

FINANCIAL TIME SERIES FORECASTING – A MACHINE LEARNING APPROACH

Alexiei Dingli and Karl Sant Fournier

Department of Artificial Intelligence, University of Malta, Malta

ABSTRACT

The Stock Market is known for its volatile and unstable nature. A particular stock could be thriving in one period and declining in the next. Stock traders make money from buying equity when they are at their lowest and selling when they are at their highest. The logical question would be: "What Causes Stock Prices To Change?". At the most fundamental level, the answer to this would be the demand and supply. In reality, there are many theories as to why stock prices fluctuate, but there is no generic theory that explains all, simply because not all stocks are identical, and one theory that may apply for today, may not necessarily apply for tomorrow. This paper covers various approaches taken to attempt to predict the stock market without extensive prior knowledge or experience in the subject area, highlighting the advantages and limitations of the different techniques such as regression and classification. We formulate both short term and long term predictions. Through experimentation we achieve 81% accuracy for future trend direction using classification, 0.0117 RMSE for next day price and 0.0613 RMSE for next day change in price using regression techniques. The results obtained in this paper are achieved using only historic prices and technical indicators. Various methods, tools and evaluation techniques will be assessed throughout the course of this paper, the result of this contributes as to which techniques will be selected and enhanced in the final artefact of a stock prediction model. Further work will be conducted utilising deep learning techniques to approach the problem. This paper will serve as a preliminary guide to researchers wishing to expose themselves to this area.

KEYWORDS

Stock Prediction, Fin Tech, Machine Learning, Time Series Forecasting, Data Science

1. INTRODUCTION

A stock is a piece of a company, and when an individual buys stock, he/she is essentially becoming a part owner of said company. Throughout this thesis, the term stock may be referred to as 'security', 'asset', 'equity' or 'share', all these terms refer to the same entity. The stock market is a community that brings together traders that want to buy and sell stocks. Among other metrics, the stock price of a company reflects its performance. The law of demand and supply holds that when the demand for a commodity increases or its supply decreases, the price of the commodity increases. On the other hand, if the demand decreases or the supply increases, the commodity price decreases. The same applies for the stock market. Stocks are traded in a stock exchange, this is an organised market through which investors can buy and/or sell securities in a public manner. The ideal scenario is to buy when there is a good chance that the stock price will increase, and to sell when there is a chance that the stock price is about to decrease, hence yielding profit from the difference between the price the stock was bought and the price the same stock was sold. In a nutshell, the role of a stockbroker is to know when to buy and when to sell stocks with the goal of earning a profit.

2. BACKGROUND

2.1. LITERATURE REVIEW

Given the nature of this topic, many researchers and traders have exploited the area of statistical analysis and machine learning in order to generate a profit from their investments. Although the initial bias would be that any successful enterprise would not share its methods in order to keep profits to themselves, there have been a good number of researchers who have published their findings.

2.1. 1 Statistical Methods

Regression analysis is a statistical tool, used to investigate relationships between variables [1].

[2] uses regression to study the relationship between news and stock price changes, in an effort to improve the performance of the conventional stock price forecasting process. The author regresses the weekly stock price changes on the news values from the previous week. Results show that there is a very significant correlation between the changes in stock price and the news values at the start of the week, such that one unit increase in news values, results in approximately 6.33 unit increase in stock price change. The author does note however that the R-squared value is minimal (0.09243). This indicates the proportion of the variability of the weekly price changes predicted by the model. This small value means that the model does not predict the price changes accurately.

Traders and researchers have exploited statistical techniques such as autoregressive integrated moving average (ARIMA) to forecast financial time series. ARIMA models work on stationarized time series. A stationary time series is one whose statistical properties are constant over time [3]. Many statistical forecasting methods are based on the assumption that the time series can be rendered stationary through the use of transformations. Once a prediction is made on the stationary time series, it can be converted back to the original series using the same transformations that made it stationary.

2.2 Machine Learning Methods

Machine learning is the science of getting computers to take decisions without being explicitly programmed to do so. This phenomenon has given positive results in experiments such as speech recognition, self-driving cars, image recognition and a number of other areas. One of the most exciting things about machine learning is that it can be applied to almost anything. In a chaotic life such as that of the 20th century, certain decisions are best delegated to machines. These decisions can range from trivial decisions to critical decisions, according to the nature of the problem at hand. Machine learning techniques can be divided into supervised, unsupervised, association and reinforcement learning. With supervised machine learning, the data being passed to the algorithm is labelled, indicating that a certain data pertains to some class. Some of these techniques include Decision Trees, Regression, Support Vector machines and Naive Bayes. On the other hand, with unsupervised machine learning, the data is not labelled, and the task of the algorithm is to draw inferences from the data without it being labelled from beforehand. Some of these unsupervised techniques include K-Means and Gaussian Mixture models. In the field of finance, various machine learning techniques have been used for predicting stock market prices. Technical indicators and market sentiment are passed as features to the various models, which output either a binary or numerical output indicating a direction or change in price respectively.

In the field of machine learning, Support Vector Machines (SVMs) [4] are becoming increasingly popular, a similar enthusiasm was shown when artificial neural networks were discovered. SVMs were first developed by Vladimir N. Vapnik and Alexey Ya Chervonenkis in 1963, these SVMs

were able to classify linear problems in a much more efficient way than ANNs. Later in 1992, along with another two scientists, Vapnik suggested a method called the kernel trick, using the kernel trick, SVMs were able to create non-linear classifiers.

[5] uses SVM to predict direction of future price changes. In the study it showed that the resulting accuracy was highly dependent on the various parameters passed to the SVM model, best results were achieved when using the Gaussian Radial Basis Function kernel as a non-linear classifier. The value of parameter C also had an effect on the results obtained. Features used where selected technical indicators and oscillators. The author compared the results of the SVM to other methods such as Neural Networks using Back-propagation and discovered that better results were achieved using the SVM approach.

Decision trees are a type of classifier which can be represented in a form of a tree. A tree can be divided into a root node, internal nodes and leaf nodes. The feature splitting under the root and internal nodes are calculated using entropy information gain. Decision trees are based on the heuristic of ‘Occam’s Razor’, which favours a short, generalised tree over a complex tree which may overfit the training data.

[6] designs a text-based decision tree for stock market forecasting. Price data is transformed into textual form according to a number of rules, such that each indicator is given a value of either up trend, down trend or no trend. The output is also based on a pre defined rule which outputs whether to buy, hold or sell the stock. Through this text based decision tree approach, the authors achieve an accuracy of 90.82% which they proved to be more effective than the normal, non text-based decision tree. The authors also noted that the model for their text based approach resulted in a simpler, shorter tree than that of the normal decision tree.

