

# Data Driven Automated Algorithmic Trading

Gabriel Gauci Maistre

**Abstract** — Various existing stock market price forecasting methods were analysed in this report. Three methods were applied towards the problem making use of Technical Analysis, these were Time Series Analysis, Machine Learning, and Bayesian Statistics. Through the results of this report, it was found that the Efficient Market Hypothesis remains true, that past data does not contain enough useful information to forecast future prices and gain an advantage over the market. However, the results proved that Technical Analysis and Machine Learning could still be used to guide an investors decision. It was also found that the Random Walk Hypothesis was not necessarily true, as some stocks showed signs of auto and partial correlation. A common application of technical analysis was demonstrated and shown to produce limited useful information in beating the market. Based on the findings, a number of automated trading algorithms were developed using machine learning and backtested to determine their effectiveness.

**Keywords** — *machine learning, time series analysis, probabilistic, bayesian, statistics, inference, algorithmic trading*

## 1. Introduction

The stock market retains its status as a prime location for investors to invest in the market and earn a profit, however this is not always easy due to the constantly thriving and changing nature which follows the stock market. Investors are constantly presented with numerous profit potential opportunities, however without intensive planning and analysis, these opportunities could easily turn into losses. This means that it is crucial for every investor to carry out stock market analysis prior to any investment by monitoring past price movements in order to forecast future trends. Even though past data is not a clear indication of future movement, it is still proven to provide some useful insight.

In this paper, we investigate the Efficient Market Hypotheses (EMH), and the Random Walk Hypothesis (RWH) in the context of stock market forecasting and trading. For the purpose of benchmarking the performance of the algorithms, a total of five stocks were randomly selected from a basket of uncorrelated stocks were selected. These were MSFT, CDE, NAVB, HRG, and HL. Our key findings are as follows:

1. We found strong auto and partial correlation in the stocks studied, disproving the RWH.
2. We found the EMH to be true, as although the automated trading algorithm made a large profit, it was not sufficient to beat the returns of the market.
3. Time series analysis was found to be a weak factor in forecasting financial stocks.

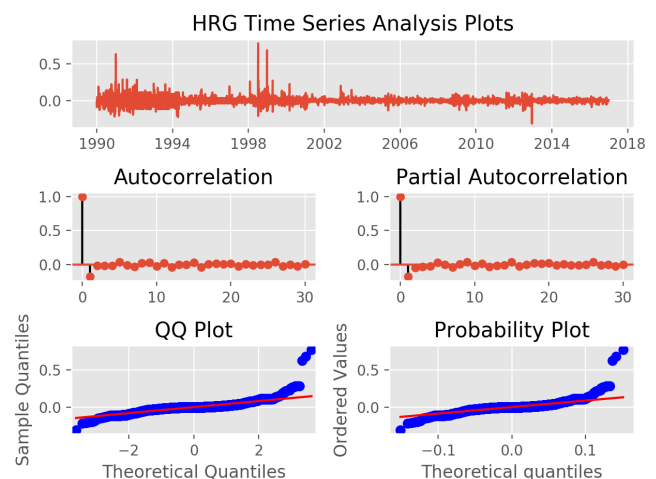
4. Bayesian statistics was a strong contender as a forecasting method.
5. Machine learning was found to be the best method of forecasting using both classification and regression methods.
6. We found regression methods of machine learning to fair better than classification methods when implemented and backtested.

## 2. Experiments

### 2.1. Time Series Analysis

The selected stocks were extracted from the data set and stored in a data frame. The log returns of the stocks were calculated by calculating the logarithm of the stock's adjusted close price divided by the following day's adjusted close price. The resulting values from the calculation were then stored in a new column in the data frame and all infinite values were dropped from the series.

The random walk theory suggests that stock price changes have the same distribution and are independent of each other, so the past movement or trend of a stock price or market cannot be used to predict its future movement. In short, this is the idea that stocks take a random and unpredictable path. As can be seen in figure 1, the correlation plots show this theory to be false and that past movement is related to future movement.



**Fig. 1.** HRG time series analysis

ARMA fared the best than the rest of the algorithms in predicting stock price returns, having the lowest difference in sharpe ratios based on the original price returns and in-sample predicted price returns, however still failed to

achieve a good fit in the in-sample tests. The series was fitted to an ARMA(p, q, r) model with an order selected based on the lowest AIC. No constants were passed to the ARMA model. The exact loglikelihood for the fit of the ARMA model was maximized via the Kalman Filter. As can be seen in the time series analysis plots, ARMA showed to have very heavy tails in the QQ and probability plots. The algorithm predicted an abnormal sharp decline at the start of the forecast and showed diminishing returns over time, as was evident in figure 2.

