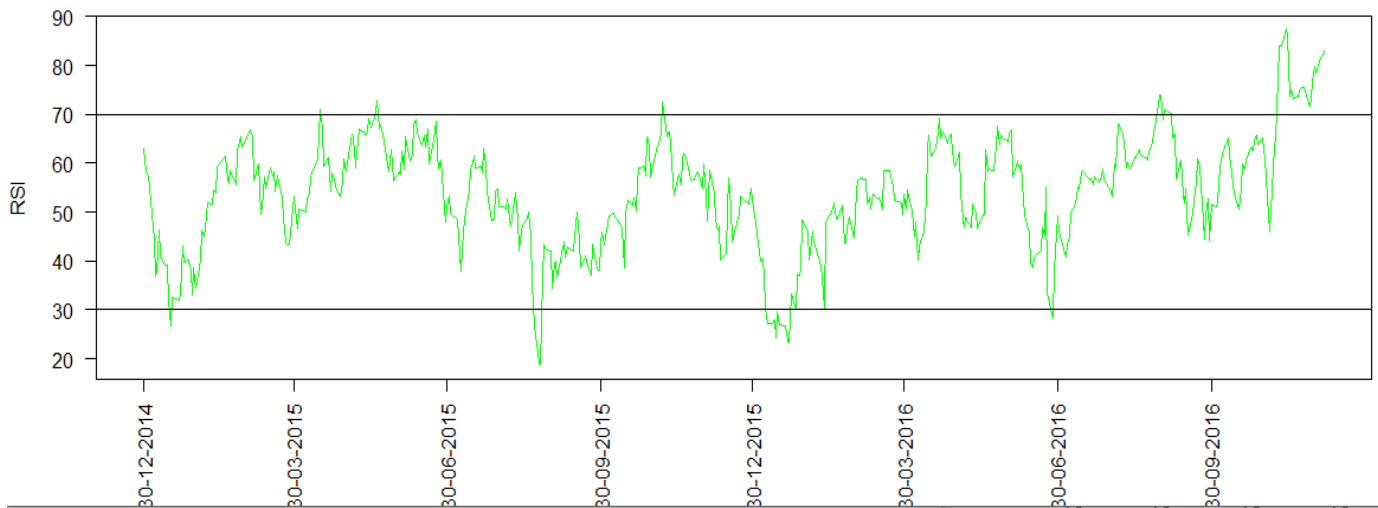
JPM Close Price Vs Exponential Moving Average (EMA) Crossovers
2014-12-30 - 2016-12-07



JPM Close Price
2014-12-30 - 2016-12-07



JPM RSI (14days)
2014-12-30 - 2016-12-07



The RSI horizontal top line represents the 70% boundary. The bottom line represents the 30% boundary.

The Sliding Time Windows

Financial time series contains historical patterns, particularly during turbulent periods such as wars, economic crisis, presidential election, etc. The aim of machine learning is to recognise these patterns to form a prediction. A frequent approach to reduce model overfitting is the cross-validation methodology (either the K-fold or Leave-One-Out). However, the re-sampling/shuffling nature of these methodologies do intrinsically modify historical patterns. Therefore, it has been ruled out for this experiment. A time window sliding approach is used instead, over a training/validation and test sets. The detail of the methodology is as follows:

- Step1
 - Create a training data set over a given time period (here a continuous 260 business days).
 - Create a validation data set, starting the business day after the end of the training data set, which represents 25% of the dates contained in the training data set (here 65 business days)
 - Create a test set that starts one day after the end of the validation data set, and last for 5 consecutive business days.
- Step2
 - Feature select on the last training set (i.e. the one closest to the end date of the dataset). A discussion of this choice is provided in the section '*Challenges and Potential Improvements*'.
- Step3
 - For each designated model (e.g. LDA, QDA, SVM, etc.):
 - Perform the model training, for a list of different attributes and modified polynomial attributes, on a 100 training/validation data sets sliding windows. In this experiment, there are 100 sliding windows of 5 days in length. In other words, when the model fitting has been performed on one sliding window, the next model fitting starts 5 business days after the previous training window start date. The process is performed 100 times. The aim is to find the set and shape of attributes that provide the best average prediction accuracy on the validation data, over 100 iterations.
 - Select the best attribute list and optimise the model on the last training/validation time window. The aim is to discover the parameters level that optimise the model on the best attribute list.
 - Train the model, its best attribute list and optimised parameters on the entire training and validation data set time frame (i.e. training start date to validation end date). Then test the model against the test data to generate the test accuracy. Run this over the 100 times windows and the average of the test accuracy is produced.

A discussion relating to the 'out of step optimisation' and choice the 100 iterations is provided in section '*Challenges and Potential Improvements*'.

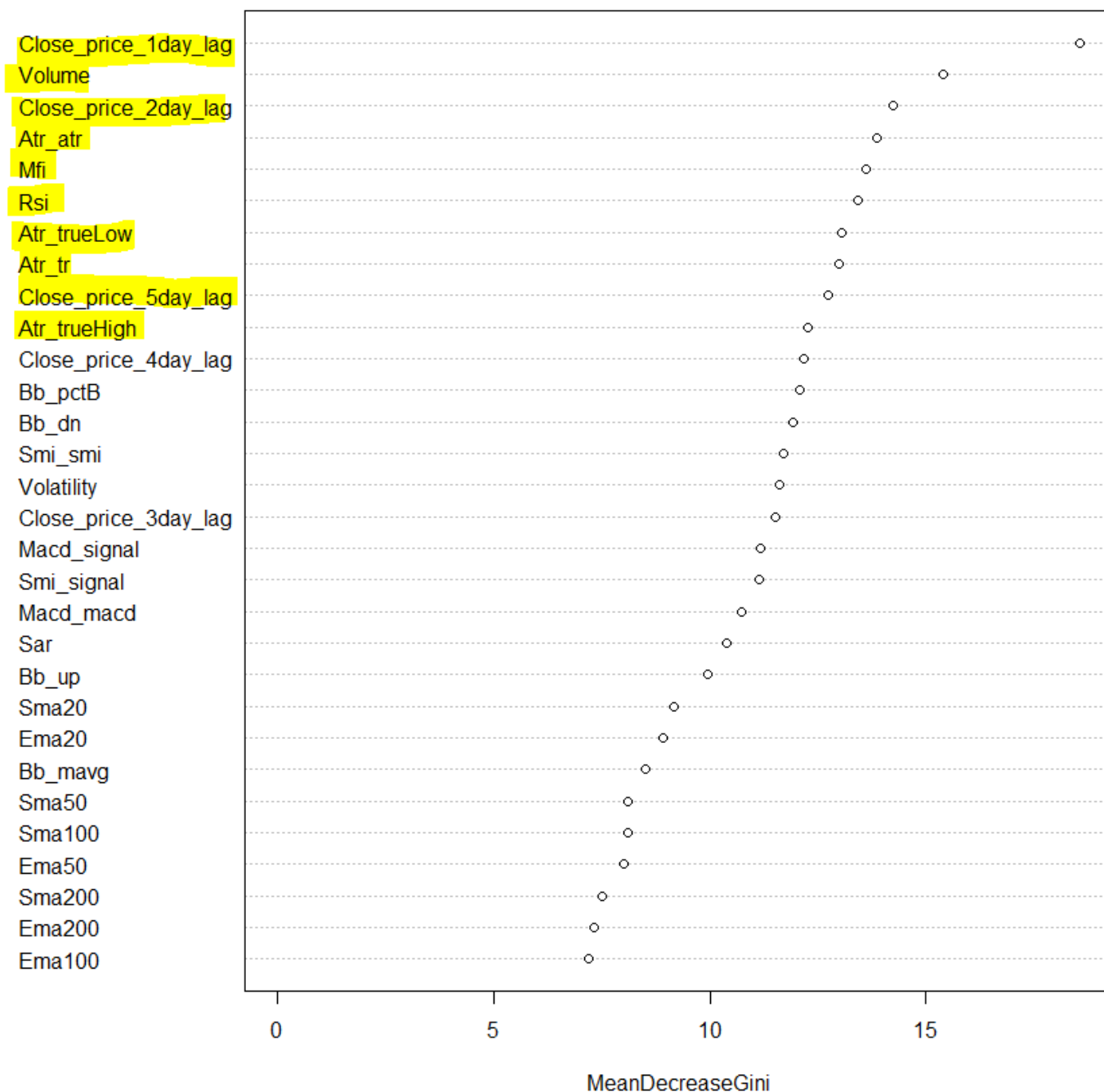
Initial experiments showed that, when tested with smaller time windows, such as 80 or 160 days; the model fitting in R breaks with the following error messages: 'Need at least two classes to do classification' or 'Some group is too small for *[the selected model]*'. Therefore, an arbitrary training time period of 260 days was selected (after experimentation). It has several advantages:

- It represents a time period slightly greater than one year. (Usually 252 days are used for the calculation of the annualized volatility [12]). Therefore, it enables to capture a larger amount of time series patterns.
- Running the LDA and QDA models on both an 80 days' time windows and a 260-day showed an increase of test performance from approximately 50% to 80% in each case.

Feature Extraction

A Random Forest Classification model is used to perform the feature selection. This model is selected as it is a popular machine learning method which couple a relatively good accuracy, robustness and ease of interpretation. The feature selection is performed only the training data of the last time window (the one closer to the dataset end date). A discussion relating to the limitation of this approach is available in section '*Challenges and Potential Improvements*'. The attributes are ordered in descending order of the '*MeanDecreaseGini*' index. The general interpretation is that a higher decreasing Gini indicates that some predictors play a greater role in classifying the data, than others. Currently, the top 10 most impacting attributes have been selected and stored in 'feature_reduction_summary.csv' file. There are also highlighted in yellow in the below picture.

Plot of the MeanDecreaseGini for each attributes for ntree 500



Normal Distribution Test

Some of the models listed below (e.g. the Linear Discriminant Analysis, the Quadratic Discriminant Analysis models, etc.) require that the explanatory variables follow a normal distribution. Although the lack of compliance should not prevent the use of these models for prediction purpose, it is interesting to establish whether these models are intrinsically statistically weak (or not). If this is the case, this may have an impact of the model stability and performance accuracy. For each of the 10 feature selected attributes, the following steps are involved:

- Create one dataset per *Direction* class type. Using the 'Close_price_1day_lag' feature as an example. It is broken in three distinct data set; one containing only rows where the *Direction* attribute is set to *Up*, a second where the rows contains the *Direction* attribute is set to *Neutral*, and the third one contains the rows where the *Direction* attribute is set to *Down*. In total, 30 datasets (attributes = 10* class types= 3) have been created.
- For each dataset, the attribute is normalized using the *scale(..) function*. It centers and scales the elements in the dataset. In other words, each element x in an attribute is normalized following this formula:

$(x - \text{avg}(x)) / \text{std}(x)$, where $\text{avg}(x)$ is the mean of x and $\text{std}(x)$ is the standard deviation of x .

- Then a Kolmogorov Smirnov Test is run on each of the dataset to establish whether H_0 can be rejected (the null hypothesis) that states; the data distribution follows a normal distribution. H_0 is rejected when the Kolmogorov Smirnov Test returns a *p-value* inferior to 0.05. The below code shows the scaling (Line 907), and the Kolmogorov Smirnov Test in action (Line 909-919).

```
904 #Produce the results of the KS test
905 kolmogorov_smirnov_normal_distribution_test = function(data,stock_name,msg) {
906   #This is the z-score scaling: (x-avg)/std as default
907   scaled_data = scale(data)
908   #Ensure all data is unique, else it breaks the ks test
909   res = ks.test(unique(scaled_data), pnorm)
910   cat(paste(stock_name, msg, sep = " "))
911   cat("\n")
912   cat("H0 = the data is normally distributed.\n", sep="")
913   if(res$p.value > 0.05){
914     cat("The ks p_value: ", res$p.value, " > 0.05 -> H0 (the null hypothesis) is NOT rejected. There is not enough evidence to reject the hypothesis that the
distribution is normal. Therefore, the data seems to follow normal distribution\n", sep="")
915   } else
916   {
917     cat("The ks p_value: ", res$p.value, " < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal
distribution.\n", sep="")
918   }
919 }
920
921 normality_check_df = final_df
```

Code Snippet 6 – Kolmogorov Smirnov Test

- QQPlots are plotted for each attribute to visually check the shape of the attribute distribution against a normally distributed QQ plot. The more the data shape match the normal QQ plot (the yellow line), the more likely is the data to follow a normal distribution. The code relating the QQPlot is shown below:

```
896 #Produce a qqplot
897 plot_qqplot = function (data, stock_name, ylab_param) {
898   data = data * 100.0
899   title = paste(stock_name, ylab_param, sep = " ")
900   qqnorm(data, main=paste(title, "\n (centered & scaled) - qqplot", sep=""), ylab=paste(ylab_param, "%", paste=""))
901   qqline(data, col="gold", lwd=2)
902 }
903
```

Code Snippet 7 –The QQplot

The below code snippet shows an end to end example with 'Close_price_1_day_lag' attribute being broken down in three datasets. For each of them the Kolmogorov Smirnov Test and QQPlot is generated. Please refer to the R code to see all the other 27 cases.