[7] use a multi-agent approach for forecasting, namely they use technical and fundamental analysis, neural networks and liquidity methods to forecast the market behaviour. They compare the performance of standard feed-forward networks, Elmand and Jordan RNNs and an architecture known as NEAT (Neuroevolution of augmenting topologies \cite{Stanley}). These architectures are evaluated on their ability to predict price changes on the Bucharest Stock Exchange (BSE). The best results are achieved using the NEAT architecture, which is a genetic algorithm developed by Ken Stanley. This methodology attempts to find a balance between the diversity of evolved solutions and their fitness.

Whilst decision trees have achieved promising results in the field of stock market prediction, some have ventured to the use of ensemble learning algorithms such as random forest. \footnote{\url{https://arxiv.org/pdf/1605.00003.pdf}} make use of the random forest ensemble to counteract against the overfitting problem of decision trees. The author states that this occurs because decision trees are very sensitive to noise in the data given that they have a very low bias and high variance. Through the Random forest algorithm, the author trains multiple decision trees on a different subspace of the feature space, causing a slight increase in bias. By taking this approach, none of the trees make use of the entire training data, but rather it is shared amongst the various trees in the ensemble. Through this method, the authors achieve an accuracy between the range of 85% and 95% for long term predictions, which is an improvement on what was achieved on the decision tree based approach of [8].

[9] make use of ANNs to forecast stock price direction for the Nikkei 225 index, using a number of popular technical indicators as features. As part of the experiment, the authors opt to divide their technical indicator features into two types. This is done with no mention of the criteria used to divide these features, constantly referring to these as 'type 1' and 'type 2'. which would have been helpful for the reader to better understand the intention for this split. Nonetheless, the authors report an accuracy of 60.87% on the 'type 1' dataset and 81.27% on the

'type 2' dataset for next day price direction forecast, which is an outstanding result given the volatile and unstable nature of day to day prices in the stock market.

3. METHODOLOGY

In this chapter, we discuss the format of the different outputs we would like to achieve, from our experiments. We then identify the various data sources used and highlight any transformations done to these datasets. Given we will be dealing with a large number of features, some of which may be deemed irrelevant, we will be performing feature selection techniques in order to identify the most important features. Once our features are established, we then discuss the various approaches taken to achieve our desired outputs, through machine learning techniques.

3.1. Structure of desired outputs

The below structure is that of which we plan to formulate our outputs, which is mainly threefold:

- Next Period Direction

This output concerns the direction of the price movement with respect to the next period. Natively, this can be considered as a classification problem, whereby we classify each record as 'Up' or 'Down'.

- Next Period Price Change

This output concerns the numerical value of change, for the price movement with respect to the next period. This can be considered to be a regression problem given that we are not dealing with predefined output classes, but with a numerical, continuous value.

- Next Period Actual Price

This output concerns the numerical value for the actual next period price. This can also be considered to be a regression problem given the nature of the output.

The above three outputs will be computed for a number of periods: Daily, Weekly, Monthly, Quarterly and Yearly. By means of this comprehensive approach we will be in a position to identify an optimal combination for forecasting.

3.2. Feature Identification

3.2.1. Historic Stock Prices and Technical Indicators

Historic stock prices are a good place to start, given that the data is publicly available through a number of sources. One of the most widely used mediums is yahoo finance, which provides unlimited historic daily prices in csv format on demand. A python API known as 'yahoo-finance' is used to retrieve these historic prices from 2003 till 2016, which are then stored in our relational database for future use. These historic prices are then used to formulate what are known as technical indicators. These indicators consist of various financial equations using historic prices which are used to paint a clearer picture as to the condition of the company for that point in time. Investopedia describe these indicators as any class of metrics whose value is derived from generic price activity. Technical indicators evaluate price levels, direction and momentum amongst others. Rather than developing custom functions for each indicator, a library called 'TA-Lib' was identified. This library openly provides 200 technical indicators, and is designed in the form of an open source API. Given that

python is the language of choice for experimentation, a Python wrapper for this API is used. The technical indicators used in our approach include a combination

of momentum, volume, volatility and cycle based indicators. As their name suggests, each class of indicator focuses on a different aspect of market activity.

3.2.2. Currency Exchange Rates

When investors purchase foreign equity, they are actually placing a bet on two elements: one on the stock itself and the other on the currency the stock trades in. In order to evaluate this hypothesis, we collect a number of the most popular exchange rates to incorporate within our predictive model. Currencies were collected from yahoo finance and Quandl data repository, and these include: EUR/USD, GBP/USD and BITCOIN/USD. As one may note, USD was taken as the benchmark since it is one of the most prominent currencies in the global stock exchange.

3.2.3. World Indices

Investopedia describes an index as a measure of change in a securities market. Indices consist of a hypothetical portfolio containing a group of securities that generally represent the performance of the overall market. Typical examples of such indices include the Standard & Poor's 500 (S&P500), Dow Jones Industrial Average (DJIA) and the Nasdaq Composite Index (NASDAQ-100). These indices are retrieved from the yahoo finance API.

3.2.4. Commodity Prices

A number of popular commodities were selected as features for our experiments. Prices for Gold and Oil commodities are retrieved from the Quandl data repository, store them in our relational database and incorporate them into our dataset as time series data.

3.3. Feature Scaling

The features identified earlier are of numeric value, however the ranges for each feature vary widely. A number of machine learning classifiers calculate the distance between points by using Euclidean distance. When adopting this methodology for a broad dataset, one of the features may have such a large range, that it automatically takes precedence over the other features, when in fact other features with smaller ranges may possibly be of more importance to the predictor. Apart from the affects on Euclidean distance, gradient descent algorithms are known to converge faster with scaled features due to the nature of the method's architecture [10]. For the purpose of this study, we will be using Z-Score scaling on all our numeric features.

3.4. Feature Selection

Now that various data sources and features have been identified, the next step is to integrate them with our dataset and evaluate which are most important for the problem at hand. In total, we have over 70 features comprising of technical indicators (Overlap, Volume, Momentum, Volatility and cycle), currency exchanges, commodity prices and world indices. Rather than passing all features to our predictive model we use feature selection techniques in order to statistically identify the most relevant features for the data at hand.

In order to gauge the effectivity of the forecast over different periods, a number of different time spans are taken into consideration, namely daily, weekly, monthly, quarterly and yearly averages. When taking this approach, all features in the dataset are averaged according to periodicity and fed into our feature selection mechanism. Given that the initial dataset comprises of one record per trading day, resulting in 252 records per year (one per official trading day as quoted by

NASDAQ), this periodic average transformation has an effect on the final size of the final dataset. For instance, if we were to take weekly averages, the number of records per year would be taken down to 52, for monthly averages this would be 12 and for quarterly we would have only four records. In machine learning, ideally a substantial amount of examples is fed to the classifier in order to be able to classify the outcome of unseen examples. To tackle this, rather than considering a single stock, numerous stocks are concatenated together in order to have a substantial dataset. The criteria for concatenated stocks is the nature of the company behind each equity. Technology companies are considered as one dataset and Finance companies form their own separate dataset. All values are scaled as discussed in the previous section, and the direction of the next day price is calculated accordingly, with each record being independent from the other.