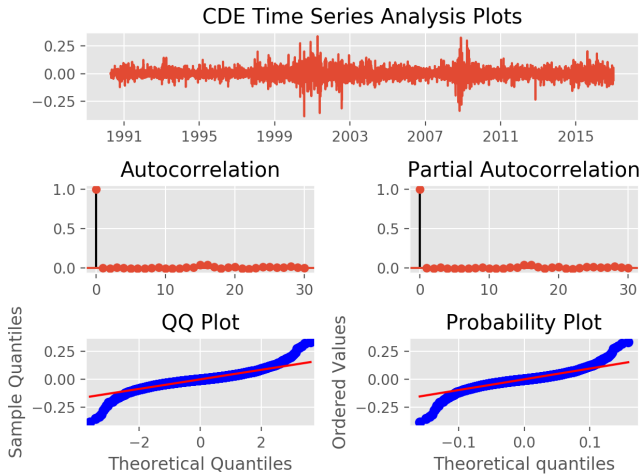


Fig. 2. CDE ARMA time series analysis

## 2.2. Machine Learning

The random forest was fit with bootstrap samples when building the trees, while all the weights associated were set to 1. The 'gini' function to measure the quality of a split, and no maximum depth of the tree was set, allowing the nodes to expand until all leaves are pure or until all leaves contain less than the minimum split samples. The number of features to consider when looking for the best split was the square root of the number of passed, and no limit on the maximum leaf nodes for growing trees was set. A threshold of  $1e-7$  was used to terminate the tree growth to determine if a node is a leaf, if the impurity of a node is below the threshold, the node is a leaf. The minimum number of samples required to be at a leaf node was set to 1, and the minimum number of samples required to split an internal node was set to 2. The minimum weighted fraction of the sum total of weights (of all the input samples) required to be at a leaf node was set to 0, and the number of trees in the forest was set to 10. The number of jobs to use for the computation was set to 1, making use of only 1 CPU core, and out-of-bag samples to estimate the generalization accuracy were not used. The verbosity of the tree building process was not controlled, and the random number generator used by the model was that of Numpy's RandomState. The model was built using a cold start by not making use of the previous call to fit and add more estimators to the ensemble, meaning a whole new forest was fit instead.

Ticker	Precision	True -	False -	True +	False +
MSFT	0.77	493	131	547	190
CDE	0.79	566	145	503	132
NAVJ	0.76	590	164	332	128
HRG	0.75	554	210	462	134
HL	0.81	581	101	509	169

Table 1  
Random Forest classification results

## 2.3. Bayesian Statistics

The last 500 rows of the selected stocks were extracted from the data set and stored into a data frame. The log returns were calculated by dividing each days adjusted close with the adjusted close of the following day, in logarithmic form. The resulting values from the said calculation were then stored in a new column in the data frame and all infinite values were dropped from the series. The sharpe ratio of the original and predicted price returns was calculated to serve as an accuracy score in the in-sample tests. The model returns were modeled with a Student-t distribution with an unknown degrees of freedom paramater, and a scale paramater determined by a latent process. The algorithm achieved a good fit in the in-sample tests, having very similar sharpe ratios based on the original price returns and in-sample predicted price returns.

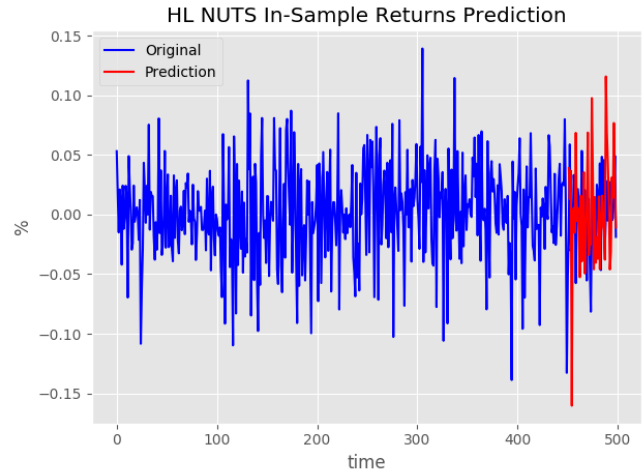


Fig. 3. HL NUTS in-sample returns prediction

## 2.4. Strategy

We evaluated both classification and regression methods for algorithmic trading, in which it was evident that regression methods were superior to those of classification. An ordered dictionary was used to store the day's close price for each stock as the backtester simulated each trading day. The algorithm was only allowed to run once there was enough price data, the amount chosen was 50. The 15 and 50 day SMAs were calculated and each stored in an array. The two arrays were converted into an array of tuples, where the i-th tu-

ple contains the  $i$ -th element from each of the argument sequences or iterables. The independent variables, the SMAs, and the dependant variables, the close prices for the particular stock, were passed to the machine learning algorithm to predict the next day's close price. An order was placed on a stock if the predicted price was higher the current day's price. A stop loss of 80% was placed on each position. When tasked with predicting rises in stock prices, the algorithm did fairly well, marking a slightly lower total return of 43%. The algorithm also underperformed immensely when tasked at also predicting stock price falls, marking a negative total return of -80.4%. To compensate for this, a stop loss was added to sell all positions if the stock in question falls below -20%, this resulted in a profit of 83.7%.