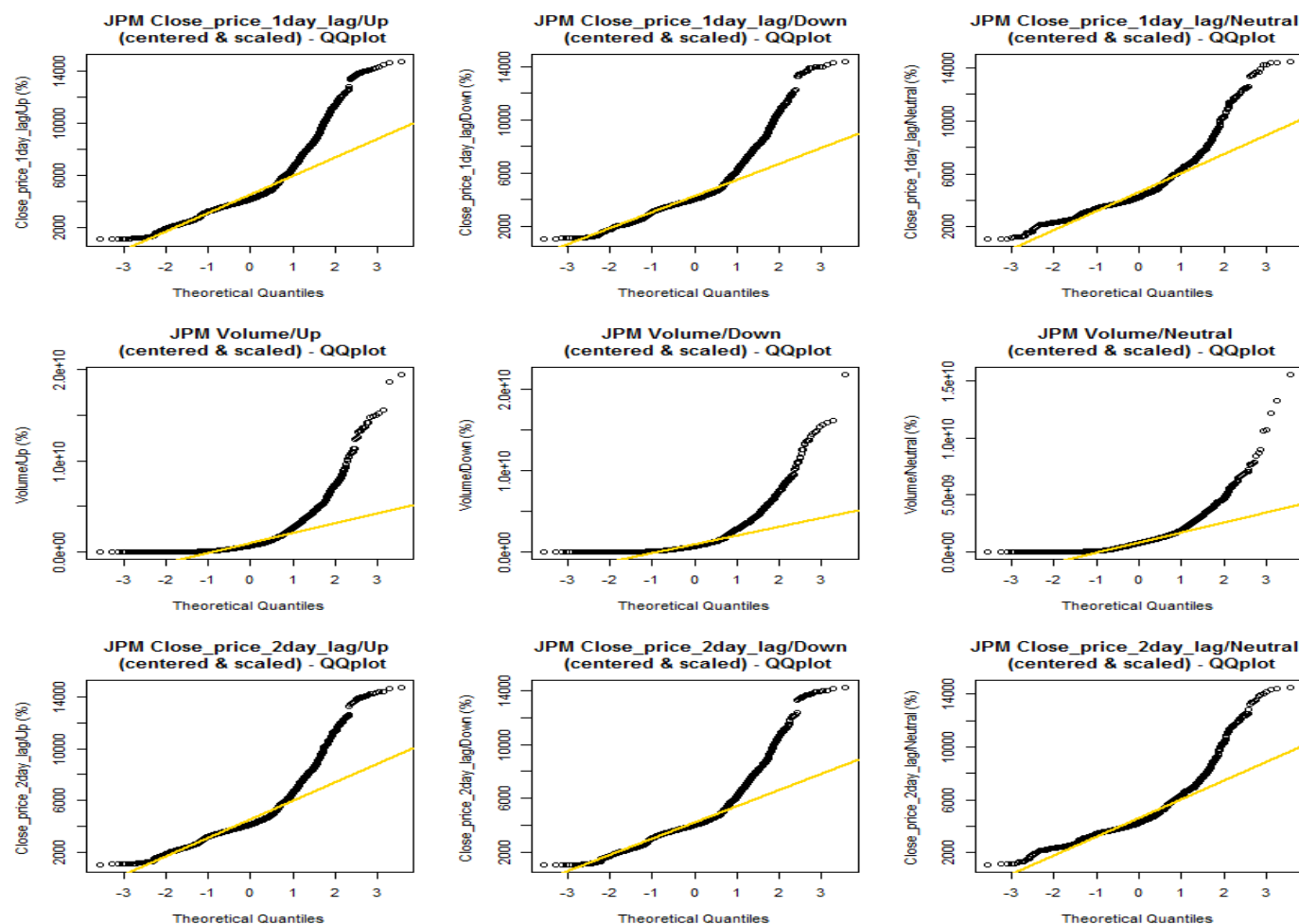
```

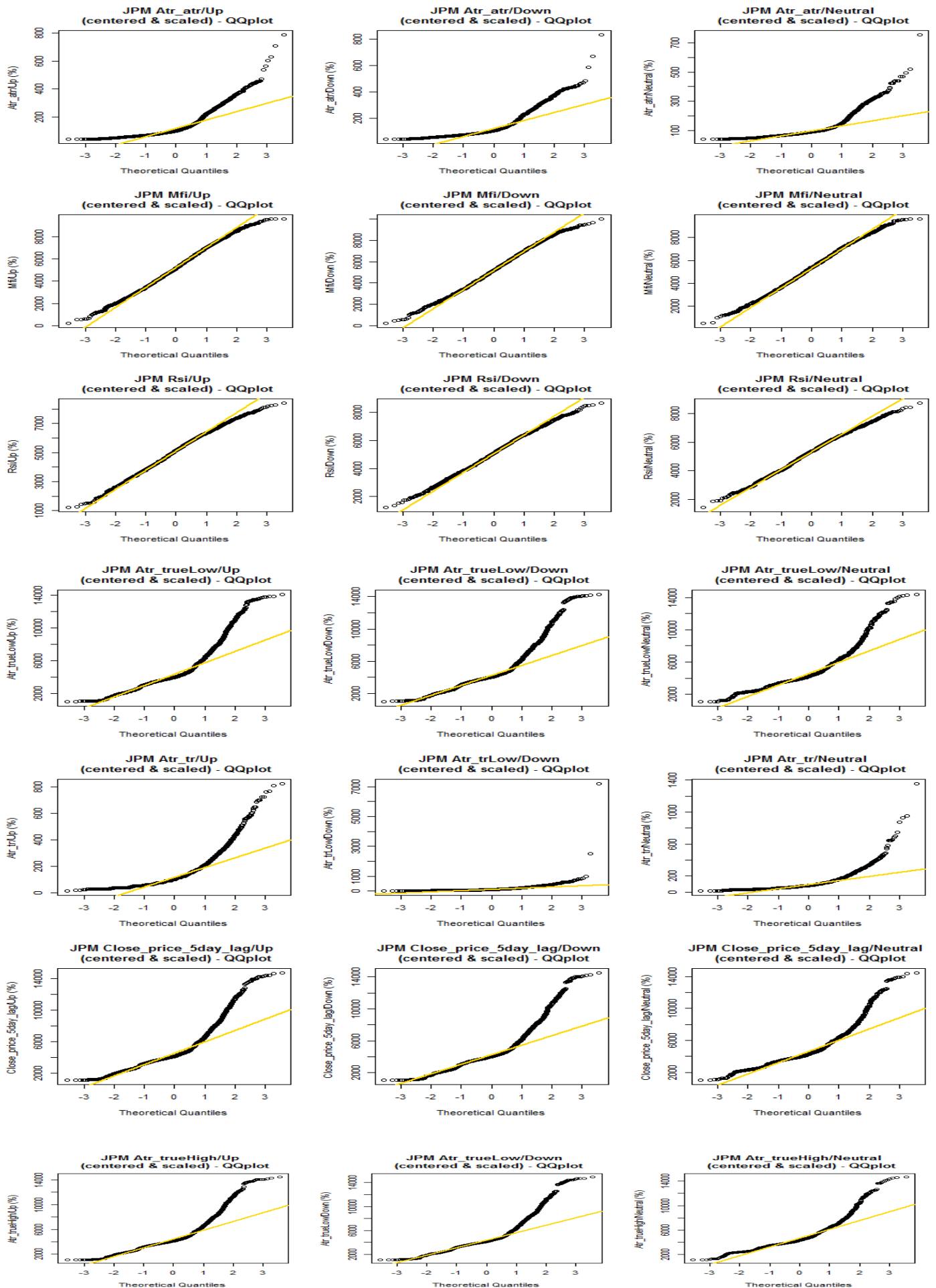
932 ##### Close_price_1day_lag KS Test #####
933 df = subset(subset(normality_check_df, select=c(Direction,Close_price_1day_lag)), Direction == 'up')
934 kolmogorov_smirnov_normal_distribution_test(df$Close_price_1day_lag, stock_name, paste("check ", stock_name, " 'close_price_1day_lag/up' is normally
distributed", sep=""))
935 plot_qqplot (df$Close_price_1day_lag, stock_name, "Close_price_1day_lag/up")
936
937 df = subset(subset(normality_check_df, select=c(Direction,Close_price_1day_lag)), Direction == 'Down')
938 kolmogorov_smirnov_normal_distribution_test(df$Close_price_1day_lag, stock_name, paste("check ", stock_name, " 'close_price_1day_lag/Down' is normally
distributed", sep=""))
939 plot_qqplot (df$Close_price_1day_lag, stock_name, "close_price_1day_lag/Down")
940
941 df = subset(subset(normality_check_df, select=c(Direction,Close_price_1day_lag)), Direction == 'Neutral')
942 kolmogorov_smirnov_normal_distribution_test(df$Close_price_1day_lag, stock_name, paste("check ", stock_name, " 'close_price_1day_lag/Neutral' is normally
distributed", sep=""))
943 plot_qqplot (df$Close_price_1day_lag, stock_name, "close_price_1day_lag/Neutral")
944
945

```

Code Snippet 8 – Kolmogorov Smirnov Test and QQ plot calling code

The detailed results of the Kolmogorov Smirnov Test can be found in Appendix A. In a nutshell, only the JPM 'Rsi/Down', 'Mfi/Up', 'Mfi/Down' and 'Mfi/Neutral' data set seem to follow a normal distribution. None of the others do. This is also visible from the below following graphs.





Models Training & Hyper Parameters Tuning

Due to the cyclical nature of the stock market time series, it is not advisable to use cross-validation, neither for the model fitting or the hyper-parameter optimisation. Each model is therefore performance tested following the methodology described in the above section named *The Sliding Time Windows*. It is worth mentioning that during the training phase, each model is tested against a different list of feature selected attributes. The first training model lists all the 10 features that have been feature selected. Each following model removes one attribute at a time; the one with the least predictive power. The last instance represents a polynomial and interaction between a few of the most predictive power attribute, in an attempt to improve model fitting.

For each model, the following list of attributes are tested (in the order presented):

- Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh
- Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag
- Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr
- Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow
- Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi
- Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi
- Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr
- Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag
- Direction ~ Close_price_1day_lag + Volume
- Direction ~ Close_price_1day_lag
- Direction ~ Close_price_1day_lag^3 + Volume^2 + Close_price_1day_lag + Volume + (Atr_atr * Mfi) + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh

Unless otherwise stated, the R default value of each hyper-parameter is used during the training phase. The optimisation phase takes the best trained model and attempt to optimise its hyper-parameters. As a general principal, when several hyperparameter configurations, for a given model, generate the same accuracy rate, the more computationally efficient parameter set is chosen. For example, when the bagging model returns an accuracy rate of 54.62% for trees number set to 1000 or 1500, and the same lambda (a.k.a. shrinkage), then the model is tested with the number of trees set to 1000.

The remaining sections i) describe each model, ii) provide details of any necessary data transformation prior to running the model, iii) list parameters that need tuning, iv) explain the code required for training/optimising and testing each model and v) offers the average test performance accuracy rate for each model.

The Ridge

Model Description

The Ridge is a method that regularise (i.e. constraints) the coefficient estimates of p predictors of a linear model. In other words, it shrinks the coefficients estimates towards zero [14]. Looking at the formula below, the Ridge regression coefficients are the one that minimise:

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2,$$

Where y_i is the expected value, β_0 is the intercept, β_j are the coefficients of each variable x_j , λ ($\lambda \geq 0$) is the tuning parameter.

The aim of the Ridge model is to fit the data and make the RSS small, by making the shrinkage penalty,

$$\lambda \sum_j \beta_j^2, \text{ as small as possible.}$$

Model Assumptions

No specific requirement.

Further data transformation

The glmnet. *glmnet(x,...)* function implementation does not require any further data transformation. It accepts the regressor as a list of character classes. The parameters can be of type character or numeric.

Parameter Tuning

The main parameter tuning is the λ . The α is set to 0 for the Ridge regression, in the glmnet package.

Code Snippet Explanation

Model Training

The *for-loop*, line 1404 iterates through the training data time windows (100 iterations). At each iteration, a training data slice and validation data sets are built for the time window in question (line 1405-1409). The column list is created line 1411 and then passed to the *run_mlr_model(...)* function, alongside a few parameters (e.g. the training and the validation datasets for the time slide period and *the α is set to 0*). On successful run, the model results are stored into an in-memory data frame named *model_comparison_summary_df*. It is tagged with the state *model_run_success = TRUE*. The training accuracy rate is computed at each iteration. The average of the training accuracy is calculated on line 1418. In case of a computational failure, the model description is added to the same in-memory data frame and tagged with the state *model_run_success = FALSE* (line 1416).

The below code snippet only shows one example of a model trained against a given set of attributes. This example is repeated for each list of attributes. An average training accuracy rate is generated for each instance. The complete list of training case is available in the source code.

```
1398     model_name = "Ridge"
1399     model_desc = "Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh"
1400     model_type = "TRAIN"
1401     the_alpha = 0
1402     uuid = UUIDgenerate()
1403
1404     for (tw_index in time_window_seq){
1405         #Get the Training and validation data for the given time window
1406         training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"training_end_index"]
1407         training_data = final_df[training_range,]
1408         validation_range = time_window_df[tw_index,"validation_start_index"]:time_window_df[tw_index,"validation_end_index"]
1409         validation_data = final_df[validation_range,]
1410
1411         columns = c("Direction","Close_price_1day_lag","Volume","Close_price_2day_lag","Atr_atr","Mfi","Rsi","Atr_trueLow","Atr_tr","Close_price_5day_lag","Atr_trueHigh")
1412
1413         possibleError = tryCatch( run_mlr_model(columns, the_alpha, NULL, uuid, stock_name, model_name, model_desc,model_type,
1414                                           model_comparison_summary_df,training_data, validation_data),
1415                               error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
1416         add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
1417     }
1418     generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
1419
```

Code Snippet 9 – Calling the Ridge model

This below function is responsible for generating the confusion matrix, which produces an accuracy rate of the predicted vs actual *Directions*. First a sparse matrix for the training and test model are created, lines 1301\1303. Then the `glmnet(..)` function is called, line 1306. The *family* parameter is set to multinomial and the α is passed as part of the function call. In this case, α is set to 0 (Ridge regression). Depending on the caller need, the *test_data_param*, could correspond to test data (testing phase) or validation data (training and optimisation phase). The *lambda_param* is also defined by the caller. When it is set to NULL (training phase), line 1309 then the minimum λ is used, else the *lambda_param* value is used (optimisation/testing phases). The prediction is performed (line 1312). The confusion matrix is then built and added as an extra row to the *model_comparison_summary_df* dataframe (line 1320).

```

1290 run_mlr_model = function(column_list, alpha_param, lambda_param, uuid, stock_name, model_name, model_desc, model_type, model_comparison_summary_df, training_data_param,
1291   test_data_param){
1292   set.seed(1)
1293
1294   #Selected the necessary columns for the training and test set
1295   mlr_training_data = training_data_param[column_list]
1296   mlr_test_data = test_data_param[column_list]
1297
1298   #Get the list of attributes col names
1299   col_names = colnames(mlr_training_data[-1])
1300   #Create a model for each set
1301   training_model = sparse.model.matrix( as.formula(paste("Direction ~", paste(col_names, sep = "", collapse=" +"))),
1302     data = mlr_training_data)
1303   test_model = sparse.model.matrix( as.formula(paste("Direction ~", paste(col_names, sep = "", collapse=" +"))),
1304     data = mlr_test_data)
1305
1306   mlr_fit = glmnet(training_model[1:nrow(mlr_training_data)], mlr_training_data$Direction, family = "multinomial", alpha=alpha_param)
1307
1308   #Generate the model prediction depending on the lambda level
1309   if (is.null(lambda_param)) {
1310     lambda_param = min(mlr_fit$lambda)
1311   }
1312   mlr.pred = predict(mlr_fit, as.matrix(test_model), type="class", s=lambda_param)
1313
1314   #Generate the confusion matrix and calculate the classification error rate on the test/validation data
1315   mlr.confusion_table = table(mlr.pred, test_data_param$Direction)
1316   mlr.accuracy_rate = accuracy_rate_perc(mlr.confusion_table)
1317   mlr.error_rate = error_rate_perc(mlr.confusion_table)
1318
1319   #Add a row in the model comparison dataframe
1320   model_comparison_summary_df <- add_row_to_model_summary( model_comparison_summary_df,
1321     uuid,
1322     stock_name,
1323     "TRUE",
1324     "",
1325     model_name, model_desc, model_type,
1326     mlr.accuracy_rate, mlr.error_rate)
1327
1328   return (model_comparison_summary_df)
1329 }
1330
1331 }

```

Code Snippet 9 – Calling the Ridge model for the training phase.

Model Optimisation

This code is similar to the training phase above. This time a list of λ are tested, line 1630/1633 on the last training window. The λ that generates the minimum accuracy error is retained for the testing phase.

```

1624 #####
1625 #####
1626 ##### The Ridge Model - Model Optimisation
1627 #####
1628 #####
1629 |
1630 lambda_list = c(0.0001, 0.0005, 0.0010, 0.0015, 0.0020, 0.0050, 0.0055, 0.0060, 0.01, 0.02, 0.03, 0.04, 0.5, 1)
1631
1632 #The optimisation is performed on the last sliding window
1633 for (the_lambda in lambda_list){
1634   model_name = "Ridge"
1635   model_desc = paste( "Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr +
1636     Close_price_5day_lag + Atr_trueHigh | lambda = ", toString(the_lambda), sep="")
1637   model_type = "OPTIMISATION"
1638   the_alpha = 0
1639   uuid = UUIDgenerate()
1640
1641   #Get the Training and Validation data for the given time window
1642   training_range = time_window_df[number_sliding_windows, "training_start_index"]:time_window_df[number_sliding_windows, "training_end_index"]
1643   training_data = final_df[training_range,]
1644   validation_range = time_window_df[number_sliding_windows, "validation_start_index"]:time_window_df[number_sliding_windows, "validation_end_index"]
1645   validation_data = final_df[validation_range,]
1646
1647   columns = c("Direction", "Close_price_1day_lag", "Volume", "Close_price_2day_lag", "Atr_atr", "Mfi", "Rsi", "Atr_trueLow", "Atr_tr",
1648     "Close_price_5day_lag", "Atr_trueHigh")
1649
1650   possibleError = tryCatch( run_mlr_model(columns, the_alpha, the_lambda, uuid, stock_name, model_name, model_desc, model_type,
1651     model_comparison_summary_df, training_data, validation_data),
1652     error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
1653   add_failed_model(model_comparison_summary_df, uuid, stock_name, model_name, model_desc, model_type, possibleError)
1654 }

```

Code Snippet 10 – Calling the Ridge model for the optimisation phase

Model Testing

The best training model is chosen alongside its hyper-parameter list. In this case *lambda_param* (line 1661) is set to the best optimised value (i.e. 0.001). It is then fitted against the entire training data set, and finally tested against the test set (i.e. the last 5 business days). The confusion matrix, that evaluates the test prediction vs the expected test data, provides the test accuracy rate. The same *run_mlr_model()* function is called as for the training phase, line 1678. This time, the test data is used in lieu of the validation data.

```
1656 #####
1657 ##### The Ridge Model - Model Testing
1658 #####
1659 #####
1660 model_name = "Ridge"
1661 the_lambda = 0.001
1662 model_desc = paste("Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr +
Close_price_5day_lag + Atr_trueHigh | lambda = ", toString(the_lambda), sep="")
1663 model_type = "TEST"
1664 the_alpha = 0
1665 uuid = UUIDgenerate()
1666
1667 #This time the model is tained on the Traning + validation data, then it is tested against the Testing data.
1668 for (tw_index in time_window_seq){
1669   #Get the Training and Validation data for the given time window
1670   training_range = time_window_df[tw_index, "training_start_index"]:time_window_df[tw_index, "validation_end_index"]
1671   training_data = final_df[training_range,]
1672   test_range = time_window_df[tw_index, "test_start_index"]:time_window_df[tw_index, "test_end_index"]
1673   test_data = final_df[validation_range,]
1674
1675   columns = c("Direction", "Close_price_1day_lag", "Volume", "Close_price_2day_lag", "Atr_atr", "Mfi", "Rsi", "Atr_trueLow", "Atr_tr",
"Close_price_5day_lag", "Atr_trueHigh")
1676
1677   possibleError = tryCatch( run_mlr_model(columns, the_alpha, the_lambda, uuid, stock_name, model_name, model_desc, model_type,
1678     model_comparison_summary_df, training_data, validation_data,
1679     error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
1680   add_failed_model(model_comparison_summary_df, uuid, stock_name, model_name, model_desc, model_type, possibleError)
1681 }
1682 generate_avg_model(model_comparison_summary_df, uuid, time_window_seq, model_type)
1683
1684 |
1685 cat ("\n")
1686 cat (paste ("Save Comparison Summary: ", model_name, "\n", sep=""))
1687 selected_model_df = model_comparison_summary_df[model_comparison_summary_df$model_name == model_name, ]
1688 print(selected_model_df)
1689 save_to_csv(selected_model_df, out_dir, paste("Model_Results_", model_name, ".csv", sep=""))
1690
```

Code Snippet 10 – Calling the Ridge model for the testing phase

Results

As shown in the below table, the Ridge model with the following configuration: *Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh* generates the highest training accurate rate at 63.17%. The optimisation model shows the best accuracy for $\lambda = 0.001$, at 73.85%. The test performance of the Ridge model, run against the selected list of attributes and the optimised λ , produces a 75.39% accuracy rate, with a standard deviation of 1.44%. The data is available in the file: *.../final results/Model_Results_Ridge.xlsx*

Row Labels	Sum of model_accuracy_rate	Std Dev
TEST		
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.001	75.39	1.44
TRAIN		
Direction ~ Direction ~ Close_price_1day_lag		28.58
Direction ~ Direction ~ Close_price_1day_lag + Volume		29.20
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag		29.37
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr		31.13
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi		32.41
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi		32.60
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow		59.08
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr		59.27
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag		60.28
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh		63.17
aggregation_type	(blank)	
Row Labels	Sum of model_accuracy_rate	
OPTIMISATION		
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.001	73.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.0015	73.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.002	73.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.005	73.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.0055	73.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.006	73.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.01	73.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.02	71.54	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.03	70.00	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.04	69.23	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.5	59.23	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.0001	62.31	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 1e-04	73.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 5e-04	73.85	

The Lasso

Model Description

Like the Ridge, the Lasso is a method that regularise (i.e. constraints) the coefficient estimates of p predictors of a linear model. Looking at the formula below, the Lasso regression coefficients are the one that minimise:

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| = \text{RSS} + \lambda \sum_{j=1}^p |\beta_j|.$$

Where y_i is the expected value, β_0 is the intercept, β_j are the coefficients of each variables x_j , λ ($\lambda \geq 0$) is the tuning parameter. This time, the ℓ_1 parameter: $|\beta_j|$ is used instead of the ℓ_2 parameter: β_j^2 (as per the Ridge mode). The advantage is that the ℓ_1 parameter has the effect of forcing the coefficient estimate to zero when λ is large enough [14].