Sklearn library known as ‘feature selection’ is used for identifying the most important and relevant features. We make use of a feature called ‘SelectKBest’ which makes use of one of two functions to evaluate each feature and assign a score accordingly. One function ‘f_classif’ is used when the target output is a predefined class. This function makes use of the ANOVA F-test and outputs a score known as the F-Value. ANOVA test outputs the proportion of variance explained by the features, the higher the ratio, the more proportion. The second function, known as ‘f_regression’ is used for regression type problems when the output is of continuous value. This function computes the cross correlation between each feature and the target output to assign scores accordingly.

Feature selection techniques are performed for each combination of Output type and periodicity. Stock types are divided into Technology and Finance; Output types that will be considered in this study are next period direction, change in price and actual price; periods considered are Daily, Weekly, Monthly, Quarterly and Yearly as mentioned earlier. Through this approach we are making no undue assumptions regarding the importance of features for our output. For instance, one feature that may be important to forecast the next day direction, may be not as important to forecast the next day price itself. A typical example of this scenario could be the importance of the previous day closing price, whereby when predicting the next day price, this feature is vital since the next day price is built as an addition or subtraction to it, it provides no indication whatsoever on the possible direction of the price movement for the following day. For this reason, all features are taken into consideration and will be evaluated accordingly. This approach considers that importance of features may vary according to whether the stock is for instance a technology stock or a financial stock. We are also making the assumption that features that are important for predicting the next day direction may differ from those that point towards the direction of the price for the following quarter. Table 1 illustrates this comprehensive approach which is performed for all features.

Table 1. Comprehensive approach taken for feature selection

Next Period:	Diretion	Change in Price	Actual Price
Features	Day	Day	Day
	Week	Week	Week
	Month	Month	Month
	Quarter	Quarter	Quarter
	Year	Year	Year

3.5. Dataset balancing, training and test split

In order to have unbiased models, it is imperative to have a well structured dataset balancing mechanism, and an adequate ratio between training and test datasets. A number of methods exist when dividing data into the various sets required. In other machine learning problems, it may be ideal to balance, shuffle the dataset and randomly split the training and test sets into 70% and 30% respectively. Such an approach ensures that bias is not present in the data. However, the nature of the data at hand is that of a time series, if this approach were to be taken, it is possible to have so called 'future instances' in our training set and 'past instances' in the test set. Although given that we are treating each record separate from the other, the latter method could still be deemed as valid. Nonetheless, we will be taking an approach which is more suited for time series data.

We will be utilising data from 2003 till 2013(11 years) as our training set and 2014 till 2016 (3 years). This ensures a substantial amount of data for training, whereby 11 years of stock movements should cover a wide range of long and short term trends. The remaining three years will be used for testing, this approach will ensure that we are forecasting and calculating our evaluation metrics on unseen, out-of-sample data.

When talking about a balanced dataset, this mainly concerns our classification task whereby our training examples can be either 'Up' or 'Down'. In order to assist the machine learning model, it is in our best interest to provide training samples which have an equal amount of examples for each class, otherwise we run the risk that during training, the classifier becomes biased towards the class with more examples. This is done by randomly shuffling the training set, identifying which subset has the least values n , retaining the first n from the largest subset, and lastly re ordering by date.

When it comes to the test set, we will not be performing any balancing so as to leave the data in its natural distribution. Like this we will be evaluating an unbiased classifier on unseen data with its native distribution.

3.6. Classification Experiments

For classification, our objective is to predict a predefined output label per example based on the variables available to us. We shall approach this task by industry, firstly focusing on the tech stock industry, followed by financial company stocks whereby we compare and contrast between the two.

Our approach consists of forecasting the price movement direction for the next period, which may be daily, weekly, monthly, quarterly or yearly. The output variable is constructed in such a way that it indicates whether the next period average price goes 'up' or 'down'. Through this, we enable the classifier to identify trends from the features identified, at the same time assessing whether trend detection is more suitable for the short term or the long term. As discussed during the literature review, some analysts state that day to day fluctuations are random, through this experiment we would like to test this assumption and compare it to detecting trends in a longer time span.

By taking the approach of evaluating multiple classifiers utilising the same dataset, we minimise the risk of focusing on a single classifier which may not be the best fit for the data at hand. Now that we've prepared our datasets, a number of algorithms can be used for classification, we will be utilising the below:

- K-Nearest Neighbors
- Logistic Regression
- Naive Bayes
- Support Vector Classifier (SVC)
- Decision Trees
- Random Forest
- Multi Layer Perceptron (MLP)
- Ada Boost
- QDA

From the output of each algorithm we can compute metrics such as accuracy, precision and recall. These metrics could give us an idea on the effectiveness of each, striving for a high number of true positives and false negatives, as against false positives and negatives. Through these results we will be in a position to compare and contrast each algorithm across various dimensions such as number of features, different periodicities and the aforementioned stock categories.

3.7. Regression Experiments

As discussed earlier, we are dividing our predictions into three forms, one of which is a classification problem, and the other two being a regression problem with a continuous output. In this section we will be discussing the approach taken to tackle the two regression tasks. We will be evaluating results by means of the Root Mean Squared Error (RMSE) and below is a list of regressor algorithms that are used for both tasks:

- Linear Regressor
- Support Vector Regressor
- Decision Tree Regressor
- Ada Boost Regressor
- Random Forest Regressor
- K-Nearest Neighbors Regressor
- Bagging Regressor

The first regression task is to forecast the next period price change. As was done in our classification task, we will be sampling the dataset into different periodicities for daily, weekly, monthly, quarterly and yearly. Given these periodicities, we would like to forecast the actual change from one period to another. The end scope here is similar to that of our classification task where we classified the direction of the next period, however here we are forecasting as to what extent the price will move up or down. The data that will be fed into the various algorithms is the output of the feature selection exercise that was done for this form of data, theoretically the features that were relevant for the direction may not be relevant for the change in price.

As mentioned earlier, this task is divided into two sub-tasks. We now attempt to predict the actual next period price itself, which is also a continuous value. At face value, we assume that the features which were important for the change, and direction should be relevant for the actual price.

However, when simply looking at the data, there is a very good chance that those features which were important for our previous two tests will not be as important for the actual price itself.

This is assumed because a volume or momentum indicator which is a calculated metric does not provide the algorithm any knowledge on the price itself. On the other hand, in our full dataset we have features such as moving averages and last period prices which do provide the predictor a base price to start on, whilst other indicators, should in theory provide an indication as to any possible future fluctuations.