Model Assumptions

No specific requirement.

Further data transformation

The `glmnet`. `glmnet(x,...)` function implementation does not require any further data transformation. It accepts the regressor as a list of character classes. The parameters can be of type character or numeric.

Parameter Tuning

The main parameter tuning is the λ . The α is set to 1 for the Lasso regression, in the `glmnet` package.

Code Snippet Explanation

Model Training

Same description as the Ridge section. The only difference is that α is set to 1 (Line 1700).

```
1692 #####
1693 #####
1694 ##### The Lasso Model - Model Training
1695 #####
1696 #####
1697 model_name = "Lasso"
1698 model_desc = "Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag +
Atr_trueHigh"
1699 model_type = "TRAIN"
1700 the_alpha = 1
1701 uuid = UUIDgenerate()
1702
1703 for (tw_index in time_window_seq){
1704   #Get the Training and Validation data for the given time window
1705   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"training_end_index"]
1706   training_data = final_df[training_range,]
1707   validation_range = time_window_df[tw_index,"validation_start_index"]:time_window_df[tw_index,"validation_end_index"]
1708   validation_data = final_df[validation_range,]
1709
1710   columns = c("Direction","close_price_1day_lag","Volume","close_price_2day_lag","Atr_atr","Mfi","Rsi","Atr_trueLow","Atr_tr","close_price_5day_lag","Atr_trueHigh")
1711
1712   possibleError = trycatch( run_mlr_model(columns, the_alpha, NULL, uuid, stock_name, model_name, model_desc,model_type,
1713                                     model_comparison_summary_df,training_data, validation_data),
1714                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
1715   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
1716 }
1717 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
1718
```

Code Snippet 11 – Calling the Lasso model (training)

Model Optimisation

Same description as the Ridge section. The only difference is that α is set to 1 (Line 1936).

```
1922 #####
1923 #####
1924 #####
1925 ##### The Lasso Model - Model Optimisation
1926 #####
1927 #####
1928
1929 lambda_list = c(0.0001, 0.0005, 0.0010,0.0015,0.0020,0.0050, 0.0055, 0.0060, 0.01, 0.02, 0.03, 0.04, 0.5,1)
1930
1931 #The optimisation is performed on the last sliding window
1932 for (the_lambda in lambda_list){
1933   model_name = "Lasso"
1934   model_desc = paste( "Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr | lambda = ",
1935   toString(the_lambda),sep="")
1936   model_type = "OPTIMISATION"
1937   the_alpha = 1
1938   uuid = UUIDgenerate()
1939
1940   #Get the Training and Validation data for the given time window
1941   training_range = time_window_df[number_sliding_windows,"training_start_index"]:time_window_df[number_sliding_windows,"training_end_index"]
1942   training_data = final_df[training_range,]
1943   validation_range = time_window_df[number_sliding_windows,"validation_start_index"]:time_window_df[number_sliding_windows,"validation_end_index"]
1944   validation_data = final_df[validation_range,]
1945
1946   columns = c("Direction","Close_price_1day_lag","Volume","Close_price_2day_lag","Atr_atr")
1947
1948   possibleError = tryCatch( run_mlr_model(columns, the_alpha, the_lambda, uuid, stock_name, model_name, model_desc,model_type,
1949   model_comparison_summary_df,training_data, validation_data),
1950   error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
1951   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
1952 }
```

Code Snippet 12– Calling the Lasso model (Optimisation)

Model Testing

The best training model is chosen alongside its hyper-parameter list. In this case lambda_param (line 1960) is set to the best optimised value (i.e. 0.005). It is then fitted against the entire training data set, and finally tested against the test set (i.e. the last 5 business days). The confusion matrix, that evaluates the test prediction vs the expected test data, provides the test accuracy rate. The same run_mlr_model() function is called as for the training phase, line 1976. This time, the test data is used in lieu of the validation data.

```
1953 #####
1954 #####
1955 ##### The Lasso Model - Model Testing
1956 #####
1957 #####
1958 model_name = "Lasso"
1959 the_alpha = 1
1960 the_lambda = 0.005
1961 model_desc = paste( "Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr +
1962 Close_price_5day_lag + Atr_trueHigh | lambda = ", toString(the_lambda),sep="")
1963 model_type = "TEST"
1964 uuid = UUIDgenerate()
1965
1966 #This time the model is tained on the Traning + Validation data, then it is tested against the Testing data.
1967 for (tw_index in time_window_seq){
1968   #Get the Training and Validation data for the given time window
1969   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"validation_end_index"]
1970   training_data = final_df[training_range,]
1971   test_range = time_window_df[tw_index,"test_start_index"]:time_window_df[tw_index,"test_end_index"]
1972   test_data = final_df[validation_range,]
1973
1974   columns = c("Direction","Close_price_1day_lag","Volume","Close_price_2day_lag","Atr_atr","Mfi","Rsi","Atr_trueLow","Atr_tr",
1975   "Close_price_5day_lag","Atr_trueHigh")
1976
1977   possibleError = tryCatch( run_mlr_model(columns, the_alpha, the_lambda, uuid, stock_name, model_name, model_desc,model_type,
1978   model_comparison_summary_df,training_data, validation_data),
1979   error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
1980   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
1981 }
1982 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
1983
1984 cat ("\n")
1985 cat (paste ("Save Comparison Summary: ", model_name, "\n", sep=""))
1986 selected_model_df = model_comparison_summary_df[model_comparison_summary_df$model_name == model_name, ]
1987 print(selected_model_df)
1988 save_to_csv(selected_model_df, out_dir,paste("Model_Results_", model_name, ".csv", sep=""))
1989 }
```

Code Snippet 12– Calling the Lasso model (Testing)

Results

As shown in the below table, the Lasso model with the following configuration: *Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh* generates the highest training accurate rate at 73.59%. The optimisation model shows the best accuracy for lambda= 0.005, at 22.31%. The test performance of the Lasso model, run against the selected list of attributes and the optimised λ , produces a 80.71% accuracy rate, with a standard deviation of 1.38%. The data is available in the file: *Model_Lasso.xlsx*

K		L	M
aggregation_type		AVG	Y
Row Labels		Sum of model_accuracy_rate	Std Dev
TEST			
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh lambda = 0.005		80.71	1.38
TRAIN			
Direction ~ Direction ~ Close_price_1day_lag		28.58	
Direction ~ Direction ~ Close_price_1day_lag + Volume		29.19	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag		29.61	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr		30.98	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi		32.62	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi		32.56	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow		71.86	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr		73.88	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag		73.62	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh		73.59	
aggregation_type		(blank)	Y
Row Labels		Sum of model_accuracy_rate	
OPTIMISATION			
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.001		20.00	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.0015		20.77	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.002		21.54	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.005		22.31	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.0055		22.31	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.006		22.31	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.01		19.23	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.02		16.92	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.03		16.15	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.04		13.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 0.5		13.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 1		13.85	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 1e-04		19.23	
Direction ~ Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr lambda = 5e-04		19.23	

Linear Discriminant Analysis (LDA)

Model Description

As described in [14], the multivariate Gaussian density function for the LDA model when $p > 1$ (p is the number of attributes) can be described as follows. We assume that $X = (x_1, x_2, \dots, x_n)$ is a vector of predictors. This model assumes that each individual predictor follows a one-dimensional normal distribution, with some correlation between each pair of predictors". Therefore, $X \sim N(\mu, \Sigma)$, where $E(X) = \mu$ is the mean of X (a vector of p components) and $Cov(X) = \Sigma$ is the covariance matrix of X . The estimates $f(x)$ are then fed into the *Bayes' Theorem* to perform the class prediction.

$$f(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right).$$

Model Assumptions

Each attribute class should follow a normal (Gaussian) distribution.

Further data transformation

The `MASS.lda(x,...)` function implementation does not require any further data transformation. It accepts the regressor as a list of character classes. The parameters can be of type character or numeric.

Parameter Tuning

There is no need for parameter tuning for this model.

Code Snippet Explanation

Model Training

The *for-loop*, line 2001 iterates through the training data time windows (100 iterations). At each iteration, the *lda(x,...)* function is called with a list of explanatory and explained variables, for a training data slice (line 2008). It produces the *lda_fit* object, i.e. the function fitting the model. The *lda_fit* model is then passed to the *run_lda_model(...)* function, alongside a few parameters. One of them is the validation set for the time slide period. On successful run, the model results are stored into an in-memory object named *model_comparison_summary_df*. It is tagged with the state *model_run_success = TRUE*. The training accuracy rate is computed at each iteration. The average of the training accuracy is calculated (line 2024). In case of a computational failure, the model description is added to the same in-memory data frame and tagged with the state *model_run_success = FALSE* (line 2022). The below code snippet only shows one example of a model trained against a given set of attributes. This example is repeated for each list of attributes. An average training accuracy rate is generated for each instance. The complete list of training case is available in the source code.

```
1990 #####
1991 #####
1992 ##### Linear Discriminant Analysis (LDA) - Model Training
1993 #####
1994 #####
1995 #####
1996 model_name = "Linear Discriminant Analysis (LDA)"
1997 model_desc = "Direction ~ Close_price_1day_lag + volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag +
1998 Atr_trueHigh"
1999 model_type = "TRAIN"
2000 uuid = UUIDgenerate()
2001 for (tw_index in time_window_seq){
2002   #Get the training and validation data for the given time window
2003   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"training_end_index"]
2004   training_data = final_df[training_range,]
2005   validation_range = time_window_df[tw_index,"validation_start_index"]:time_window_df[tw_index,"validation_end_index"]
2006   validation_data = final_df[validation_range,]
2007
2008   lda_fit = lda(Direction ~ Close_price_1day_lag +
2009                 volume +
2010                 Close_price_2day_lag +
2011                 Atr_atr +
2012                 Mfi +
2013                 Rsi +
2014                 Atr_trueLow +
2015                 Atr_tr +
2016                 Close_price_5day_lag +
2017                 Atr_trueHigh,
2018                 data=training_data)
2019
2020   possibleError = tryCatch( run_lda_model(lda_fit,uuid,stock_name,model_name,model_desc,model_type,model_comparison_summary_df,validation_data),
2021                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
2022   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
2023 }
2024 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
2025
```

Code Snippet 13– Calling the LDA model (Training)

This below function is responsible for generating the confusion matrix, which produces an accuracy rate of the predicted vs actual *Directions*. The validation data set is used for this purpose. The accuracy result is stored for each iteration in the in-memory table (Line 1114/1120).

```

1102 #the test_data_param could be a training/validation or test dataset
1103 run_lda_model = function(fit_param, uuid, stock_name, model_name, model_desc, model_type, model_comparison_summary_df, test_data_param, save_confusion_matrix=FALSE){
1104   lda.fit = fit_param
1105
1106   #Generate the confusion matrix and calculate the training/validation error rate
1107   lda.pred = predict(lda.fit, test_data_param)$class
1108   lda.confusion_table = table(lda.pred, test_data_param$Direction)
1109   lda.accuracy_rate = accuracy_rate_perc(lda.confusion_table)
1110   lda.error_rate = error_rate_perc(lda.confusion_table)
1111
1112   #Add a row in the model comparison dataframe
1113   model_comparison_summary_df <- add_row_to_model_summary( model_comparison_summary_df,
1114                                                         uuid,
1115                                                         stock_name,
1116                                                         "TRUE",
1117                                                         "",
1118                                                         model_name, model_desc, model_type,
1119                                                         lda.accuracy_rate, lda.error_rate)
1120
1121   if (save_confusion_matrix == TRUE){
1122     save_to_csv(lda.confusion_table, out_dir, paste("lda_confusion_matrix_", UUIDgenerate(), ".csv", sep=""))
1123   }
1124   return (model_comparison_summary_df)
1125 }
1126
1127
1128

```

Code Snippet 14– The *run_lda_model()* function implementation

Model Testing

The best training model is chosen and is fitted against the entire training data set. It is then tested against the test set (i.e. the last 5 business days). The confusion matrix that evaluates the test prediction vs the expected test data provides the test accuracy rate. The same *run_lda_model()* function is called as for the training phase. This time, the test data is used in lieu of the validation data.

```

2296 model_name = "Linear Discriminant Analysis (LDA)"
2297 model_desc = "Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr"
2298 model_type = "TEST"
2299 uuid = UUIDgenerate()
2300
2301 #This time the model is trained on the Training + Validation data, then it is tested against the Testing data.
2302 for (tw_index in time_window_seq){
2303   #Get the Training and Validation data for the given time window
2304   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"validation_end_index"]
2305   training_data = final_df[training_range,]
2306   test_range = time_window_df[tw_index,"test_start_index"]:time_window_df[tw_index,"test_end_index"]
2307   test_data = final_df[validation_range,]
2308
2309   lda_fit = lda (Direction ~ Close_price_1day_lag +
2310                 volume +
2311                 Close_price_2day_lag +
2312                 Atr_atr +
2313                 Mfi +
2314                 Rsi +
2315                 Atr_trueLow +
2316                 Atr_tr,
2317                 data=training_data)
2318
2319   possibleError = tryCatch( run_lda_model(lda_fit,uuid,stock_name,model_name,model_desc,model_type,model_comparison_summary_df,test_data, TRUE),
2320                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
2321   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
2322 }
2323 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
2324
2325 cat ("\n")
2326 cat (paste ("Save Comparison Summary: ", model_name, "\n", sep=""))
2327 selected_model_df = model_comparison_summary_df[model_comparison_summary_df$model_name == model_name, ]
2328 print(selected_model_df)
2329 save_to_csv(selected_model_df, out_dir, paste("Model_Results_", model_name, ".csv", sep=""))
2330

```

Code Snippet 15– Calling the LDA model (Testing)

Results

As shown in the below table, the LDA model with the following configuration: *Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr* generates the highest training accurate ate at 72.98%. The test performance of the LDA model, run against the selected list of attributes, produces an 80.20% accuracy rate, with standard deviation of 1.06%. The data is available in the file: *Model_Results_Linear Discriminant Analysis (LDA).xlsx*

Row Labels	Sum of model_accuracy_rate	Std Dev
TEST		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr	80.20	1.06
TRAIN		
Direction ~ Close_price_1day_lag	39.45	
Direction ~ Close_price_1day_lag + Volume	41.52	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag	42.32	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr	41.38	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi	39.99	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi	39.68	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow	69.98	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr	72.98	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag	72.41	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh	72.41	

Quadratic Discriminant Analysis (QDA)

Model Description

As described in [14], the QDA classifier model follows the same assumption as the LDA model, and is also dependent on the Bayes' theorem to perform class predictions. The main difference resides in the assumption that each class has its own covariance matrix. Therefore, $X \sim N(\mu_k, \Sigma_k)$ where k represents the k th class. The class separator becomes quadratic instead of being linear, hence the name.

Model Assumptions

Each class should follow a normal (Gaussian) distribution.

Further data transformation

The R *MASS.qda* (x, \dots) function implementation does not require any further data transformation. It accepts the regressor as a list of character classes. The parameters can be of type character or numeric.

Parameter Tuning

There is no need for parameter tuning for this model.

Code Snippet Explanation

Model Training

The *for-loop*, line 2372 iterates through the training data time windows (100 iterations). At each iteration, the *qda*(x, \dots) function is called with a list of explanatory and explained variables, for a training data slice (line 2379). It produces the *qda_fit* object, i.e. the function fitting the model. The *qda_fit* model is then passed to the *run_qda_model*(\dots) function, alongside a few parameters. One of them is the validation set for the time slide period. On successful run, the model results are stored into an in-memory object named *model_comparison_summary_df*. It is tagged with the state *model_run_success* = *TRUE*. The training accuracy rate is computed at each iteration. The average of the training accuracy is calculated on line 2394. In case of a computational failure, the model description is added to the same in-memory data frame and tagged with the state *model_run_success* = *FALSE* (line 2392). The below code snippet only shows one example of a model trained against a given set of attributes. This example is repeated for each list of attributes. An average training accuracy rate is generated for each instance. The complete list of training case is available in the source code.

```
2367 model_name = "Quadratic Discriminant Analysis (QDA)"
2368 model_desc = "Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag"
2369 model_type = "TRAIN"
2370 uuid = UUIDgenerate()
2371
2372 for (tw_index in time_window_seq){
2373   #Get the Training and Validation data for the given time window
2374   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"training_end_index"]
2375   training_data = final_df[training_range,]
2376   validation_range = time_window_df[tw_index,"validation_start_index"]:time_window_df[tw_index,"validation_end_index"]
2377   validation_data = final_df[validation_range,]
2378
2379   qda_fit = qda (Direction ~ Close_price_1day_lag +
2380                 Volume +
2381                 Close_price_2day_lag +
2382                 Atr_atr +
2383                 Mfi +
2384                 Rsi +
2385                 Atr_trueLow +
2386                 Atr_tr +
2387                 Close_price_5day_lag,
2388                 data=training_data)
2389
2390   possibleError = trycatch( run_qda_model(qda_fit,uuid,stock_name,model_name,model_desc,model_type,model_comparison_summary_df,validation_data),
2391                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
2392   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
2393 }
2394 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
2395 }
```

Code Snippet 16 – Calling the QDA model (Training)

This below function is responsible for generating the confusion matrix, which produces an accuracy rate of the predicted vs actual *Directions*. The validation data set is used for this purpose. The accuracy result is stored for each iteration in the in-memory table (line 1140/1146).

```

1129 run_qda_model = function(fit_param, uuid, stock_name, model_name, model_desc,model_type, model_comparison_summary_df,test_data_param, save_confusion_matrix=FALSE){
1130
1131   qda.fit = fit_param
1132
1133   #Generate the confusion matrix and calculate the training/validation error rate
1134   qda.pred = predict(qda.fit, test_data_param)$class
1135   qda.confusion_table = table(qda.pred, test_data_param$Direction)
1136   qda.accuracy_rate = accuracy_rate_perc(qda.confusion_table)
1137   qda.error_rate = error_rate_perc(qda.confusion_table)
1138
1139   #Add a row in the model comparison dataframe
1140   model_comparison_summary_df <-> add_row_to_model_summary( model_comparison_summary_df,
1141                                                            uuid,
1142                                                            stock_name,
1143                                                            "TRUE",
1144                                                            "",
1145                                                            model_name, model_desc,model_type,
1146                                                            qda.accuracy_rate, qda.error_rate)
1147
1148   if (save_confusion_matrix == TRUE){
1149     save_to_csv(qda.confusion_table, out_dir,paste("qda_confusion_matrix_", UUIDgenerate(), ".csv",sep=""))
1150   }
1151
1152   return (model_comparison_summary_df)
1153 }
1154

```

Code Snippet 17 – The *run_qda_model()* function implementation

Model Testing

The best training model is chosen and is fitted against the entire training data set. It is then tested against the test set (i.e. the last 5 business days). The confusion matrix, that evaluates the test prediction vs the expected test data, provides the test accuracy rate. The same *run_qda_model()* function is called as for the training phase. This time, the test data is used in lieu of the validation data.

```

2629 #####
2630 #####
2631 ##### Quadratic Discriminant Analysis (QDA) - Model Testing
2632 #####
2633 #####
2634
2635 model_name = "Quadratic Discriminant Analysis (QDA)"
2636 model_desc = "Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag"
2637 model_type = "TEST"
2638 uuid = UUIDgenerate()
2639
2640 #This time the model is tained on the Traning + Validation data, then it is tested against the Testing data.
2641 for (tw_index in time_window_seq){
2642   #Get the Training and Validation data for the given time window
2643   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"validation_end_index"]
2644   training_data = final_df[training_range,]
2645   test_range = time_window_df[tw_index,"test_start_index"]:time_window_df[tw_index,"test_end_index"]
2646   test_data = final_df[validation_range,]
2647
2648   qda_fit = qda (Direction ~ Close_price_1day_lag +
2649                 Volume +
2650                 Close_price_2day_lag +
2651                 Atr_atr +
2652                 Mfi +
2653                 Rsi +
2654                 Atr_trueLow +
2655                 Atr_tr +
2656                 Close_price_5day_lag,
2657                 data=training_data)
2658
2659   possibleError = tryCatch( run_qda_model(qda_fit,uuid,stock_name,model_name,model_desc,model_type,model_comparison_summary_df,test_data, TRUE),
2660                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
2661   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
2662 }
2663 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
2664

```

Code Snippet 18– Calling the QDA model (Testing)

Results

As shown in the below table, the LDA model with the following configuration: *Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag* generates the highest training accurate rate at 68.48%. The test performance of the LDA model run against the selected list of attributes produces a 80.90% accuracy rate, with a standard deviation of 1.64%. The data is available in the file: *Model_Results_Quadratic Discriminant Analysis (QDA).xlsx*

aggregation_type	AVG	
Row Labels	Sum of model_accuracy_rate	
TEST		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag	80.90	1.64
TRAIN		
Direction ~ Close_price_1day_lag	36.97	
Direction ~ Close_price_1day_lag + Volume	38.71	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag	40.96	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr	40.52	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi	40.24	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi	39.32	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow	66.49	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr	67.92	
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag	68.48	

Decision Tree

Model Description

The goal is to find regions R_1, \dots, R_n that minimise the classification error rate. It corresponds to the fraction of the training observations in that region that do not belong to the most common class. It is given by the equation from [14], $p_{\text{Hat}_{mk}}$ represents the portion of training observations in the m th region that are from the k th class.

$$E = 1 - \max_k (\hat{p}_{mk}).$$

Model Assumptions

No specific requirements.

Further data transformation

The `tree.tree(x,...)` function implementation requires that the regressor variable is encoded as a factor (i.e. an enumerated type). Therefore, the `factor()` function has been applied on the *Direction* regressor. None of the other parameters need to be adapted.

Parameter Tuning

The tree size parameter is optimised.

Code Snippet Explanation

Model Training

The *for-loop*, line 3144 iterates through the training data time windows (100 iterations). At each iteration, the `tree(x,...)` function is called with a list of explanatory and explained variables, for a training data slice (line 3151). It produces the `tree_fit` object, i.e. the function fitting the model. The `tree_fit` model is then passed to the `run_tree_model(...)` function, alongside a few parameters. One of them is the validation set for the time slide period. On successful run, the model results are stored into an in-memory object named `model_comparison_summary_df`. It is tagged with the state `model_run_success = TRUE`. The training accuracy rate is computed at each iteration. The average of the training accuracy is calculated on line 3167. In case of a computational failure, the model description is added to the same in-memory data frame and tagged with the state `model_run_success = FALSE` (line 3165). The below code snippet only shows one example of a model trained against a given set of attributes. This example is repeated for each list of attributes. An average training accuracy rate is generated for each instance. The complete list of training case is available in the source code.

```
3134 #####
3135 #####
3136 ##### Decision Tree - Model Training
3137 #####
3138 #####
3139 model_name = "Decision Tree"
3140 model_desc = "Factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag +
Atr_trueHigh"
3141 model_type = "TRAIN"
3142 uuid = UUIDgenerate()
3143
3144 for (tw_index in time_window_seq){
3145   #Get the Training and Validation data for the given time window
3146   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"training_end_index"]
3147   training_data = final_df[training_range,]
3148   validation_range = time_window_df[tw_index,"validation_start_index"]:time_window_df[tw_index,"validation_end_index"]
3149   validation_data = final_df[validation_range,]
3150
3151   tree_fit = tree (factor(Direction) ~ Close_price_1day_lag +
3152                   Volume +
3153                   Close_price_2day_lag +
3154                   Atr_atr +
3155                   Mfi +
3156                   Rsi +
3157                   Atr_trueLow +
3158                   Atr_tr +
3159                   Close_price_5day_lag +
3160                   Atr_trueHigh,
3161                   data=training_data)
3162
3163   possibleError = tryCatch( run_tree_model(tree_fit,uuid,stock_name,model_name,model_desc,model_type,model_comparison_summary_df,validation_data),
3164                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
3165   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
3166 }
3167 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
3168
```

Code Snippet 19 – Calling the Tree model (Training)

This below function is responsible for generating the confusion matrix, which produces an accuracy rate of the predicted vs actual *Directions*. The validation data set is used for this purpose. The accuracy result is stored for each iteration in the in-memory table (line 1217/1124).

```

1203- run_tree_model = function(fit_param,uuid, stock_name, model_name, model_desc,model_type, model_comparison_summary_df,test_data_param, can_prune=TRUE){
1204
1205     set.seed(1)
1206
1207     tree.fit = fit_param
1208     #Type="class" is selected as we are dealing with a classification problem
1209     tree.pred = predict(tree.fit, test_data_param, type="class")
1210
1211     #Generate the confusion matrix and calculate the training/validation error rate
1212     tree.confusion_table = table(tree.pred, test_data_param$Direction)
1213     tree.accuracy_rate = accuracy_rate_perc(tree.confusion_table)
1214     tree.error_rate = error_rate_perc(tree.confusion_table)
1215
1216     #Add a row in the model comparison dataframe
1217     model_comparison_summary_df <- add_row_to_model_summary( model_comparison_summary_df,
1218                                                            uuid,
1219                                                            stock_name,
1220                                                            "TRUE",
1221                                                            "",
1222                                                            model_name, model_desc,model_type,
1223                                                            tree.accuracy_rate, tree.error_rate)
1224
1225     return (model_comparison_summary_df)
1226 }

```

Code Snippet 20 – Calling the *run_tree_model()* function

Model Optimisation

This code is similar to the training phase above. This time the tree minimum size list (*min_size_list*) is provided (line 1434). The optimisation is run over the last training period and generate the accuracy for each tree minimum size element in the list. The highest accuracy is retained. The minimum tree size is selected for this accuracy level.

```

1434 min_size_list = c(1,2,5,10,20,40,60,80,100)
1435
1436 #The optimisation is performed on the last sliding window
1437- for (the_min_size in min_size_list){
1438     model_name = "Decision Tree"
1439     model_desc = paste( "factor(Direction) ~ Close_price_1day_lag + volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr +
Close_price_5day_lag + Atr_trueHigh | min size = ", toString(the_min_size),sep="" )
1440     model_type = "OPTIMISATION"
1441     uuid = UUIDgenerate()
1442
1443     #Get the Training and validation data for the given time window
1444     training_range = time_window_df[number_sliding_windows,"training_start_index"]:time_window_df[number_sliding_windows,"training_end_index"]
1445     training_data = final_df[training_range,]
1446     validation_range = time_window_df[number_sliding_windows,"validation_start_index"]:time_window_df[number_sliding_windows,"validation_end_index"]
1447     validation_data = final_df[validation_range,]
1448
1449     tree_fit = tree (factor(Direction) ~ Close_price_1day_lag +
1450                                     volume +
1451                                     Close_price_2day_lag +
1452                                     Atr_atr +
1453                                     Mfi +
1454                                     Rsi +
1455                                     Atr_trueLow +
1456                                     Atr_tr +
1457                                     Close_price_5day_lag +
1458                                     Atr_trueHigh,
1459                       data=training_data)
1460
1461     tree.pruned = prune.misclass(tree_fit, best = the_min_size)
1462
1463     possibleError = tryCatch( run_tree_model(tree.pruned,uuid,stock_name,model_name,model_desc,model_type,model_comparison_summary_df,validation_data),
1464                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
1465     add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
1466 }
1467

```

Code Snippet 21 – Calling the Tree model for the optimisation phase

Model Testing

The best training model is chosen alongside its hyper-parameter list. It is then fitted against the entire training data set, and finally tested against the test set (i.e. the last 5 business days). The confusion matrix, that evaluates the test prediction vs the expected test data, provides the test accuracy rate. The same *run_tree_model()* function is called as for the training phase, line 3496. This time, the test data is used in lieu of the validation data.

Results

As shown in the below table, the Ridge model with the following configuration: *Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh* generates the highest training accurate rate at 44.33%. The optimisation model shows the best min tree size = 2. The test performance of the Tree model run against the selected list of attributes and the optimised min tree size produces a 70.80% accuracy rate, with a standard deviation of 21.45%. The data is available in the file: *Model_Results_Tree.xlsx*

K		L	M	N
Row Labels		Sum of model_accuracy_rate	Std Dev	
TEST				
Direction ~ Close_price_1day_lag + Close_price_4day_lag + Close_price_2day_lag + Atr_trueHigh		40.80	21.45	
TRAIN				
factor(Direction) ~ Close_price_1day_lag		38.07		
factor(Direction) ~ Close_price_1day_lag + Volume		38.18		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag		38.31		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr		37.22		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi		37.39		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi		38.90		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow		42.19		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr		41.52		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag		43.79		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh		44.33		
aggregation_type		(blank)		
Row Labels		Sum of model_accuracy_rate		
OPTIMISATION				
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh min size = 10		54.62		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh min size = 100		54.62		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh min size = 2		54.62		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh min size = 20		54.62		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh min size = 40		54.62		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh min size = 5		54.62		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh min size = 60		54.62		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh min size = 80		54.62		

Bagging

Model Description

Bagging is a procedure that reduces the statistical learning method variance [14]. This is achieved by making repeated samples from the training data set (bags) and averaging the prediction accuracy over the bags number.

$$\hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x).$$

With B, the number of bag and $\hat{f}^{*b}(x)$ function, the prediction for each bag.

This methodology can be applied to classification trees, where the class predicted is the most commonly occurring class for a bag. Bagging is a special case randomForest where the number of variables randomly sampled as candidates at each split is equal to the number of attributes in the model (a.k.a. *mtry*).

Model Assumptions

No specific requirements.

Further data transformation

The *randomForest.randomForest(x,...)* function implementation requires that the regressor variable is encoded as a factor (i.e. an enumerated type). Therefore, the *factor()* function has been applied on the *Direction* regressor.

Parameter Tuning

The tree size parameter is optimised.

Code Snippet Explanation

Model Training

The *for-loop*, line 3520 iterates through the training data time windows (100 iterations). At each iteration, the the *randomForest(x,...)* function is called with a list of explanatory and explained variables, for a training data slice (line 3527). It produces the *bagging_fit* object, i.e. the function fitting the model. The *bagging_fit* model is then passed to the *run_tree_model(...)* function, alongside a few parameters. One of them is the validation set for the time slide period. On successful run, the model results are stored into an in-memory object named *model_comparison_summary_df*. It is tagged with the state *model_run_success = TRUE*.

The training accuracy rate is computed at each iteration. The average of the training accuracy is calculated on line 3549. In case of a computational failure, the model description is added to the same in-memory data frame and tagged with the state *model_run_success = FALSE* (line 3547).

The below code snippet only shows one example of a model trained against a given set of attributes. This example is repeated for each list of attributes. An average training accuracy rate is generated for each instance. The complete list of training case is available in the source code.

```
3508 #####
3509 #####
3510 ##### Bagging Model
3511 ##### Specialised Random Forest where mtry = p (number of attributes)
3512 #####
3513 #####
3514 model_name = "Bagging"
3515 model_desc = "factor(Direction) ~ Close_price_1day_lag + volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag +
Atr_trueHigh"
3516 model_type = "TRAIN"
3517 uuid = UUIDgenerate()
3518 m_try=10
3519
3520 for (tw_index in time_window_seq){
3521   #Get the Training and validation data for the given time window
3522   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"training_end_index"]
3523   training_data = final_df[training_range,]
3524   validation_range = time_window_df[tw_index,"validation_start_index"]:time_window_df[tw_index,"validation_end_index"]
3525   validation_data = final_df[validation_range,]
3526
3527   bagging_fit = randomForest( factor(Direction) ~ Close_price_1day_lag +
3528                               volume +
3529                               Close_price_2day_lag +
3530                               Atr_atr +
3531                               Mfi +
3532                               Rsi +
3533                               Atr_trueLow +
3534                               Atr_tr +
3535                               Close_price_5day_lag +
3536                               Atr_trueHigh,
3537                               data=training_data,
3538                               mtry=m_try,
3539                               importance=TRUE)
3540
3541   possibleError = tryCatch( run_rdmForest_model( bagging_fit,n_trees,
3542                                                 uuid,stock_name,
3543                                                 model_name,model_desc,
3544                                                 model_type,model_comparison_summary_df,
3545                                                 validation_data),
3546                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
3547   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
3548 }
3549 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
3550
```

Code Snippet 20 – Calling the Bagging model (Training)

This below function is responsible for generating the confusion matrix, which produces an accuracy rate of the predicted vs actual *Directions*. The validation data set is used for this purpose. The accuracy result is stored for each iteration in the in-memory table (line 1280/1286).

```

1265- run_rdmForest_model = function(fit_param, n_trees_param, uuid, stock_name, model_name, model_desc,model_type, model_comparison_summary_df,data_param){
1266
1267   set.seed(1)
1268
1269   rdmForest.fit = fit_param
1270
1271   #Generate the model prediction.
1272   rdmForest.pred = predict(rdmForest.fit, data_param, n.trees = n_trees_param, type="response")
1273
1274   #Generate the confusion matrix and calculate the classification error rate on the test/validation data
1275   rdmForest.confusion_table = table(rdmForest.pred, data_param$Direction)
1276   rdmForest.accuracy_rate = accuracy_rate_perc(rdmForest.confusion_table)
1277   rdmForest.error_rate = error_rate_perc(rdmForest.confusion_table)
1278
1279   #Add a row in the model comparison dataframe
1280   model_comparison_summary_df <- add_row_to_model_summary( model_comparison_summary_df,
1281                                                         uuid,
1282                                                         stock_name,
1283                                                         "TRUE",
1284                                                         "...",
1285                                                         model_name, model_desc,model_type,
1286                                                         rdmForest.accuracy_rate, rdmForest.error_rate)
1287   return (model_comparison_summary_df)
1288 }
1289