4. RESULTS AND EVALUATION

4.1. Feature Selection Results

In this section, we will be discussing the outcome obtained from our feature selection methodology. Discussion and evaluation of feature selection results will be done per stock type. Feature selection scores are tabulated as per Table 2 for each combination of forecast type and periodicity. One can note, that the number of examples indicated in the second row of the table, decreases with each periodicity sampling as discussed earlier. Given that the table is considerably large due to the number of features, in order to easily understand the contents, all feature scores which are above the global average for each forecast type are highlighted in green. These cells are to be considered as highest statistical importance, with reference to the output class. Towards the right of each section, a sum is calculated for each feature, indicating the overall relevance to the industry as a whole rather than to a specific period. The cells highlighted in yellow signify a total of above average for each forecast type.

As discussed earlier, all features are taken into consideration for the tech industry, as we are relying on our feature selection methodology to outline the most important features. Table 2 shows the feature selection results for the Tech industry. As one may note, the number of shaded cells varies from one periodicity to another, signifying that relevance of certain features differs when evaluating the same dataset with different periodicity.

Table 2. Feature Selection Scores for the Tech Industry

Forecast Type	Direction					S					Change					S					Price					S				
	Periodicity		D	W	M	Q	Y	U	D	W	M	Q	Y	U	D	W	M	Q	Y	U	D	W	M	Q	Y	U				
	No. of Examples	22661	4699	1079	359	89	M	22661	4699	1079	359	89	M	22661	4699	1079	359	89	M	22661	4699	1079	359	89	M					
AD	0.16	6.73	0.27	2.54	0.04	9.74	0.47	0.11	0.04	0.00	0.26	0.89	2103.56	4401.72	971.76	283.91	52.25	26753.21												
ADOSC	1.40	27.94	48.34	16.18	20.56	114.42	6.90	23.11	47.33	9.59	8.76	95.69	36.23	8.99	3.56	1.70	0.77	51.26												
ADX	0.02	0.80	0.48	0.04	1.78	0.37	0.69	1.50	1.83	3.49	7.88	105.29	20.55	5.82	3.07	2.74	135.48													
ADXR	0.80	0.00	1.66	0.26	0.69	3.40	1.73	3.15	2.63	2.44	3.56	13.50	118.10	24.05	5.65	3.20	2.85	153.85												
APO	0.14	4.17	34.77	17.80	19.85	76.72	0.00	8.24	53.17	28.15	51.27	140.83	482.33	100.05	32.96	21.16	23.58	660.08												
AROONOSC	1.74	20.33	92.36	40.92	24.00	179.34	0.22	19.27	108.47	31.50	24.71	184.17	68.39	16.19	7.63	5.88	4.90	102.98												
ATR	0.12	6.98	5.13	6.06	2.21	20.53	0.94	2.11	1.65	0.75	0.61	6.07	35482.81	7339.01	1763.02	567.22	83.55	45215.61												
Adj Close	0.00	4.75	0.49	3.56	0.44	9.25	0.92	0.04	0.19	0.25	0.04	1.44	14504440.00	1002745.00	50783.42	5055.73	228.99	15563253.13												
BITUSD_ROCP	3.64	4.03	4.58	0.02	0.01	12.27	1.36	2.29	5.95	0.50	0.10	10.21	0.27	0.10	0.53	0.83	2.09	3.82												
BOP	3.77	426.22	125.66	48.58	15.95	620.37	0.01	467.89	147.66	35.22	21.49	672.27	8.03	6.89	6.69	5.25	1.08	27.93												
CCI	0.21	121.85	152.30	57.02	30.72	362.11	0.48	123.49	185.73	45.80	30.16	385.66	26.62	8.07	4.08	3.63	3.43	45.83												
CMO	0.10	74.40	87.77	34.13	36.50	232.90	3.45	72.59	89.43	23.39	32.85	221.70	132.86	31.11	11.03	5.91	6.47	187.38												
ChangeLastPeriod	2.94	91.49	60.80	3.54	0.68	159.45	1.69	198.29	108.62	7.01	2.89	316.71	24.05	33.44	36.39	28.67	33.28	155.83												
Close	0.00	4.75	0.49	3.56	0.44	9.25	0.92	0.04	0.19	0.25	0.04	1.44	14504440.00	1002745.00	50783.42	5055.73	228.99	15563253.13												
DEMA	0.02	7.65	1.34	4.11	0.46	13.61	0.98	0.38	0.05	0.13	0.02	1.07	3628307.00	458521.40	36153.31	4678.75	220.02	4126160.46												
DJ ROC	0.05	37.63	39.05	12.43	9.59	98.77	0.44	38.49	38.82	6.84	4.78	89.37	72.99	18.37	8.29	5.71	2.05	107.40												
DX	1.05	8.73	0.60	0.11	0.09	10.59	1.04	1.93	0.12	0.52	3.04	6.65	36.86	7.76	3.09	2.47	2.11	52.41												
EMA	0.01	8.57	1.93	4.93	0.75	16.20	0.51	0.64	0.21	0.03	1.40	182484.00	279532.10	28215.88	4057.99	203.38	2138493.35													
EURUSD_ROCP	2.34	2.98	9.63	4.56	1.77	21.27	2.23	0.45	8.70	5.60	0.05	17.02	27.50	7.10	1.75	0.59	5.88	42.81												
FTSE ROC	0.05	29.87	32.24	20.76	9.41	92.32	0.13	22.67	37.42	10.70	2.40	73.32	6.91	1.71	0.92	1.58	0.01	11.14												
GBPUSD_ROCP	2.53	8.39	4.95	3.56	8.12	27.59	4.43	4.98	3.17	2.80	1.20	18.57	10.91	2.86	1.35	0.69	16.50													
GOLD_ROCP	0.42	7.35	0.00	0.03	0.84	8.63	5.13	8.24	0.28	0.09	8.03	21.77	149.27	38.33	16.90	14.39	42.21	261.11												
HT_DCPERIOD	0.02	1.97	0.04	0.56	0.76	3.35	0.74	1.71	1.21	4.48	1.82	9.96	22.01	4.00	1.43	1.22	2.14	30.80												
HT_DCPHASE	0.00	0.01	4.66	25.24	25.03	55.14	1.12	0.57	4.53	14.40	30.25	50.87	38.52	9.41	4.66	3.55	5.16	61.31												
HT_PHASOR_INPHASE	0.82	0.54	84.35	31.00	18.79	135.51	0.12	0.00	151.85	49.99	51.57	253.55	97.13	31.01	24.27	21.53	18.05	191.99												
HT_PHASOR_QUADRATURE	3.63	0.54	26.89	0.78	0.25	32.08	10.44	0.55	63.76	3.08	0.01	77.83	0.33	0.56	0.51	1.29	0.00	2.69												
HT_SINE_LEADSINE	0.91	0.00	4.44	20.52	27.43	53.20	2.30	1.15	4.37	11.97	25.86	45.65	54.47	11.55	5.34	4.27	5.82	81.44												
HT_SINE_SINE	0.17	3.51	38.67	0.25	6.64	49.25	2.91	0.01	45.50	2.64	3.24	54.28	11.57	2.37	2.74	2.15	4.46	23.29												
HT_TRENDLINE	0.02	8.73	2.15	5.07	1.73	16.69	0.40	0.68	0.36	0.02	0.04	1.46	1855358.00	246512.60	26647.95	4035.65	201.28	1850933.48												
HT_TRENDSIDE	7.10	5.13	0.00	0.11	0.35	12.69	0.26	1.09	1.49	0.67	2.14	5.65	138.88	43.34	15.71	13.28	9.83	218.04												
High	0.00	5.41	0.59	3.65	0.45	10.10	0.78	0.03	0.11	0.22	0.04	1.17	11159410.00	885254.00	486939.99	4992.80	226.35	12086825.13												
IXIC ROC	0.58	48.46	41.22	18.25	10.91	119.49	1.10	51.33	45.95	10.09	5.48	113.93	68.82	16.55	8.48	7.03	20.96	102.96												
KAMA	0.01	8.53	1.91	4.83	0.72	15.99	0.61	0.73	0.24	0.04	0.00	1.62	1618055.00	254322.90	27328.18	4031.29	208.64	190346.00												
Low	0.00	5.26	0.51	3.59	0.44	9.80	0.76	0.00	0.17	0.25	0.04	1.21	11103190.00	911513.60	49965.12	5035.09	228.08	12068001.86												
MA	0.01	8.99	2.29	5.13	0.77	17.19	0.52	0.80	0.41	0.02	0.00	1.75	1421415.00	228495.70	25428.19	3950.46	202.55	1679491.91												
MFI	0.20	26.39	77.97	35.63	22.23	162.42	0.13	35.84	64.89	28.21	20.54	179.61	41.91	9.90	4.32	3.55	1.41	61.09												
MIDPOINT	0.01	7.72	1.50	4.40	0.58	14.21	0.49	0.34	0.08	0.08	0.01	1.01	3265047.00	429033.00	34892.04	4466.84	216.23	373655.62												
MIDPRICE	0.01	7.88	1.52	4.40	0.59	14.40	0.53	0.40	0.09	0.08	0.01	1.11	3082744.00	412132.00	34559.92	4456.90	215.99	3542488.81												
MINUS_DI	0.38	46.24	78.13	21.01	13.83	159.69	0.37	46.71	56.29	7.57	9.03	119.97	0.11	0.11	0.09	0.16	0.00	0.47												
MINUS_DM	0.75	19.79	24.38	9.99	4.36	69.26	0.66	20.28	23.11	3.01	2.70	49.77	9881.65	2170.66	584.03	235.00	52.20	12923.54												
MOM	4.92	44.26	129.29	31.95	24.76	235.18	18.33	74.02	279.52	54.79	60.54	490.93	246.10	64.26	22.46	28.36	22.08	21.45	380.22											
NATR	0.09	2.99	11.20	0.46	1.41	16.15	0.31	7.01	7.45	2.87	3.46	21.11	2586.95	536.59	130.51	49.70	14.55	3318.10												
OBV	2.32	7.79	2.69	3.82	2.08	18.69	4.46	1.81	3.42	4.60	4.69	18.98	511.88	1069.07	221.17	84.66	9.63	6576.50												
OIL_ROCP	1.03	11.04	24.26	5.69	2.91	44.92	0.56	5.69	15.68	3.53</																				