```

Code Snippet 21 – Calling the *run_rdmForest_model()* function

Model Optimisation

This code is similar to the training phase above. This time the number of tree(s) to grow (*n_tree_list*) is provided (line 3890). The optimisation is run over the last training period and generates the accuracy for each number of tree element in the list. The highest accuracy is retained. The minimum tree to grow is selected for this accuracy level.

```

3884 #####
3885 #####
3886 ##### Bagging Model - optimisation
3887 ##### Specialised Random Forest where mtry = p (number of attributes)
3888 #####
3889 #####
3890 n_trees_list = c(500,1000,2000,3000,4000, 5000)
3891
3892- for (n_trees in n_trees_list){
3893   model_name = "Bagging"
3894   model_desc = paste("Factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow | ntrees = ", n_trees,
3895   paste="")
3896   model_type = "OPTIMISATION"
3897   uuid = UUIDgenerate()
3898   m_try=7
3899
3900   #Get the Training and validation data for the given time window
3901   training_range = time_window_df[number_sliding_windows,"training_start_index"]:time_window_df[number_sliding_windows,"training_end_index"]
3902   training_data = final_df[training_range,]
3903   validation_range = time_window_df[number_sliding_windows,"validation_start_index"]:time_window_df[number_sliding_windows,"validation_end_index"]
3904   validation_data = final_df[validation_range,]
3905
3906   bagging_fit = randomForest( factor(Direction) ~ Close_price_1day_lag +
3907                               Volume +
3908                               Close_price_2day_lag +
3909                               Atr_atr +
3910                               Mfi +
3911                               Rsi +
3912                               Atr_trueLow,
3913                               data=training_data,
3914                               ntree=n_trees,
3915                               mtry=m_try,
3916                               importance=TRUE)
3917
3918   possibleError = tryCatch( run_rdmForest_model( bagging_fit,n_trees,
3919   uuid,stock_name,
3920   model_name,model_desc,
3921   model_type,model_comparison_summary_df,
3922   validation_data),
3923   error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
3924   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)

```

Code Snippet 22 – Calling the Tree model for the optimisation phase

Model Testing

The best training model is chosen alongside its hyper-parameter list. In this case *n_tree* (line 3934) is set to the best optimised value (i.e. 500). It is then fitted against the entire training data set, and finally tested against the test set (i.e. the last 5 business days). The confusion matrix, that evaluates the test prediction vs the expected test data, provides the test accuracy rate. The same `run_rdmForest_model()` function is called as for the training phase, line 3496. This time, the test data is used in lieu of the validation data.

```
3927 ~ # #####
3928 ~ # #####
3929 ~ # ##### Bagging Model - Testing
3930 ~ # ##### Specialised Random Forest where mtry = p (number of attributes)
3931 ~ # #####
3932 ~ # #####
3933
3934 n_trees = 500
3935 model_name = "Bagging"
3936 model_desc = paste("factor(Direction) ~ c(Close_price_1day_lag + Volume + c(Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + c(Close_price_5day_lag + Atr_trueHigh | n_trees = ", n_trees, paste=""))
3937 model_type = "TEST"
3938 uuid = UUIDgenerate()
3939 m_try=7
3940
3941 for (tw_index in time_window_seq){
3942   #Get the Training and Validation data for the given time window
3943   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"validation_end_index"]
3944   training_data = final_df[training_range,]
3945   test_range = time_window_df[tw_index,"test_start_index"]:time_window_df[tw_index,"test_end_index"]
3946   test_data = final_df[test_range,]
3947
3948   bagging_fit = randomForest( factor(Direction) ~ c(Close_price_1day_lag +
3949                                                     Volume +
3950                                                     Close_price_2day_lag +
3951                                                     Atr_atr +
3952                                                     Mfi +
3953                                                     Rsi +
3954                                                     Atr_trueLow,
3955                                                     data=training_data,
3956                                                     ntree=n_trees,
3957                                                     mtry=m_try,
3958                                                     importance=TRUE)
3959
3960   possibleError = tryCatch( run_rdmForest_model( bagging_fit,n_trees,
3961                                                 uuid,stock_name,
3962                                                 model_name,model_desc,
3963                                                 model_type,model_comparison_summary_df,
3964                                                 test_data),
3965                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
3966   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
3967 }
3968 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
3969
```

Code Snippet 20 – Calling the Bagging model (Testing)

Results

As shown in the below table, the Bagging model with the following configuration: *Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow* generates the highest training accurate rate at 45.68%. The optimisation model shows the best tree to grow level at 500. The test performance of the Bagging model, run against the selected list of attributes and the optimised min tree size, produces a 37.00% accuracy rate, with a standard deviation of 27.47%. The data is available in the file: *Model_Results_Bagging.xlsx*

K		L	M
Row Labels		Sum of model_accuracy_rate	Std Dev
TEST			
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow ntrees = 500		37.00	27.47
TRAIN			
factor(Direction) ~ Close_price_1day_lag		33.72	
factor(Direction) ~ Close_price_1day_lag + Volume		35.45	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag		33.31	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr		32.08	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi		34.99	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi		37.32	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow		45.68	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr		44.09	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag		42.40	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh		44.50	
aggregation_type		(blank)	.Y
Row Labels		Sum of model_accuracy_rate	
OPTIMISATION			
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow ntrees = 500		54.62	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow ntrees = 1000		52.31	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow ntrees = 2000		52.31	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow ntrees = 3000		52.31	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow ntrees = 4000		50.77	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow ntrees = 5000		50.77	

RandomForest

Model Description

Bagging is a special case of Random Forest where the number of variables randomly sampled as candidates at each split is equal to the number of attributes in the model (a.k.a. *mtry*). Usually Random Forest have an *mtry* set approximately to $\text{SQRT}(p)$ or $p/2$. Please refer to section *Bagging* for detailed information on the Bagging algorithm.

Model Assumptions

No specific requirements.

Further data transformation

The `randomForest.randomForest(x,...)` function implementation requires that the regressor variable is encoded as a factor (i.e. an enumerated type). Therefore, the `factor()` function has been applied on the *Direction* regressor. None of the other parameters need to be adapted.

Parameter Tuning

The tree to grow parameter is optimised.

Code Snippet Explanation

Model Training

Same comment as for the Bagging implementation, however this time there are two cases: the training case with $mtry = \text{SQRT}(p)$ and $mtry = p/2$.

```
3977 #####
3978 #####
3979 ##### Random Forest Model
3980 ##### mtry approx equal to SQRT(p)
3981 #####
3982 #####
3983 model_name = "Random Forest (mtry = SQRT(p))"
3984 model_desc = "factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag +
Atr_trueHigh"
3985 model_type = "TRAIN"
3986 uuid = UUIDgenerate()
3987 m_try=3
3988
3989- for (tw_index in time_window_seq){
3990-   #Get the training and validation data for the given time window
3991-   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"training_end_index"]
3992-   training_data = final_df[training_range,]
3993-   validation_range = time_window_df[tw_index,"validation_start_index"]:time_window_df[tw_index,"validation_end_index"]
3994-   validation_data = final_df[validation_range,]
3995-
3996-   bagging_fit = randomForest( factor(Direction) ~ Close_price_1day_lag +
3997-                               Volume +
3998-                               Close_price_2day_lag +
3999-                               Atr_atr +
4000-                               Mfi +
4001-                               Rsi +
4002-                               Atr_trueLow +
4003-                               Atr_tr +
4004-                               Close_price_5day_lag +
4005-                               Atr_trueHigh,
4006-                               data=training_data,
4007-                               mtry=m_try,
4008-                               importance=TRUE)
4009-
4010-   possibleError = trycatch( run_rdmForest_model( bagging_fit,n_trees,
4011-                                                  uuid,stock_name,
4012-                                                  model_name,model_desc,
4013-                                                  model_type,model_comparison_summary_df,
4014-                                                  validation_data),
4015-                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
4016-   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
4017- }
4018- generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
```

Code Snippet 21a – Calling the Random Forest model, where $m_try = \text{SQRT}(p)$, c.f. line 3987 (Training)

```

1450 #####
1451 #####
1452 ##### Random Forest Model
1453 ##### mtry approx equal to p/2
1454 #####
1455 #####
1456 model_name = "Random Forest (mtry = p_div_2)"
1457 model_desc = "factor(Direction) ~ Close_price_1day_lag + volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag +
Atr_trueHigh"
1458 model_type = "TRAIN"
1459 uuid = UUIDgenerate()
1460 m_try=5
1461
1462 for (tw_index in time_window_seq){
1463   #Get the Training and Validation data for the given time window
1464   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"training_end_index"]
1465   training_data = final_df[training_range,]
1466   validation_range = time_window_df[tw_index,"validation_start_index"]:time_window_df[tw_index,"validation_end_index"]
1467   validation_data = final_df[validation_range,]
1468
1469   bagging_fit = randomForest( factor(Direction) ~ Close_price_1day_lag +
1470                               volume +
1471                               Close_price_2day_lag +
1472                               Atr_atr +
1473                               Mfi +
1474                               Rsi +
1475                               Atr_trueLow +
1476                               Atr_tr +
1477                               Close_price_5day_lag +
1478                               Atr_trueHigh,
1479                               data=training_data,
1480                               mtry=m_try,
1481                               importance=TRUE)
1482
1483   possibleError = tryCatch( run_rdmForest_model( bagging_fit,n_trees,
1484                                                  uuid,stock_name,
1485                                                  model_name,model_desc,
1486                                                  model_type,model_comparison_summary_df,
1487                                                  validation_data),
1488                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
1489   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
1490 }
1491 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)

```

Code Snippet 22b – Calling the Random Forest model, where $m_try = p/2$, c.f. line 1460 (Training)

Model Optimisation

Same comments as for the bagging case. The n_tree optimisation is performed for both $m_try = p/2$ and $m_try = \text{SRQT}(p)$.

```

352 #####
353 #####
354 ##### Random Forest Model - Optimisation
355 ##### mtry approx equal to SQRT(p)
356 #####
357 #####
358 n_trees_list = c(500,1000,2000,3000,4000, 5000)
359
360 for (n_trees in n_trees_list){
361   model_name = "Random Forest (mtry = SQRT(p))"
362   model_desc = paste("factor(Direction) ~ Close_price_1day_lag + volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr +
Close_price_5day_lag + Atr_trueHigh | n_trees = ", n_trees, paste="")
363   model_type = "OPTIMISATION"
364   uuid = UUIDgenerate()
365   m_try=3
366
367   #Get the Training and Validation data for the given time window
368   training_range = time_window_df[number_sliding_windows,"training_start_index"]:time_window_df[number_sliding_windows,"training_end_index"]
369   training_data = final_df[training_range,]
370   validation_range = time_window_df[number_sliding_windows,"validation_start_index"]:time_window_df[number_sliding_windows,"validation_end_index"]
371   validation_data = final_df[validation_range,]
372
373   bagging_fit = randomForest( factor(Direction) ~ Close_price_1day_lag +
374                               volume +
375                               Close_price_2day_lag +
376                               Atr_atr +
377                               Mfi +
378                               Rsi +
379                               Atr_trueLow +
380                               Atr_tr +
381                               Close_price_5day_lag +
382                               Atr_trueHigh,
383                               data=training_data,
384                               ntree=n_trees,
385                               mtry=m_try,
386                               importance=TRUE)
387
388   possibleError = tryCatch( run_rdmForest_model( bagging_fit,n_trees,
389                                                  uuid,stock_name,
390                                                  model_name,model_desc,
391                                                  model_type,model_comparison_summary_df,
392                                                  validation_data),
393                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
394   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)

```

Code Snippet 23b – Calling the Random Forest model, where $m_try = \text{SQRT}(p)$, c.f. line 365 (Optimisation)

```

827 ##### Random Forest Model - Optimisation
828 ##### mtry approx equal to p_div_2
829 #####
830 #####
831 n_trees_list = c(500,1000,2000,3000,4000, 5000)
832
833 for (n_trees in n_trees_list){
834   model_name = "Random Forest (mtry = p_div_2)"
835   model_desc = paste("factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr +
Close_price_5day_lag + Atr_trueHigh | ntrees = ", n_trees, paste="")
836   model_type = "OPTIMISATION"
837   uuid = UUIDgenerate()
838   m_try=4
839
840   #Get the Training and Validation data for the given time window
841   training_range = time_window_df[number_sliding_windows,"training_start_index"]:time_window_df[number_sliding_windows,"training_end_index"]
842   training_data = final_df[training_range,]
843   validation_range = time_window_df[number_sliding_windows,"validation_start_index"]:time_window_df[number_sliding_windows,"validation_end_index"]
844   validation_data = final_df[validation_range,]
845
846   bagging_fit = randomForest( factor(Direction) ~ Close_price_1day_lag +
847                                     Volume +
848                                     Close_price_2day_lag +
849                                     Atr_atr +
850                                     Mfi +
851                                     Rsi +
852                                     Atr_trueLow +
853                                     Atr_tr +
854                                     Close_price_5day_lag +
855                                     Atr_trueHigh,
856                                     data=training_data,
857                                     ntree=n_trees,
858                                     mtry=m_try,
859                                     importance=TRUE)
860
861   possibleError = tryCatch( run_rdmForest_model( bagging_fit,n_trees,
862                                     uuid,stock_name,
863                                     model_name,model_desc,
864                                     model_type,model_comparison_summary_df,
865                                     validation_data),
866                                     error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
867   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
868

```

Code Snippet 23b – Calling the Random Forest model, where $m_try = p/2$, c.f. line 838 (Optimisation)

Model Testing

The best training_fit model is chosen and is fitted against the entire training data set. It is then tested against the test set (i.e. the last 5 business days).

- For the $m_try = \text{SQRT}(p)$ case, the n_tree is set to the best optimised value (1000).
- For the $m_try = p/2$ case, the n_tree is set to the best optimised value (4000).

The confusion matrix, that evaluates the test prediction vs the expected test data, provides the test accuracy rate. The same `run_rdmForest_model()` function is called as for the training phase, line 3496. This time, the test data is used in lieu of the validation data.

```

4399 ##### Random Forest Model - Testing
4400 ##### mtry approx equal to SQRT(p)
4401 #####
4402 #####
4403 #####
4404 n_trees = 1000
4405 model_name = "Random Forest (mtry = SQRT(p))"
4406 model_desc = paste("factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag
+ Atr_trueHigh | ntrees = ", n_trees, paste="")
4407 model_type = "TEST"
4408 uuid = UUIDgenerate()
4409 m_try=3
4410
4411 for (tw_index in time_window_seq){
4412   #Get the Training and Validation data for the given time window
4413   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"validation_end_index"]
4414   training_data = final_df[training_range,]
4415   test_range = time_window_df[tw_index,"test_start_index"]:time_window_df[tw_index,"test_end_index"]
4416   test_data = final_df[test_range,]
4417
4418   bagging_fit = randomForest( factor(Direction) ~ Close_price_1day_lag +
4419                                     Volume +
4420                                     Close_price_2day_lag +
4421                                     Atr_atr +
4422                                     Mfi +
4423                                     Rsi +
4424                                     Atr_trueLow +
4425                                     Atr_tr +
4426                                     Close_price_5day_lag +
4427                                     Atr_trueHigh,
4428                                     data=training_data,
4429                                     ntree=n_trees,
4430                                     mtry=m_try,
4431                                     importance=TRUE)
4432
4433   possibleError = tryCatch( run_rdmForest_model( bagging_fit,n_trees,
4434                                     uuid,stock_name,
4435                                     model_name,model_desc,
4436                                     model_type,model_comparison_summary_df,
4437                                     test_data),
4438                                     error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
4439   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
4440 }
4441 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
4442

```

Code Snippet 23a – Calling the Random Forest model, where $m_try = \text{SQRT}(p)$, c.f. line 4404 (Testing)

```

4870 #####
4871 #####
4872 #####
4873 ##### Random Forest Model - Testing
4874 ##### mtry approx equal to p/2
4875 #####
4876 #####
4877 n_trees = 4000
4878 model_name = "Random Forest (mtry = p_div_2)"
4879 model_desc = paste("factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh | ntrees = ", n_trees, paste="")
4880 model_type = "TEST"
4881 uuid = UUIDgenerate()
4882 m_try=4
4883
4884 for (tw_index in time_window_seq){
4885   #Get the Training and Validation data for the given time window
4886   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"validation_end_index"]
4887   training_data = final_df[training_range,]
4888   test_range = time_window_df[tw_index,"test_start_index"]:time_window_df[tw_index,"test_end_index"]
4889   test_data = final_df[test_range,]
4890
4891   bagging_fit = randomForest( factor(Direction) ~ Close_price_1day_lag +
4892                               Volume +
4893                               Close_price_2day_lag +
4894                               Atr_atr +
4895                               Mfi +
4896                               Rsi +
4897                               |
4898                               Atr_trueLow +
4899                               Atr_tr +
4900                               Close_price_5day_lag +
4901                               Atr_trueHigh,
4902                               data=training_data,
4903                               ntree=n_trees,
4904                               mtry=m_try,
4905                               importance=TRUE)
4906
4907   possibleError = tryCatch( run_rdmForest_model( bagging_fit,n_trees,
4908                                                 uuid,stock_name,
4909                                                 model_name,model_desc,
4910                                                 model_type,model_comparison_summary_df,
4911                                                 test_data),
4912                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
4913   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
4914 }

```

Code Snippet 23b – Calling the Random Forest model, where $m_try = p/2$, c.f. line 4877 (Testing)

Results

The $m_try = \text{SQRT}(p)$ Case

As shown in the below table, the Random Forest model with the following configuration: $\text{factor}(\text{Direction}) \sim \text{Close_price_1day_lag} + \text{Volume} + \text{Close_price_2day_lag} + \text{Atr_atr} + \text{Mfi} + \text{Rsi} + \text{Atr_trueLow} + \text{Atr_tr} + \text{Close_price_5day_lag} + \text{Atr_trueHigh}$ generates the highest training accurate rate at 45.08%. The optimisation model shows the best tree to grow level is at 1000. The test performance of the Tree model run against the selected list of attributes and the optimised min tree size produces a 34.80% accuracy rate, with a standard deviation of 26.30%. The data is available in the file: *Model_Results_Random Forest (mtry = SQRT(p)).xlsx*

aggregation_type	AVG	Std Dev
Row Labels	Sum of model_accuracy_rate	Std Dev
TEST		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 1000	34.80	26.30
TRAIN		
factor(Direction) ~ Close_price_1day_lag	33.72	
factor(Direction) ~ Close_price_1day_lag + Volume	35.14	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag	33.18	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr	31.94	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi	34.87	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi	37.22	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow	44.12	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr	42.28	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag	41.69	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh	45.08	
aggregation_type	(blank)	
Row Labels	Sum of model_accuracy_rate	
OPTIMISATION		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 500	46.15	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 1000	50.00	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 2000	49.23	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 3000	48.46	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 4000	48.46	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 5000	48.46	

The m try = SQRT(p) Case

As shown in the below table, the Random Forest model with the following configuration: *factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh* generates the highest training accurate rate at 45.08%. The optimisation model shows the best tree to grow level is at 4000. The test performance of the Tree model run against the selected list of attributes and the optimised min tree size produces a 35.00% accuracy rate, with a standard deviation of 25.33%. The data is available in the file: *Model_Results_Random Forest (mtry = p_div_2).xlsx*

K		L	M
Row Labels		Sum of model_accuracy_rate	Std Dev
TEST			
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 4000		35.00	25.33
TRAIN			
factor(Direction) ~ Close_price_1day_lag		33.72	
factor(Direction) ~ Close_price_1day_lag + Volume		35.14	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag		33.39	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr		31.94	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi		34.87	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi		37.15	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow		44.12	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr		42.72	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag		41.93	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh		44.70	
aggregation_type		(blank)	
Row Labels		Sum of model_accuracy_rate	
OPTIMISATION			
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 1000		49.23	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 2000		50.77	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 3000		51.54	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 4000		52.31	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 500		49.23	
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh ntrees = 5000		50.00	

Boosting

Model Description

Boosting is similar to the Bagging approach, except that the trees are grown sequentially: each tree is grown using information from previously grown trees [14].

Model Assumptions

No specific requirements.

Further data transformation

No data transformation is required.

Parameter Tuning

The tree to fit and the shrinkage parameters are optimised.

Code Snippet Explanation

Model Training

The *for-loop*, line 2683 iterates through the training data time windows (100 iterations). At each iteration, the *gbm* (x, \dots) function is called with a list of explanatory and explained variables, for a training data slice (line 2692). It produces the *boost_fit* object, i.e. the function fitting the model. The *boost_fit* model is then passed to the *run_boosting_model(...)* function, alongside a few parameters. One of them is the validation set for the time slide period. On successful run, the model results are stored into an object named *model_comparison_summary_df*. It is tagged with the state *model_run_success = TRUE*. The training accuracy rate is computed at each iteration. The average of the training accuracy is calculated on line 2712. In case of a computational failure, the model description is added to the same in-memory data frame and tagged with the state *model_run_success = FALSE* (line 2710). The below code snippet only shows one example of a model trained against a given set of attributes. This example is repeated for each list of attributes. An average training accuracy rate is generated for each instance. The complete list of training case is available in the source code.

```
2671 #####
2672 #####
2673 ##### Boosting Model
2674 #####
2675 #####
2676
2677 model_name = "Boosting"
2678 model_desc = "Direction ~ Close_price_1day_lag + volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag +
Atr_trueHigh"
2679 model_type = "TRAIN"
2680 uuid = UUIDgenerate()
2681 n_trees = 500
2682
2683 - for (tw_index in time_window_seq){
2684   #Get the Training and Validation data for the given time window
2685   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"training_end_index"]
2686   training_data = final_df[training_range,]
2687   validation_range = time_window_df[tw_index,"validation_start_index"]:time_window_df[tw_index,"validation_end_index"]
2688   validation_data = final_df[validation_range,]
2689
2690   #This is a multinomial classification problem. Therefore, the distribution is set to "multinomial".
2691   #The gbm() function fit the Direction variable against all other selected predictors on the training set.
2692   boost_fit = gbm(Direction ~ Close_price_1day_lag +
2693                   volume +
2694                   Close_price_2day_lag +
2695                   Atr_atr +
2696                   Mfi +
2697                   Rsi +
2698                   Atr_trueLow +
2699                   Atr_tr +
2700                   Close_price_5day_lag +
2701                   Atr_trueHigh,
2702                   data = training_data, distribution = "multinomial", n.trees = n_trees)
2703
2704   possibleError = tryCatch( run_boosting_model( boost_fit,uuid,
2705                                               n_trees,stock_name,
2706                                               model_name,model_desc,
2707                                               model_type,model_comparison_summary_df,
2708                                               validation_data),
2709                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
2710   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
2711 }
2712 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
2713
```

Code Snippet 24 – Calling the Boosting model (Training)

This below function is responsible for generating the confusion matrix, which produces an accuracy rate of the predicted vs actual *Directions*. The validation data set is used for this purpose. The accuracy result is stored for each iteration in the in-memory table (line 1253/1259).

```

1229 run_boosting_model = function(fit_param, uuid, n_trees_param, stock_name, model_name, model_desc, model_type, model_comparison_summary_df, test_data_param){
1230
1231   set.seed(1)
1232
1233   boost_fit = fit_param
1234
1235   #Generate the model prediction. Type=response returns the class prediction as probability
1236   boost_pred = predict(boost_fit, test_data_param, n.trees = n_trees_param, type="response")
1237
1238   #We know need to find the most probable class (the class showing the highest probability) for each row of the boost_pred
1239   #see http://stackoverflow.com/questions/29454883/in-gbm-multinomial-dist-how-to-use-predict-to-get-categorical-output
1240   boost_proba = apply(boost_pred, 1, which.max)
1241   #And now turn the class id (1,2 or 3) into the class name
1242   boost_pred_class = colnames(boost_pred)[boost_proba]
1243   #necessary as a non-compatible shrinkage could generate unpredictable classes
1244   boost_accuracy_rate = "NA"
1245   boost_error_rate = "NA"
1246   if (length(boost_pred_class) > 0){
1247     boost_confusion_table = table(boost_pred_class, test_data_param$direction)
1248     boost_accuracy_rate = accuracy_rate_perc(boost_confusion_table)
1249     boost_error_rate = error_rate_perc(boost_confusion_table)
1250   }
1251
1252   #Add a row in the model comparison dataframe
1253   model_comparison_summary_df <- add_row_to_model_summary( model_comparison_summary_df,
1254                                                         uuid,
1255                                                         stock_name,
1256                                                         "TRUE",
1257                                                         "",
1258                                                         model_name, model_desc, model_type,
1259                                                         boost_accuracy_rate, boost_error_rate)
1260
1261
1262   return (model_comparison_summary_df)
1263 }
1264

```

Code Snippet 25 – Calling the *run_boosting_model()* function

Model Optimisation

This code is similar to the training phase above. The number of tree(s) to fit (*tree_list*) and their shrinkage (*lambdas*) parameters are provided (line 3056/3058). The optimisation of these two parameters is run over the last training period and generate the accuracy. The highest accuracy is retained. The minimum tree/shrinkage level are selected for this accuracy level.

```

3044 #####
3045 ##### Boosting optimisation
3046 #####
3047 #####
3048 #####
3049 #####
3050 #The optimisation is performed on the last sliding window
3051 training_range = time_window_df[number_sliding_windows,"training_start_index"]:time_window_df[number_sliding_windows,"training_end_index"]
3052 training_data = final_df[training_range,]
3053 validation_range = time_window_df[number_sliding_windows,"validation_start_index"]:time_window_df[number_sliding_windows,"validation_end_index"]
3054 validation_data = final_df[validation_range,]
3055
3056 tree_list = c (500,1000,1500,2000,3000,5000)
3057 #Generate a vector of sequence of starting from 0.1 and going to 1.0, with a 0.1 step
3058 lambdas = seq(0.1, 1, by = 0.1)
3059 len = length(lambdas)
3060 for (the_tree in tree_list){
3061   for (i in 1:len){
3062     model_name = "Boosting"
3063     model_desc = paste("Direction ~ Close_price_1day_lag + volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr | n.trees =", the_tree, " | shrinkage =", lambdas[i], sep = " ")
3064     model_type = "OPTIMISATION"
3065     uuid = UUIDgenerate()
3066
3067     boost_fit = gbm(Direction ~ Close_price_1day_lag +
3068                    volume +
3069                    Close_price_2day_lag +
3070                    Atr_atr +
3071                    Mfi +
3072                    Rsi +
3073                    Atr_trueLow +
3074                    Atr_tr
3075                    ,data = training_data, distribution = "multinomial", n.trees = the_tree, shrinkage = lambdas[i])
3076
3077     possibleError = tryCatch( run_boosting_model( boost_fit,uuid,
3078                                                  n_trees=the_tree,stock_name,
3079                                                  model_name,model_desc,
3080                                                  model_type,model_comparison_summary_df,
3081                                                  validation_data),
3082                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
3083     add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
3084   }
3085 }
3086

```

Code Snippet 26 – Calling the Boosting model for the optimisation phase

Model Testing

The best training model is chosen alongside its hyper-parameter list. In this case the n_tree (line 1092) is set to the best optimised value (1000), with its shrinkage level at 0.1 (line1093). It is then fitted against the entire training data set, and finally tested against the test set (i.e. the last 5 business days). The confusion matrix, that evaluates the test prediction vs the expected test data, provides the test accuracy rate. The same `run_mlr_model()` function is called as for the training phase, line 1678. This time, the test data is used in lieu of the validation data.

```
3087 #####
3088 #####
3089 ##### Boosting Model - Testing
3090 #####
3091 #####
3092 n_trees = 500
3093 the_lambda = 0.1
3094 model_name = "Boosting"
3095 model_desc = paste("Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr | n_trees =", the_tree, " | shrinkage =", the_lambda, sep = " ")
3096
3097 model_type = "TEST"
3098 uuid = UUIDgenerate()
3099
3100 for (tw_index in time_window_seq){
3101   #Get the Training and Validation data for the given time window
3102   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"validation_end_index"]
3103   training_data = final_df[training_range,]
3104   test_range = time_window_df[tw_index,"test_start_index"]:time_window_df[tw_index,"test_end_index"]
3105   test_data = final_df[test_range,]
3106
3107   boost_fit = gbm(Direction ~ Close_price_1day_lag +
3108     Volume +
3109     Close_price_2day_lag +
3110     Atr_atr +
3111     Mfi +
3112     Rsi +
3113     Atr_trueLow +
3114     Atr_tr
3115     ,data = training_data, distribution = "multinomial", n.trees = n_trees, shrinkage = the_lambda)
3116
3117   possibleError = tryCatch( run_boosting_model( boost_fit,uuid,
3118     n_trees,stock_name,
3119     model_name,model_desc,
3120     model_type,model_comparison_summary_df,
3121     test_data),
3122     error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
3123   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
3124 }
3125 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
3126
3127 #####
```

Code Snippet 27 – Calling the Boosting model (Testing)

Results

As shown in the below table, the Boosting model with the following configuration: *Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi+Atr_trueLow + Atr_tr* generates the highest training accurate rate at 35.53%. The optimisation model shows the number of trees and shrinkage should be set to respectively 1000 and 0.1. The test performance of the Boosting model, run against the selected list of attributes and the optimised parameters, produces a 47.60% accuracy rate, with a standard deviation of 25.71%. The data is available in the file: *Model_Results_Boosting.xlsx*

Row Labels	K	L	M	
TEST	Sum of model_accuracy_rate	Std Dev		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 1		47.60	25.71	
TRAIN				
Direction ~ Close_price_1day_lag + Volume		31.00		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag		30.94		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr		32.76		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi		32.08		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi		31.94		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow		32.37		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr		35.53		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag		33.83		
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh		34.02		
factor(Direction) ~ Close_price_1day_lag^3 + Volume^2 + Close_price_1day_lag + Volume + (Atr_atr * Mfi) + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr		31.42		
aggregation_type	(blank)			
Row Labels	Sum of model_accuracy_rate			
OPTIMISATION				
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 0.1	54.62			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 0.2	52.31			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 0.3	50.00			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 0.4	50.77			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 0.5	53.85			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 0.6	50.00			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 0.7	0.00			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 0.8	0.00			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 0.9	0.00			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1000 shrinkage = 1	0.00			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1500 shrinkage = 0.1	54.62			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1500 shrinkage = 0.2	52.31			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1500 shrinkage = 0.3	50.00			
Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr n_trees = 1500 shrinkage = 0.4	50.77			

Support Vector Machine (SVM)

Model Description

The SVM is a generalisation of the classifier method named the maximal margin classifier [14]. The SVM finds a plane that separates the classes in feature space. It can accommodate non-linear class boundaries and multinomial classes.

Model Assumptions

No specific requirements.

Further data transformation

The R `e1071.svm(x,...)` function implementation requires that the regressor variable is encoded as a factor (i.e. an enumerated type). Therefore, the `factor()` function has been applied on the *Direction* regressor. None of the other parameters need to be adapted.

Parameter Tuning

There are several parameters that are optimised: the kernel type (used for linear or nonlinear classification learning), the gamma (a parameter used for kernels) and the cost (i.e. 'C'-constant of the regularization term in the Lagrange formulation).

Code Snippet Explanation

Model Training

The *for-loop*, line 933 iterates through the training data time windows (100 iterations). At each iteration, the `svm(x,...)` function is called with a list of explanatory and explained variables, for a training data slice (line 940). It produces the `svm_fit` object, i.e. the function fitting the model. The `svm_fit` model is then passed to the `run_svm_model(...)` function, alongside a few parameters. One of them is the validation set for the time slide period. On successful run, the model results are stored into an object named `model_comparison_summary_df`. It is tagged with the state `model_run_success = TRUE`. The training accuracy rate is computed at each iteration. The average of the training accuracy is calculated on line 957. In case of a computational failure, the model description is added to the same in-memory data frame and tagged with the state `model_run_success = FALSE` (line 953). The below code snippet only shows one example of a model trained against a given set of attributes. This example is repeated for each list of attributes. An average training accuracy rate is generated for each instance. The complete list of training case is available in the source code.

```
922
923 #####
924 #####
925 ##### Support Vector Machine (SVM) - Model Training
926 #####
927 #####
928 model_name = "Support Vector Machine (SVM)"
929 model_desc = "factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag +
Atr_trueHigh"
930 model_type = "TRAIN"
931 uuid = UUIDgenerate()
932
933 + for (tw_index in time_window_seq){
934     #Get the Training and validation data for the given time window
935     training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"training_end_index"]
936     training_data = final_df[training_range,]
937     validation_range = time_window_df[tw_index,"validation_start_index"]:time_window_df[tw_index,"validation_end_index"]
938     validation_data = final_df[validation_range,]
939
940     svm_fit = svm (factor(Direction) ~ Close_price_1day_lag +
941                   Volume +
942                   Close_price_2day_lag +
943                   Atr_atr +
944                   Mfi +
945                   Rsi +
946                   Atr_trueLow +
947                   Atr_tr +
948                   Close_price_5day_lag +
949                   Atr_trueHigh,
950                   probability = TRUE, #indicates the model should allow for probability predictions
951                   data=training_data)
952
953     possibleError = trycatch( run_svm_model(svm_fit,uuid,stock_name,model_name,model_desc,model_type,model_comparison_summary_df,validation_data),
954                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
955     add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
956 }
957 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
958
```

Code Snippet 28 – Calling the SVM model (Training)