Table 3. Feature Selection Scores for the Finance Industry

Forecast Type Periodicity No. of Examples	Direction					S					Change					S					Price					S				
	D	W	M	Q	Y	U	D	W	M	Q	Y	U	D	W	M	Q	Y	U	D	W	M	Q	Y	U	W	M	Q	Y	U	
	22932	4753	1092	364	91	M	22932	4753	1092	364	91	M	22932	4753	1092	364	91	M	22932	4753	1092	364	91	M	22932	4753	1092	364	91	
AD	0.23	0.62	0.29	2.62	1.28	5.04	0.18	0.24	0.37	0.07	0.22	1.09	779.54	172.85	38.46	16.86	2.62	1010.31												
ADOSC	24.14	1.85	36.41	6.82	3.50	72.72	73.18	0.08	29.59	2.65	2.80	108.30	4.30	0.88	0.10	0.02	0.15	5.45												
ADX	0.38	0.07	1.10	1.32	0.07	2.95	1.90	2.64	3.60	0.76	1.31	10.20	30.20	6.97	3.96	3.48	0.29	44.90												
ADXR	0.09	0.24	0.88	1.96	0.09	3.28	2.26	4.08	4.47	1.89	1.07	13.77	43.76	10.70	4.73	3.99	0.21	63.39												
APO	0.11	4.57	41.91	32.31	59.84	138.74	0.12	16.46	79.91	52.10	64.91	213.50	489.46	105.99	34.19	19.37	26.03	675.05												
AROONOSC	10.28	5.16	100.55	49.93	66.47	232.40	6.66	6.54	109.18	55.86	50.46	228.70	313.61	76.95	37.46	30.18	51.11	509.32												
ATR	12.31	12.29	18.15	16.78	14.11	73.64	20.33	47.29	59.53	51.87	31.53	211.15	104.06	199.90	31.23	4.47	0.08	1280.72												
Adj Close	0.04	0.02	0.06	1.94	0.79	3.75	2.95	2.31	1.94	4.58	6.81	24.02	7212672.00	537854.48	30418.29	2776.98	126.54	7783848.29												
BITUSD_ROCP	12.01	6.83	5.51	1.22	3.31	28.88	8.59	3.75	2.69	0.04	3.17	18.25	1.82	0.25	0.43	0.59	3.61	6.89												
BOP	52.16	321.25	132.03	41.91	39.06	586.40	37.62	304.10	133.19	28.24	27.99	531.14	45.08	48.47	44.02	32.77	42.64	212.98												
CCI	6.69	94.74	182.46	57.73	6.24	407.86	7.75	9.30	205.00	62.50	53.85	423.18	253.75	67.96	38.74	27.25	40.41	428.11												
CMO	7.95	61.77	110.43	48.60	73.14	301.89	11.64	64.21	118.05	52.69	58.52	305.12	681.56	157.69	51.76	28.39	33.24	952.64												
ChangeLastPeriod	25.90	50.64	66.44	19.54	5.22	167.74	70.96	94.21	130.37	31.14	4.36	330.04	23.61	32.60	35.70	31.27	38.41	161.59												
Close	0.04	0.92	0.06	1.94	0.79	3.75	2.95	2.31	1.94	4.58	6.81	24.02	7212672.00	537854.48	30418.29	2776.98	126.54	7783848.29												
DEMA	0.08	0.14	0.72	2.32	0.82	4.08	4.84	5.23	4.98	5.27	6.91	21.23	2067934.00	271586.28	21772.17	2607.14	125.54	2364025.14												
DJ_ROCP	20.84	20.28	51.05	18.11	24.52	134.81	31.32	28.65	117.78	34.21	19.75	231.70	65.62	14.14	9.33	7.62	2.04	98.75												
DX	0.02	0.65	0.00	0.14	0.25	1.05	1.99	0.63	0.07	0.21	2.02	4.92	4.63	1.47	1.91	0.30	10.22													
EMA	0.01	0.03	1.50	3.56	1.52	6.62	5.08	6.80	7.38	12.76	12.00	2.51	46.20	29.08	8.18	6.52	9.01	0.01	52.79											
EURUSD_ROCP	0.10	0.58	5.84	2.88	8.43	17.85	7.89	11.05	12.76	12.00	2.51	46.20	29.08	8.18	6.52	9.01	0.01	52.79												
FTSE_ROCP	8.03	6.56	46.46	42.46	37.68	141.99	34.56	3.54	81.11	54.33	25.99	199.52	16.73	3.29	6.02	3.34	32.63													
GBPUSD_ROCP	3.61	7.26	2.89	7.36	22.44	43.56	0.00	0.58	12.05	30.61	13.25	56.49	59.15	12.86	8.20	8.89	0.48	89.58												
GOLD_ROCP	7.34	10.51	6.77	3.05	0.06	27.73	0.06	10.33	16.58	0.25	4.10	31.32	0.02	0.18	0.39	0.20	1.27	2.07												
HT_DCPERIOD	0.32	0.02	0.42	0.18	0.80	1.74	0.02	0.15	0.00	0.65	6.89	7.71	1.82	0.75	0.02	2.42	12.18													
HT_DCPHASE	1.27	0.45	13.48	33.94	45.59	94.74	2.51	0.12	11.45	42.10	32.14	88.33	249.06	71.96	35.93	32.66	37.37	426.99												
HT_PHASOR_INPHASE	5.23	1.25	101.88	45.92	66.71	220.9	1.50	0.03	204.71	69.00	74.37	349.60	95.88	33.42	29.00	18.34	27.22	203.86												
HT_PHASOR_QUADRATURE	2.84	5.54	56.73	11.59	0.28	71.99	42.95	1.16	80.05	6.61	0.57	131.35	0.00	0.01	0.27	0.25	3.35	3.89												
HT_SINE_SLEADINE	11.94	0.02	37.73	5.05	13.92	96.66	4.37	0.09	43.39	7.93	19.76	75.55	10.21	3.17	1.68	7.75	18.67	41.48												
HT TRENDLINE	0.03	0.01	1.73	3.43	1.40	6.60	4.75	7.15	8.04	7.43	8.43	35.80	935920.30	152345.16	15790.07	2199.69	111.69	1106365.91												