The below code generates the confusion matrix and accuracy rate for the SVM model. Line 1162/1164 deals with predicting each class index (1,2,3) probabilities. This is necessary as the SVM is used in a multinomial mode. The class index with the highest probability is selected (line 1666/1176). Lines 1177/1183 deal with converting the class index into class names (*Up/Down/Neutral*). Then the confusion matrix is generated, which produces an accuracy rate of the predicted vs actual *Directions*. The validation data set is used for this purpose. The accuracy result is stored for each iteration in the in-memory table (Line 1192/1198).

```

1154
1155 #the test_data_param could be a training/validation or test dataset
1156 run_svm_model = function(svm_param, uuid, stock_name, model_name, model_desc, model_type, model_comparison_summary_df, test_data_param){
1157
1158     svm.fit = svm_param
1159
1160     #Generate the confusion matrix and calculate the training/validation error rate
1161     #decision.values = TRUE, probability = TRUE are necessary when multi class probabilities are required
1162     pred_prob = predict(svm.fit, test_data_param, decision.values = TRUE, probability = TRUE)
1163     #Get the class probabilities
1164     pred_prob_attr = attr(pred_prob, "probabilities")
1165     #Get the most probable class per row
1166     pred_class_as_factor = apply(pred_prob_attr, 1, which.max)
1167     pred_class_as_name = pred_class_as_factor
1168     #Get the name of each class
1169     col_name = colnames(attr(pred_prob, "probabilities"))
1170     index = 1
1171     map = NULL
1172     #Create a list that map an index to a class name
1173     for (cn in col_name){
1174         map = append(map, list(id = index, name = cn))
1175         index = index + 1
1176     }
1177     #Find the predicted class name, given the predicted class factor, from the map list
1178     index = 1
1179     for(class_as_factor in pred_class_as_factor){
1180         pos = match(class_as_factor, map)
1181         pred_class_as_name[index] = map[pos+1]
1182         index = index + 1
1183     }
1184     #Generate a dataframe for the predicted class names
1185     svm.pred_class_as_name_df = data.frame(matrix(unlist(pred_class_as_name)), stringsAsFactors=FALSE)
1186     #Generate the confusion matrix
1187     svm.confusion_table = table(svm.pred_class_as_name_df$matrix.unlist.pred_class_as_name., test_data_param$Direction)
1188     svm.accuracy_rate = accuracy_rate_perc(svm.confusion_table)
1189     svm.error_rate = error_rate_perc(svm.confusion_table)
1190
1191     #Add a row in the model comparison dataframe
1192     model_comparison_summary_df <- add_row_to_model_summary( model_comparison_summary_df,
1193                                                             uuid,
1194                                                             stock_name,
1195                                                             "TRUE",
1196                                                             "",
1197                                                             model_name, model_desc, model_type,
1198                                                             svm.accuracy_rate, svm.error_rate)
1199
1200     return (model_comparison_summary_df)
1201 }

```

Code Snippet 29 – Calling the `run_svm_model()` function

Model Optimisation

This code is similar to the training phase above. However, there are three parameters to optimise, namely the cost, gamma and kernel type (c.f. line 5237/8239). The optimisation of these three parameters, via 3 inner loops (line 5241/5239), is run over the last training period and generate the accuracy. The hyperparameter list with the highest accuracy is retained.

```
5230 #####
5231 #####
5232 ##### Support Vector Machine (SVM) - Model Optimisation (cost/gamma and kernel)
5233 #####
5234 #####
5235
5236 #The model above model will be run each cost in the below cost list:
5237 cost_list = c(0.001 , 0.01 , 0.1 , 1 , 100 )
5238 gamma_list = c(0.01,0.03,0.05,0.5,1,100)
5239 | kernel_list = c("linear", "radial", "sigmoid", "polynomial")
5240
5241 + for (the_kernel in kernel_list){
5242 +   for (the_cost in cost_list){
5243 +     for (the_gamma in gamma_list){
5244 +       model_name = "Support Vector Machine (SVM)"
5245 +       model_desc = paste( "factor(Direction) ~ Close_price_1day_lag + volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh +
Atr_trueLow | kernel = ", toString(the_kernel), " | cost = ", toString(the_cost), " | gamma = ", toString(the_gamma), sep="")
5246 +       model_type = "OPTIMISATION"
5247 +       uuid = UUIDgenerate()
5248
5249 +       #The optimisation is performed on the last sliding window
5250 +       training_range = time_window_df[number_sliding_windows,"training_start_index"]:time_window_df[number_sliding_windows,"training_end_index"]
5251 +       training_data = final_df[training_range,]
5252 +       validation_range = time_window_df[number_sliding_windows,"validation_start_index"]:time_window_df[number_sliding_windows,"validation_end_index"]
5253 +       validation_data = final_df[validation_range,]
5254
5255 +       svm_fit = svm (factor(Direction) ~ Close_price_1day_lag +
5256 +                     Volume +
5257 +                     Atr_tr +
5258 +                     Close_price_2day_lag +
5259 +                     Rsi +
5260 +                     Mfi +
5261 +                     Close_price_5day_lag +
5262 +                     Atr_trueHigh +
5263 +                     Atr_trueLow,
5264 +                     kernel = the_kernel,
5265 +                     cost = the_cost,
5266 +                     gamma = the_gamma,
5267 +                     probability = TRUE, #indicates the model should allow for probability predictions
5268 +                     data=training_data)
5269
5270 +       possibleError = tryCatch( run_svm_model(svm_fit,uuid,stock_name,model_name,model_desc,model_type,model_comparison_summary_df,validation_data),
5271 +                                error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
5272 +       add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
5273 +     }
5274 +   }
5275 }
```

Code Snippet 30 – Calling the SVM model for the optimisation phase

Model Testing

The best training model is chosen alongside the optimised parameters (Kernel = 'linear', gamma = 0.01, cost = 1). This model is fitted against the entire training data set. It is then tested against the test set (i.e. the last 5 business days). The confusion matrix, that evaluates the test prediction vs the expected test data, provides the test accuracy rate. The same `run_svm_model()` function is called as for the training phase, line 5313. This time, the test data is used in lieu of the validation data.

```
5277 #####
5278 #####
5279 ##### Support Vector Machine (SVM) - Model Testing
5280 #####
5281 #####
5282 the_kernel = "linear"
5283 the_cost = 1
5284 the_gamma = 0.01
5285 model_name = "Support Vector Machine (SVM)"
5286 model_desc = paste("factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh | kernel = ", toString(the_kernel), " | cost = ", toString(the_cost), " | gamma = ", toString(the_gamma), sep="")
5287 model_type = "TEST"
5288 uuid = uuidgenerate()
5289
5290 for (tw_index in time_window_seq){
5291   #Get the Training and Validation data for the given time window
5292   training_range = time_window_df[tw_index,"training_start_index"]:time_window_df[tw_index,"validation_end_index"]
5293   training_data = final_df[training_range,]
5294   test_range = time_window_df[tw_index,"test_start_index"]:time_window_df[tw_index,"test_end_index"]
5295   test_data = final_df[test_range,]
5296
5297   svm_fit = svm (factor(Direction) ~ Close_price_1day_lag +
5298                 Volume +
5299                 Close_price_2day_lag +|
5300                 Atr_atr +
5301                 Mfi +
5302                 Rsi +
5303                 Atr_trueLow +
5304                 Atr_tr +
5305                 Close_price_5day_lag +
5306                 Atr_trueHigh,
5307                 kernel = the_kernel,
5308                 cost = the_cost,
5309                 gamma = the_gamma,
5310                 probability = TRUE, #indicates the model should allow for probability predictions
5311                 data=training_data)
5312
5313   possibleError = tryCatch( run_svm_model(svm_fit,uuid,stock_name,model_name,model_desc,model_type,model_comparison_summary_df,test_data),
5314                             error = function(e) print(paste("MODEL ERROR: ", e, sep="")))
5315   add_failed_model(model_comparison_summary_df,uuid, stock_name, model_name,model_desc,model_type, possibleError)
5316 }
5317 generate_avg_model(model_comparison_summary_df,uuid, time_window_seq,model_type)
5318
5319 cat ("")
```

Code Snippet 30 – Calling the SVM model (Testing)

Results

As shown in the below table, the SVM model with the following configuration: *Direction ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh* generates the highest training accurate rate at 49.85%. The optimisation model shows the kernel, cost and gamma to be respectively set to 'linear', 1 and 0.01. The test performance of the SVM model run against the selected list of attributes and the optimised parameters produces a 68.40% accuracy rate, with a standard deviation of 26.39%. The data is available in the file: *Model_Results_Boosting.xlsx*

aggregation_type	AVG	Std Dev
Row Labels	Sum of model_accuracy	
TEST		
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh kernel = linear cost = 1 gamma = 0.01	68.40	26.39
TRAIN		
factor(Direction) ~ Close_price_1day_lag		28.74
factor(Direction) ~ Close_price_1day_lag + Volume		28.43
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag		28.23
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr		23.08
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi		32.42
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi		32.45
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow		47.38
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr		48.44
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag + Atr_trueHigh		49.85
factor(Direction) ~ Close_price_1day_lag + Volume + Close_price_2day_lag + Atr_atr + Mfi + Rsi + Atr_trueLow + Atr_tr + Close_price_5day_lag		47.35
aggregation_type	(blank)	
Row Labels	Sum of model_accuracy_rate	
OPTIMISATION		
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.001 gamma = 0.01		60.77
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.001 gamma = 0.03		63.85
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.001 gamma = 0.05		60.77
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.001 gamma = 0.5		60.00
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.001 gamma = 1		63.08
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.001 gamma = 100		66.15
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.01 gamma = 0.01		62.31
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.01 gamma = 0.03		64.62
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.01 gamma = 0.05		65.38
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.01 gamma = 0.5		60.77
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.01 gamma = 1		66.15
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.01 gamma = 100		60.77
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.1 gamma = 0.01		74.62
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.1 gamma = 0.03		73.95
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.1 gamma = 0.05		73.95
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.1 gamma = 0.5		73.95
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.1 gamma = 1		73.95
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 0.1 gamma = 100		73.08
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 1 gamma = 0.01		80.00
factor(Direction) ~ Close_price_1day_lag + Volume + Atr_tr + Close_price_2day_lag + Rsi + Mfi + Close_price_5day_lag + Atr_trueHigh + Atr_trueLow kernel = linear cost = 1 gamma = 0.03		80.00

Evaluation

The below table summarises the test performance obtained for each model, in decreasing order of average test performance accuracy. Initial experiments on a smaller time slice window of approximately 100 days (vs 260days here), and without class rebalancing, showed the test accuracy for LDA and QDA at approximately 50%. Bagging was approximately at 63%. With the current settings, the situation is reversed. The three best performing algorithms are the QDA, Ridge and LDA with approximately an 80% test accuracy rate and a standard deviation at approximately 1.4%. The Lasso model follows in the tail of the first three ones, with a test performance accuracy of 75.39% and a low standard deviation at 1.44%. All the other model show test performance accuracy rates below 70%, with high variance, around 25% standard deviation.

The surprising fact is that most of attributes do not follow a normal distribution (c.f. *Normal Distribution Test* section), but at the time, the models that require attributes normal distribution, i.e. QDA, Ridge, LDA and Lasso, are showing the best average test performance accuracies and stabilities.

Model Name	Average Test Performance
QDA	Avg: 80.90% - Std Dev:1.64%
Ridge	Avg: 80.71% - Std Dev:1.38%
LDA	Avg: 80.20% - Std Dev:1.06%
Lasso	Avg: 75.39% - Std Dev:1.44%
SVM	Avg: 68.40% - Std Dev: 26.39%
Boosting	Avg: 47.60% - Std Dev: 25.71%
Decision Tree	Avg: 40.80% - Std Dev:21.45%
Bagging (Random Forest where mtry= p)	Avg: 37.00% - Std Dev:27.47%
Random Forest mtry= p/2 mtry= SQRT(p)	Avg: 35.00% - Std Dev: 25.33% Avg: 34.80% - Std Dev: 26.30%

The next step involves looking at some of test confusion matrix measures for the model with the best test performance results, i.e. QDA. The aim is to establish whether other patterns can be uncovered. The description of the measures and formulas are provided below. The full implementation is available in the file named: *qda_confusion_matrix.xlsx*. The sensitivity levels showed in the below QDA table, indicates that the model is better at discovering *Neutral* trends (sensitivity =86.63%), and then it scores higher in discovering *Up* trends (sensitivity=76.45%) than *Down* trends (sensitivity=68.91%). The precision and F1 measure shows the same pattern for *Neutral* trends, but this time, precision for *Down* trend scores higher than for *Up* trends.

Total Population	Prediction Positive	Prediction Negative
Expected Positive	True Positive (TP)	False Negative (FN)
Expected Negative	False Positive (FP)	True Negative (TN)

Measure	Formula	Description
Sensitivity (a.k.a. recall)	$TP/(TP+FN)$	It measures the proportion of positives that are correctly identified (i.e. how good is a test at detecting the positives).
Specificity	$TN/(TN+FP)$	It measures the proportion of negatives that are correctly identified (i.e. how good is a test at detecting false alarms).
Precision	$TP/(TP+FP)$	It measures the proportion of positives that were relevant.
F1-Score	$2*TP/(2*TP+FP+ FN)$	It is a measure of the test accuracy. F1-Score is always between 0 (worst) and 1(best).
Accuracy	$(TP+TN)/(TP+FN+FP+TN)$	It is a measure of the statistical bias. Accuracy is between 0 (maximum bias) and 1 (no bias).

QDA Measures / Scenarios	Down Vs (Neutral + Up)	Neutral Vs (Down + Up)	Up Vs (Down + Neutral)
Sensitivity/Recall (i.e. True Positive Rate)	68.91%	86.63%	76.45%
Specificity (i.e. true negative rate)	98.59%	77.64%	88.89%
Precision (i.e. positive predicted values)	91.72%	83.91%	68.71%
F1-Score (i.e. is the harmonic mean of precision and sensitivity)	78.69%	85.25%	72.37%

Challenges & Potential Improvements

- The Skewness reduction uses a constant set to 0.0025. This represents a +/-25bps buffer zone around zero, a.k.a. the ϵ range. This is quite a large number in trading terms, as spot trading margins are usually a few basis points. Furthermore, it is a synthetic buffer fabricated solely to rebalance the class instances. Indeed, the *Neutral* direction instances are usually underrepresented, and can cause models to break (e.g. QDA). It can also provoke skew in the training/test accuracy measures. The advantage of the approach is that rebalancing classes improve the model overall accuracy. However, the drawback relates to the potential for ignoring slow and continuous daily trend increment (decrement) over several days. Assuming there is a slow daily continuous trend increment of ϵ -0.01 on a daily basis for 10days, the data would be tagged with 10 days of *Neutral* directions. In reality, the cumulated effect of the 10days increase(decrease) induces a *Up (Down)* trend, over the time period. This could have a financial impact, as trading margins are much smaller than 25bps. This situation is unlikely but not entirely impossible. A better solution could be to use the daily asset volatility multiplied by a 'fudge' factor (e.g. 0.25) to try to stick closer to the asset real directionality.
- The feature selection was performed on the last time window, and was then used for all other time windows during the model training, optimisation and the test phases. This is a strong hypothesis that may need to be refined, as the data pattern carried in the last time window may not necessarily reflect the data pattern in each time window. More research needs to be carried out in this space.
- It would be interesting to investigate the test performance of neural networks such as multilayer perceptron (MLP), recursive neural networks (RNN), etc.
- The Ridge ($\alpha = 0$) and Lasso ($\alpha = 1$) models are subset of the Elastic Net regression. It would be interesting to see the impact of different levels of α (between 0 and 1) on the test accuracy rates.
- Initially, the `nnet.multinom()` function was selected to perform multinomial linear regression (due to the simplicity of its API). However, it requires that the training and validation/test sets are of the same dimension. This did not fit with the current sliding time window set-up. Therefore, it was dismissed in favour of the more complex `glmnet.glmnet()` function.
- Each model has been run against a number of different attributes. This is currently very a manual procedure. More investigation should be carried out (perhaps using the Caret library) to see whether code reduction could be implemented in this area. Also, dynamic formula creation should also be investigated [13].
- Currently the sliding time window supports 100 iterations. In order to improve the test performance accuracy, the time window should be increased (e.g. to 200, 300, etc. iterations).
- Some algorithms, such as the SVM are very slow when running the model with a matrix of optimisation factors, over multiple time windows. This could mean running multiple inner loops over the selected number of time windows. The number of loops depends on the number of factor to optimise. Therefore, it was decided to perform the optimisation on the last time window (the one closer to dataset end date). This may not be the best approach, as the optimisation should be run across all time windows for a given model. This can only be achieved with code refactoring, as explained below, as we need to move from a synchronous to a parallel model, to benefit fully from the server multi-cores.
- The current code has grown organically and has been built in a monolithic way. Each model is run after one another. This is not ideal for two reasons: i) running a model requires all the raw data to be load and massaged before any models can be run, ii) models cannot be run in parallel and iii) when a model fail, all following models may not run. The code refactoring should produce the following three independent modules:
 - An R module specialized in i) loading the raw dataset for an asset name (e.g. JPM), ii) removing the missing data row, iii) generating the derived data (e.g. technical indicators), iv) saving modified dataset into a file containing the asset name. This operation should be performed once (or anytime the derived data needs to be regenerated).
 - A module specialised in carrying out data analysis (e.g. the Kolmogorov Smirnov Normal Distribution test) and the generation of visuals.
 - A R module per model that deals with both model training, optimisation and test performance/accuracy measures.