HT TRENDMODE	0.13	0.31	0.30	1.03	4.55	6.32	0.44	0.00	0.16	0.10	0.84	1.54	128.74	43.06	25.14	20.17	10.55	226.76												
High	0.00	0.62	0.14	2.17	0.91	3.83	7.43	3.18	2.52	5.14	7.22	25.49	626075.00	469501.81	28253.98	2663.58	122.78	6126617.15												
IXIC_ROCP	19.05	14.57	54.60	20.24	25.82	134.28	20.67	23.37	103.86	51.84	18.98	218.73	29.66	5.44	4.32	5.65	0.74	45.81												
KAMA	0.02	0.05	1.67	3.81	1.69	7.24	4.82	6.57	7.29	7.30	8.81	34.00	795678.20	135453.22	15514.13	2116.88	108.06	948860.49												
Low	0.00	0.00	0.04	1.78	0.06	3.31	7.62	2.55	1.64	4.09	6.43	23.44	614838.00	404920.19	31278.60	2867.49	130.12	594763.40												
MA	0.02	0.01	1.96	3.67	1.53	7.18	4.95	7.56	7.93	8.71	8.71	39.71	809382.40	134706.81	14545.95	2115.72	109.41	960914.29												
MFI	13.47	19.49	108.24	41.53	59.07	247.16	8.05	26.80	124.65	57.66	37.37	254.54	224.65	54.90	24.82	18.10	34.83													
MIDPOINT	0.03	0.12	0.99	2.79	1.09	5.02	5.11	5.42	5.78	6.17	7.63	30.12	1760452.00	245458.48	20591.37	2445.42	118.70	2029065.90												
MIDPRICE	0.05	0.14	0.95	2.75	1.07	4.97	4.83	5.11	5.54	6.01	7.60	29.08	127996.00	245970.78	20900.76	2474.62	119.03	199461.20												
MINUS_DI	1.44	28.56	75.40	39.68	57.08	202.09	2.84	30.52	85.52	40.49	34.16	193.53	609.86	128.73	36.81	17.08	16.87	809.36												
MINUS_DM	3.11	18.30	38.06	25.91	23.31	107.69	1.28	56.85	100.44	65.73	40.18	264.48	262.18	48.20	5.73	0.44	0.71	317.27												
MOM	15.94	28.20	157.67	53.13	6.08	321.02	43.15	37.29	342.28	85.50	75.69	583.90	247.88	58.51	30.77	18.33	24.10	375.98												
NATR	1.86	2.87	6.92	1.25	5.40	18.29	2.24	3.90	12.19	1.77	6.69	26.79	3204.05	669.31	180.23	57.99	19.16	4130.74												
OBV	8.89	6.46	1.20	1.45	2.56	20.55	9.90	4.04	1.33	2.52	2.17	19.95	248.46	40.48	9.03	5.46	8.51	311.94												
OIL_ROCP	0.89	4.12	16.15	6.22	6.75	34.13	0.43	5.55	41.31	37.63	1.62	86.54	13.09	4.34	4.35	4.93	0.00	26.72</												

4.2. Classification Results

We attempt to find the optimal combination between the classifier and number of features that are fed to it. We design an experiment whereby we attempt each classifier for 70 times on identical datasets, each iteration increasing the number of features till we reach 70.

Table 4 shows the top 15 results obtained from this test for the monthly dataset, with the best combination being the MLP Classifier with 28 features. It is important to note that the top 15 results, all obtain an accuracy of 69% and over, with a marginal difference. As one may note, these results only feature three classifiers, other classifiers were evaluated in the same test, but failed to reach the top 15 list, hence are omitted from this table.

Table 4. Table showing the top 15 results for optimal classifier and number of features for the Monthly Tech Industry Dataset

Method	NoFeatures	Accuracy
MLPClassifier	28	0.696640
LogisticRegression	21	0.695652
MLPClassifier	20	0.695652
SVC	18	0.693676
SVC	26	0.692688
LogisticRegression	28	0.692688
LogisticRegression	20	0.691700
SVC	16	0.691700
LogisticRegression	24	0.691700
LogisticRegression	22	0.691700
LogisticRegression	23	0.691700
LogisticRegression	30	0.690711
SVC	12	0.690711
SVC	17	0.690711
LogisticRegression	42	0.690711

In order to better visualise this experiment, Figure 1 shows a graph derived from the full results for these three classifiers, illustrating the of Number of Features plotted against the accuracy rate. From this visualisation, we can see that the accuracy peaks when utilising between 20 and 30 features, and starts to overfit beyond 30 features. One can see that accuracy deterioration occurs in different rates for each classifier. In particular, Logistic Regression and SVC maintain a relatively constant accuracy when compared to the Multi Layer Perceptron line. This substantiates the fact that each classifier has its own merits when preventing overfitting of the training data or even dealing with noisy data.