Final Note:

- While running the different model with the polynomial and interaction between variables, two error messages (depending on the model type) were consistently displayed: i) *'Error in [model_name].default(x, grouping, ...) : rank deficiency in group Down'* or ii) *'Error in tree(factor(Direction) ~ Close_price_1day_lag^3 + Volume^2 + : trees cannot handle interaction terms'*. The experiment relating to mixing polynomial and attribute interactions was therefore commented out for each model that threw this exception.
- Another error message appeared, in the case of the Boosting model. The model could not be trained against only one attribute as the code threw the following error: *'Boosting with Close_price_1day_lagas only attribute => Error in x[1:nTrain, , drop = FALSE] : incorrect number of dimensions'*. In that instance, the experiment was also removed.

Conclusion

This experiment showed that for a sliding window of 100 iterations, over a 260 days' training, 65 days' validation and 5 days' testing contiguous data sets; the daily JPM stock trend can be explained endogenously at average test performance accuracy of approximately 80% (and standard deviation of 1.4%) for the Quadratic Discriminant Analysis (QDA), Ridge and Linear Discriminant Analysis (LDA) models. This means that approximately 20% of the price move is dependent on other factors. Therefore, it would be interesting to assess the impact of other variables such as fundamental or sentiment factors on the price trend. It would also be interesting to test the stability of the model accuracy over longer sliding windows and/or extra iterations, and to use other models such as neural network models (MPL, RNN, etc.).

References

- [1] Technical Analysis [Online], Available at: <https://www.tradingview.com/chart/technicalanalysis/>, [Accessed 27 March 2017]
- [2] Technical Indicator[Online], Available at: <http://www.investopedia.com/terms/t/technicalindicator.asp>, [Accessed 27 March 2017]
- [3] Vangie B., API - application program interface [Online], Available at: <http://www.webopedia.com/TERM/A/API.html>, [Accessed 27 March 2017]
- [4] Fundamental Analysis [Online], Available at: <http://www.investopedia.com/terms/f/fundamentalanalysis.asp>, [Accessed 27 March 2017]
- [5] Available at: <http://ichart.finance.yahoo.com/table.csv?s=JPM>, [Accessed 14 December 2016]
- [6] Bruder B., Dao T.L., Richard J.C, Roncalli T. (December 2011), *Trend Filtering Methods for Momentum Strategies*, Available at: <http://www.thierry-roncalli.com/download/lwp-tf.pdf>, [Accessed 27 march 2017]
- [7] Murphy J. (1986), *Technical Analysis of the Financial Markets: A Comprehensive Guide to Trading Methods and Applications*, New York Institute of Finance, pp.225-262.
- [7bis] Quantmod Quantitative Financial Modelling & Trading Framework for R [Online], Available at: <http://www.quantmod.com/>, [Accessed 27 March 2017]
- [8] Indicator Reference [Online], Available at: <http://www.fmlabs.com/reference/default.htm>, [Accessed 27 March 2017]
- [9] Moving Averages – Simple and Exponential [Online], Available at: http://stockcharts.com/school/doku.php?id=chart_school:technical_indicators:moving_averages , [Accessed 27 March 2017]
- [10] Jeni L.A., Cohn J.F., De La Torre F, *Facing Imbalanced Data: Recommendations for the Use of Performance Metrics* [Online], Available at: <http://www.pitt.edu/~jeffcohn/skew/PID2829477.pdf>, [Accessed 27 March 2017]
- [11] Morales R., Di Matteo T., Aste T [Online], Available at: Dependency structure and scaling properties of financial time series are related, <http://www.nature.com/articles/srep04589>, [Accessed 27 March 2017]
- [12] Volatility (finance) [Online], Available at: [https://en.wikipedia.org/wiki/Volatility_\(finance\)](https://en.wikipedia.org/wiki/Volatility_(finance)), [Accessed 27 March 2017]
- [13] Dynamic formula creation in R? [Online], Available at: <http://stackoverflow.com/questions/29711599/dynamic-formula-creation-in-r>, [Accessed 27 March 2017]
- [14] Gareth James , Daniela Witten , Trevor Hastie , Robert Tibshirani, *An Introduction to Statistical Learning: with Applications in R*, Springer Publishing Company, Incorporated, 2014, pp. 143, 312

Appendices

Appendix A – The Kolmogorov Smirnov Test Details

Check JPM 'Close_price_1day_lag/Up' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $2.648992e-13 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Close_price_1day_lag/Down' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $3.620437e-13 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Close_price_1day_lag/Neutral' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $0 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Volume/Up' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $0 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Volume/Down' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $0 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Volume/Neutral' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $0 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Close_price_2day_lag/Up' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $9.880985e-15 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Close_price_2day_lag/Down' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $1.338374e-12 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Close_price_2day_lag/Neutral' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $0 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_atr/Up' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $0 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_atr/Down' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $0 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_atr/Neutral' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $0 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Mfi/Up' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $0.3247685 > 0.05 \rightarrow H_0$ (the null hypothesis) is NOT rejected. There is not enough evidence to reject the hypothesis that the distribution is normal. Therefore, the data seems to follow normal distribution

Check JPM 'Mfi/Down' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 0.2228885 > 0.05 -> H0 (the null hypothesis) is NOT rejected. There is not enough evidence to reject the hypothesis that the distribution is normal. Therefore, the data seems to follow normal distribution

Check JPM 'Mfi/Neutral' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 0.3501365 > 0.05 -> H0 (the null hypothesis) is NOT rejected. There is not enough evidence to reject the hypothesis that the distribution is normal. Therefore, the data seems to follow normal distribution

Check JPM 'Rsi/Up' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 0.01189858 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Rsi/Down' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 0.08154547 > 0.05 -> H0 (the null hypothesis) is NOT rejected. There is not enough evidence to reject the hypothesis that the distribution is normal. Therefore, the data seems to follow normal distribution

Check JPM 'Rsi/Neutral' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 0.02306936 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_trueLow/Up' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 1.699751e-12 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_trueLow/Down' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 3.83249e-13 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_trueLow/Neutral' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 0 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_tr/Up' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 0 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_tr/Down' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 0 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_trLow/Neutral' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 0 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Close_price_5day_lag/Up' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 1.598721e-14 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Close_price_5day_lag/Down' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 6.695755e-13 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Close_price_5day_lag/Neutral' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 0 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_trueHigh/Up' is normally distributed

H0 = the data is normally distributed.

The ks p_value: 2.742251e-14 < 0.05 -> H0 (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_trueHigh/Down' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $4.746203e-13 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Check JPM 'Atr_trueHigh/Neutral' is normally distributed

H0 = the data is normally distributed.

The ks p_value: $0 < 0.05 \rightarrow H_0$ (the null hypothesis) is rejected. The data distribution does not seem to follow a normal distribution.

Appendix B – Code for Graphics Generation

```
255 #This function plots the different moving averages
256 plot_simple_moving_averages = function (stock_df, from_date, to_date){
257   #reset display to one graph per window
258   par(mfrow=c(1,1))
259
260   step_by = "6 mon"
261   date_format = "%d-%m-%Y"
262
263   from_date = get_from_date(stock_df, from_date);
264   to_date = get_to_date(stock_df, to_date);
265
266   #Only plot when dates are consistent
267   if(to_date > from_date){
268     #Reduce the data range to the required time range
269     stock_df = stock_df[stock_df$DateAsDate >= from_date & stock_df$DateAsDate <= to_date, ]
270
271     #Get the highest y value for the plot
272     max_height = max(stock_df$Sma20, stock_df$Sma50, stock_df$Sma100, stock_df$Sma200, na.rm = TRUE)
273     min_height = min(stock_df$Sma20, stock_df$Sma50, stock_df$Sma100, stock_df$Sma200, na.rm = TRUE)
274     #Plot
275     plot( stock_df$DateAsDate,
276           stock_df$Close,
277           col="black",
278           type="l",
279           xlab= "",
280           xlim=c(from_date, to_date),
281           ylab= "Stock Price ($)",
282           ylim=c(min_height, max_height),
283           main=paste(stock_name, "Close Price vs Simple Moving Average (SMA) Crossovers\n", from_date, "-", to_date, sep=" "),
284           las=2,
285           xaxt='n') #this removes the default x-axis label
286     axis.Date(1, at=seq(from_date, to_date, by=step_by), format=date_format, las=2)
287     lines(stock_df$DateAsDate, stock_df$Sma20, col="green", type="l", ylab= "Sma20")
288     lines(stock_df$DateAsDate, stock_df$Sma50, col="blue", type="l", ylab= "Sma50")
289     lines(stock_df$DateAsDate, stock_df$Sma100, col="orange", type="l", ylab= "Sma100")
290     lines(stock_df$DateAsDate, stock_df$Sma200, col="red", type="l", ylab= "Sma200")
291     #inset=c(-0.2,0) to have the legend outside of the plot
292     legend( x= "topleft", legend=c("Close", "Sma20", "Sma50", "Sma100", "Sma200"),
293           col=c("black", "green", "blue", "orange", "red"), lty=1, cex=0.8, horiz = TRUE)
294   }
295 }
296 }
```

```

298 #This function plots the different moving averages
299 - plot_exponential_moving_averages = function (stock_df, from_date, to_date){
300   #reset display to one graph per window
301   par(mfrow=c(1,1))
302
303   step_by = "6 mon"
304   date_format = "%d-%m-%Y"
305
306   from_date = get_from_date(stock_df, from_date);
307   to_date = get_to_date(stock_df, to_date);
308
309   #only plot when dates are consistent
310 - if(to_date > from_date){
311     #Reduce the data range to the required time range
312     stock_df = stock_df[stock_df$DateAsDate >=from_date & stock_df$DateAsDate <=to_date, ]
313
314     #Get the highest y value for the plot
315     max_height = max(stock_df$Ema20,stock_df$Ema50,stock_df$Ema100,stock_df$Ema200, na.rm = TRUE)
316     min_height = min(stock_df$Ema20,stock_df$Ema50,stock_df$Ema100,stock_df$Ema200, na.rm = TRUE)
317     #Plot
318     plot( stock_df$DateAsDate,
319           stock_df$close,
320           col="black",
321           type="l",
322           xlab= "",
323           xlim=c(from_date,to_date),
324           ylab= "Stock Price ($)",
325           ylim=c(min_height,max_height),
326           main=paste(stock_name,"Close Price Vs Exponential Moving Average (EMA) Crossovers\n", from_date, "-", to_date, sep=" "),
327           las=2,
328           xaxt='n') #this removes the default x-axis label
329     axis.Date(1, at=seq(from_date,to_date, by=step_by), format=date_format, las=2)
330     lines(stock_df$DateAsDate, stock_df$Ema20,col="green", type="l", ylab= "Ema20")
331     lines(stock_df$DateAsDate, stock_df$Ema50,col="blue", type="l", ylab= "Ema50")
332     lines(stock_df$DateAsDate, stock_df$Ema100,col="orange", type="l", ylab= "Ema100")
333     lines(stock_df$DateAsDate, stock_df$Ema200,col="red", type="l", ylab= "Ema200")
334     #inset=c(-0.2,0) to have the legend outside of the plot
335     legend( x= "topleft", legend=c("close", "Ema20","Ema50", "Ema100", "Ema200"),
336           col=c("black", "green", "blue", "orange", "red"), lty=1, cex=0.8, horiz = TRUE)
337
338   }
339 }
340

```

```

341 #This function plots the relative strength index (rsi)
342 plot_relative_strength_index = function (stock_df, from_date, to_date){
343   step_by = "3 mon"
344   date_format = "%d-%m-%Y"
345
346   from_date = get_from_date(stock_df, from_date);
347   to_date = get_to_date(stock_df, to_date);
348
349   #Display the graphs (one under the other), i.e. in a 2 rows/1 col grid system
350   par(mfrow=2:1)
351
352   #Only plot when dates are consistent
353   if(to_date > from_date){
354     #Reduce the data range to the required time range
355     stock_df = stock_df[stock_df$DateAsDate >=from_date & stock_df$DateAsDate <=to_date, ]
356
357     #Get the highest y value for the plot
358     max_height = max(stock_df$Close, na.rm = TRUE)
359     min_height = min(stock_df$Close, na.rm = TRUE)
360     #Plot
361     plot( stock_df$DateAsDate,
362           stock_df$Close,
363           col="black",
364           type="l",
365           xlab= "",
366           xlim=c(from_date,to_date),
367           ylab= "Stock Price ($)",
368           ylim=c(min_height,max_height),
369           main=paste(stock_name,"Close Price\n", from_date, "-", to_date, sep=" "),
370           las=2,
371           xaxt='n') #this removes the default x-axis label
372     axis.Date(1, at=seq(from_date,to_date, by=step_by), format=date_format, las=2)
373
374     max_height = max(stock_df$Rsi, na.rm = TRUE)
375     min_height = min(stock_df$Rsi, na.rm = TRUE)
376     plot( stock_df$DateAsDate,
377           stock_df$Rsi,
378           col="green",
379           type="l",
380           xlab= "",
381           xlim=c(from_date,to_date),
382           ylab= "RSI",
383           ylim=c(min_height,max_height),
384           main=paste(stock_name,"RSI (14days)\n", from_date, "-", to_date, sep=" "),
385           las=2,
386           xaxt='n') #this removes the default x-axis label
387     axis.Date(1, at=seq(from_date,to_date, by=step_by), format=date_format, las=2)
388
389     #Two vertical lines, at ordinate = 30 and 70, to show the RSI signal trigger
390     abline(h=70)
391     abline(h=30)
392
393     #reset display to one graph per window
394     par(mfrow=c(1,1))
395   }
396 }
397

```



```

399 #This function plots the relative strength index (rsi)
400 plot_prices_volume = function (stock_df, from_date, to_date){
401   step_by = "3 mon"
402   date_format = "%d-%m-%Y"
403
404   from_date = get_from_date(stock_df, from_date);
405   to_date = get_to_date(stock_df, to_date);
406
407   #Display the graphs (one under the other), i.e. in a 2 rows/1 col grid system
408   par(mfrow=2:1)
409
410   #Only plot when dates are consistent
411   if(to_date > from_date){
412     #Reduce the data range to the required time range
413     stock_df = stock_df[stock_df$DateAsDate >=from_date & stock_df$DateAsDate <=to_date, ]
414
415     #Get the highest y value for the plot
416     max_height = max(stock_df$close, na.rm = TRUE)
417     min_height = min(stock_df$close, na.rm = TRUE)
418     #Plot
419     plot( stock_df$DateAsDate,
420           stock_df$close,
421           col="black",
422           type="l",
423           lwd=1.5,
424           xlab= "",
425           xlim=c(from_date,to_date),
426           ylab= "Stock Price ($)",
427           ylim=c(min_height,max_height),
428           main=paste(stock_name,"Close Price\n", from_date, "-", to_date, sep=" "),
429           las=2,
430           xaxt='n') #this removes the default x-axis label
431     axis.Date(1, at=seq(from_date,to_date, by=step_by), format=date_format, las=2)
432     lines(stock_df$DateAsDate, stock_df$close_price_1day_lag,col="orange", type="l", ylab= "1Day Price Lag")
433     lines(stock_df$DateAsDate, stock_df$close_price_4day_lag,col="blue", type="l", ylab= "4Day Price Lag")
434     #inset=c(-0.2,0) to have the legend outside of the plot
435     legend( x= "topleft", legend=c("Close", "Close 1-Day Lag", "Close 4-Day Lag"),
436            col=c("black", "orange", "blue"), lty=1, cex=0.8, horiz = TRUE)
437
438     max_height = max(stock_df$volume, na.rm = TRUE)
439     min_height = min(stock_df$volume, na.rm = TRUE)
440     plot( stock_df$DateAsDate,
441           stock_df$volume,
442           col="green",
443           type="l",
444           xlab= "",
445           xlim=c(from_date,to_date),
446           ylab= "Volume",
447           ylim=c(min_height,max_height),
448           main=paste(stock_name,"Volume\n", from_date, "-", to_date, sep=" "),
449           las=2,
450           xaxt='n') #this removes the default x-axis label
451     axis.Date(1, at=seq(from_date,to_date, by=step_by), format=date_format, las=2)
452
453     #reset display to one graph per window
454     par(mfrow=c(1,1))
455   }
456 }

```

```

459 #This function plots the relative strength index (rsi)
460 plot_log_returns = function (stock_df, from_date, to_date){
461   #reset display to one graph per window
462   par(mfrow=c(1,1))
463
464   step_by = "3 mon"
465   date_format = "%d-%m-%Y"
466
467   from_date = get_from_date(stock_df, from_date);
468   to_date = get_to_date(stock_df, to_date);
469
470   #Only plot when dates are consistent
471   if(to_date > from_date){
472     #Reduce the data range to the required time range
473     stock_df = stock_df[stock_df$DateAsDate >=from_date & stock_df$DateAsDate <=to_date, ]
474
475     #Get the highest y value for the plot
476     max_height = max(stock_df$Log_returns, na.rm = TRUE)
477     min_height = min(stock_df$Log_returns, na.rm = TRUE)
478     #Plot
479     plot( stock_df$DateAsDate,
480           stock_df$Log_returns,
481           col="black",
482           type="l",
483           xlab= "",
484           xlim=c(from_date,to_date),
485           ylab= "Log Returns ($)",
486           ylim=c(min_height,max_height),
487           main=paste(stock_name,"Log Returns (from close prices)\n", from_date, "-", to_date, sep=" "),
488           las=2,
489           xaxt='n') #this removes the default x-axis label
490     axis.Date(1, at=seq(from_date,to_date, by=step_by), format=date_format, las=2)
491   }
492 }
493 }
494

```

Code snippet relating to the generation of the Times Series graphs.