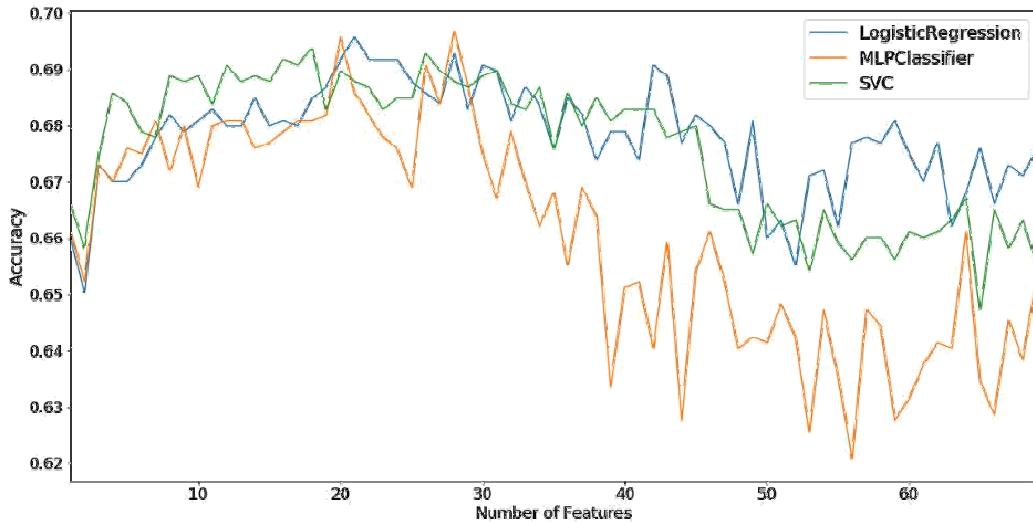


Figure 1. Graph showing accuracy changing according to the number of features passed for each classifier for the Tech industry

Table 5 illustrates the final, optimal combination of classifier and features for each Industry and periodicity, achieving the best results on the yearly datasets for both finance and tech industries.

Table 5. Optimal combination of classifier and features for the Next Period Direction Forecast

Next Period Direction				
Industry	Periodicity	Method	Features	Accuracy
Finance	Daily	<i>AdaBoostClassifier</i>	8	0.52
	Weekly	<i>LogisticRegression</i>	59	0.65
	Monthly	<i>QuadraticDiscriminantAnalysis</i>	2	0.71
	Quarterly	<i>GaussianNB</i>	21	0.64
	Yearly	<i>LogisticRegression</i>	23	0.81
Tech	Daily	<i>LogisticRegression</i>	10	0.52
	Weekly	<i>LogisticRegression</i>	34	0.67
	Monthly	<i>MLPClassifier</i>	28	0.69
	Quarterly	<i>AdaBoostClassifier</i>	3	0.69
	Yearly	<i>KNeighborsClassifier</i>	5	0.77

4.3. Regression Results

We now attempt to achieve the best possible results for each combination by iterating over each regressor, training and testing each one for 70 times, increasing the features that are passed to the model with each iteration. This is similar to what was done in the classification test, with the difference that we are evaluating regressors, and their success is calculated by means of error (RMSE), which is at its optimum when it is lowest.

Table 6 tabulates the optimal combination of algorithm and number of features. One can note that in contrast to our classification results where we had 7 different algorithms that feature in the optimal results, here we have the dominance of Linear Regression and Support Vector Regression, and one instance where KNN proved best.

Table 6. Table showing the optimal regressors for the Next period Actual Price Forecast

Next Period Actual Price				
Industry	Periodicity	Method	Features	RMSE
Finance	Daily	<i>LinearRegression</i>	1	0.011706
	Weekly	<i>LinearRegression</i>	68	0.02047
	Monthly	<i>LinearRegression</i>	8	0.043636
	Quarterly	<i>KNeighborsRegressor</i>	15	0.100451
	Yearly	<i>LinearRegression</i>	26	0.161316
Tech	Daily	<i>LinearRegression</i>	3	0.012995
	Weekly	<i>LinearRegression</i>	67	0.020943
	Monthly	<i>LinearRegression</i>	59	0.04753
	Quarterly	SVR	49	0.086884
	Yearly	SVR	61	0.165858

With regards to the number of features, we can note that the daily forecast required least features to reach optimal accuracy. Moving down the table, we can see that the optimal number of features unlike the RMSE, is not in any way linear to the periodicity.

Table 7 tabulates the optimal combination of algorithm and number of features for our second regression problem, that of forecasting the change in price rather than the price itself. Whilst the error rates are noticeably higher than those for predicting the actual price itself, they follow a similar pattern when it comes to the number of features. It also follows a similar pattern when comparing the RMSE to the periodicity, where a linear relationship seems to exist. Errors from the Finance and Tech industry are relatively similar with the exception of the Monthly and Yearly datasets, where Finance industry stocks achieve considerably better results on the quarterly and yearly datasets.

Table 7. Table showing the optimal regressors for the Next period Change in Price Forecast

Next Period Change in Price				
Industry	Periodicity	Method	Features	RMSE
Finance	Daily	<i>LinearRegression</i>	1	0.061304
	Weekly	<i>LinearRegression</i>	38	0.097223
	Monthly	<i>LinearRegression</i>	20	0.132902
	Quarterly	SVR	29	0.188626
	Yearly	<i>LinearRegression</i>	54	0.157614
Tech	Daily	SVR	1	0.063599
	Weekly	<i>LinearRegression</i>	69	0.104637
	Monthly	<i>LinearRegression</i>	33	0.152169
	Quarterly	SVR	27	0.191362
	Yearly	SVR	66	0.217411

In order to better visualise the forecasts and errors done by the regressor algorithms, Figure 2 shows a plot of the normalised close price and the next day predicted close price against the test period (2014-2016). Figure 3 illustrates the same visualisation but for the weekly tech dataset, here we can see that the plot lines are not as much in line with each other as the daily dataset, this represents the increase in error between one period and another which was discussed earlier. Figure 4 illustrates the same representation, however this time we are dealing with the Monthly dataset, the increase in error here is more prevalent than those discussed up till now.

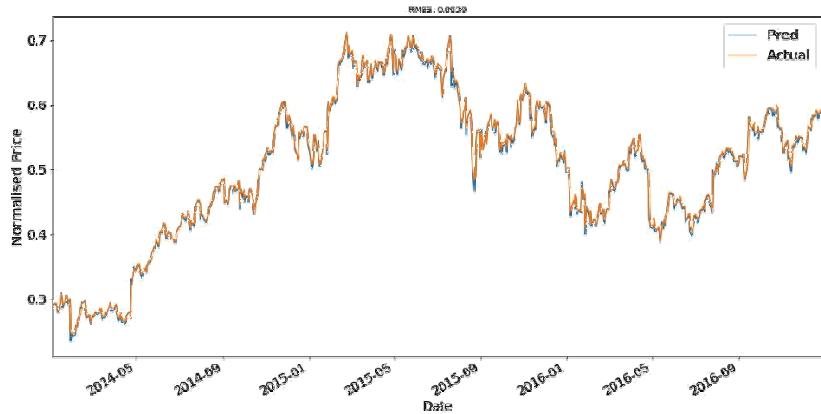


Figure 2. Next Day Actual Price forecast for AAPL Stock using the optimal regressor

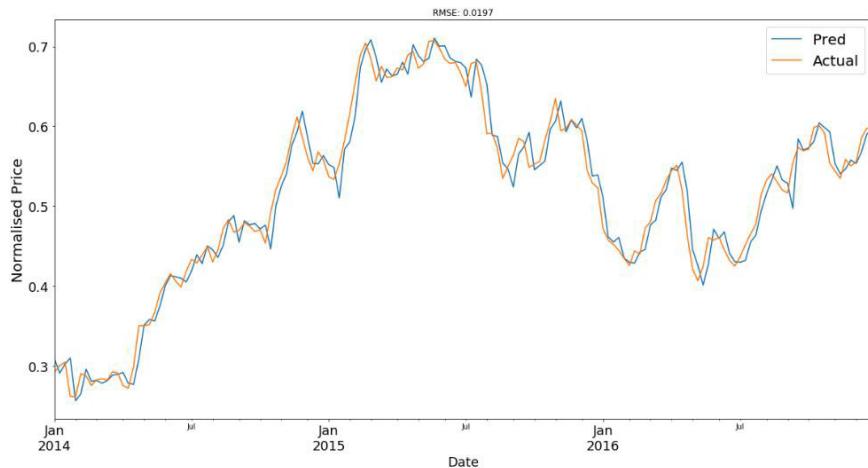


Figure 3. Next Week Actual Price forecast for AAPL Stock using the optimal regressor

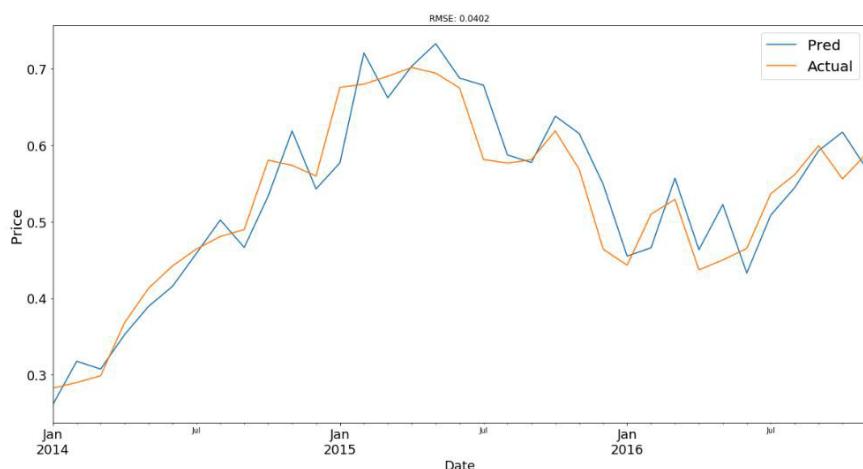


Figure 4. Next Month Actual Price forecast for AAPL Stock using the optimal regressor

5. CONCLUSION AND FUTURE WORKS

This paper gives a basic introduction to the stock market and various approaches taken to predict price movements. Data sources are identified and features are formulated from this data. Through experimentation we achieve 81% accuracy for future trend direction using classification, 0.0117 RMSE for next day price and 0.0613 RMSE for next day change in price using regression techniques. Although many consider movements in the stock market to be somewhat random, our results prove that there are elements that can be learnt, after all the job of a stockbroker or analyst is far from random. By means of identifying relevant features for each stock, a trend can be detected. One of Spencer's Laws of Data holds that anyone can make a decision given enough facts. Whilst extra-ordinary events such as natural disasters are difficult to predict, certain external data such as common market opinions can be gauged through various sentiment analysis techniques. Building on the works discussed in this paper, our approach will be enhanced using deep learning techniques such as Recurrent Neural Networks and Convolutional Neural Networks.

REFERENCES

- [1] Sykes A. O., "An Introduction to Regression Analysis", The Inaugural Coase Lecture, 1993
- [2] Yue Xu S., "Stock Price Forecasting Using Information from Yahoo Finance and Google Trend", UC Berkeley, 2012
- [3] Duke, "Stationarity and differencing", [Online], Available: <http://people.duke.edu/~rnau/411diff.htm> [Accessed January 2017]
- [4] Vapnik V. N., "An Overview of Statistical Learning Theory" *IEEE TRANSACTIONS ON NEURAL NETWORKS*, 1999.
- [5] Kim K. "Financial time series forecasting using support vector machines", Department of Information Systems, Dongguk University 2003
- [6] Panigrahi S. S. and Mantri J. K., "A text based Decision Tree model for stock market forecasting," Green Computing and Internet of Things (ICGCIoT), 2015 International Conference on, Noida, 2015
- [7] G. Iuhasz, M. Tirea and V. Negru, "Neural Network Predictions of Stock Price Fluctuations," Symbolic and Numeric Algorithms for Scientific Computing (SYNASC), 2012 14th International Symposium on, Timisoara, 2012
- [8] Siripurapu A., "Convolutional Networks for Stock Trading", Stanford University, Department of Computer Science, 2014
- [9] Qiu M. and Song Y. "Predicting the Direction of Stock Market Index Movement Using an Optimized Artificial Neural Network", Department of Systems Management, Fukuoka Institute of Technology, Fukuoka, Japan, 2016
- [10] S. and Szegedy C., "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", 2015

Authors

Prof. Alexiee Dingli is the Head of the Department of Artificial Intelligence within the Faculty of ICT at the University of Malta. He is also one of the founder members of the ACM student chapter in Malta, founder member of the Web Science Research, founder member of the International Game Developers Association (IGDA) Malta and of the Gaming group at the same University. He pursued his Ph.D. on the Semantic Web at the University of Sheffield in the UK under the supervision of Professor Yorick Wilks.



Karl Sant Fournier is a Bachelor's degree graduate in Business and Computing and is currently a part time student reading for a Masters degree with the department of Artificial Intelligence. He works full time as a Business Intelligence Developer at a local Bank. Given his knowledge of the financial industry he is focusing his thesis on Financial Time Series forecasting, using Artificial Intelligence techniques