Starting Capital	\$100,000
Total Capital Used	\$229,547.84
Sharpe Ratio	0.420
Portfolio Value	\$183,714.616
Algorithm Period Return	0.837
Benchmark Period Return	1.008
Algorithm Volatility	0.373
Benchmark Volatility	0.156

Table 2  
Machine Learning Regression strategy with stop loss

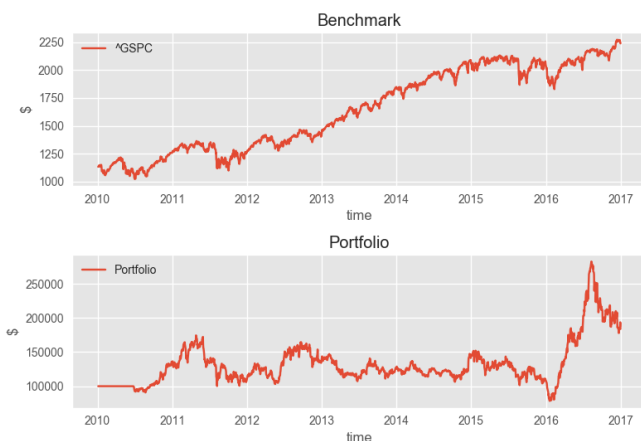


Fig. 4. Machine Learning Regression strategy with stop loss

### 3. Related work

There is a growing demand for forecasting interest rates, as financial researchers, economists, and players in the fixed income markets seek to find the best method to get ahead of the market. A study was carried out to develop an appropriate model for forecasting the short-term interest rates, implicit yield on 91 day treasury bill, overnight MIBOR rate, and call money rate. The short-term interest rates are forecasted using univariate models such as the Random Walk, ARIMA, ARMA-GARCH, and ARMA-EGARCH. The appropriate model for forecasting is determined considering a

six-year period from 1999. Radha et al. showed evidence that GARCH models are best suited for forecasting when applied towards time series having volatility clustering effects. It was their firm belief that ARIMA-EGARCH is the most appropriate forecasting model for these circumstances.

Darrat et al. set out to investigate with the use of new daily data, whether prices in the two Chinese stock exchanges (Shanghai and Shenzhen) follow a random walk process as required by market efficiency. Two different approaches were applied, the standard variance-ratio test, and a model-comparison test that compares the ex post forecasts from a naive model with those obtained from several alternative models such as ARIMA, GARCH, and ANNs. To evaluate ex post forecasts, Darrat et al. made use of several procedures including root-mean-square error (RMSE), mean absolute error (MAE), uncertainty coefficient, and encompassing tests. It was concluded that the model-comparison approach yielded results which were quite strongly rejected the RWH in both Chinese stock markets when compared with the variance-ratio test. Darret et al. recommended the use of ANNs, as their results showed strong support for the model as a potentially useful factor for forecasting stock prices in emerging markets.

A vast majority of academics tend to predict the price of stocks in financial markets, however most models used are flawed and only focus on the accurate forecasting of the levels of the underlying stock index. There is a lack of studies examining the predictability of the direction of stock index movement. Given the notion that a prediction with little forecast error does not necessarily translate into capital gain, the authors of this research attempt to predict the direction of the S&P CNX NIFTY Market Index of the National Stock Exchange, one of the fastest growing financial exchanges in developing Asian countries. Machine learning models such as random forest and SVMs, differ widely from other models, and are making strides in predicting the financial markets. Kumar et al. tested classification models to predict the direction of the markets, by applying models such as linear discriminant analysis, logistic regression, ANNs, random forest, and SVM. Their evidence shows that SVMs outperform the other classification methods in terms of predicting the direction of the stock market movement, and that the random forest model outperforms other models such as ANNs, discriminant analysis, and logistic regression.

Creamer et al. developed an automated trading algorithm making use of multiple stocks relying on a layered structure consisting of a machine learning algorithm, an online learning utility, and a risk management overlay. The machine learning algorithm which they made use of was an Alternating Decision Tree (ADT) implemented with Logitboost. Their algorithm was able to select the best combination of rules derived from well known technical analysis indicators, and the best parameters of the indicators in question. Additionally, their online learning layer was also able to combine the output of several ADTs, suggesting a short or long position. Finally, the risk management layer

in which they implemented, was able to validate the trading signal once it exceeds a specified non-zero threshold and limit the application of their trading strategy when it is not profitable. They tested the expert weighting algorithm with data of 100 randomly selected companies of the S&P 500 index during the period 2003-2005. They found that their algorithm generated abnormal returns during the test period. Their experiments show that the boosting approach was able to improve the predictive capacity when indicators were combined and aggregated as a single predictor. Furthermore, their results indicated that the combination of indicators of different stocks was adequate in order to reduce the use of computational resources, while still maintaining an adequate predictive capacity.

Hoffman et al. introduce the No-U-Turn-Sampler, an extension to the Hamiltonian Monte Carlo (HMC), which is an MCMC algorithm that avoids the random walk behaviour and sensitivity to correlated parameters that plague many MCMC methods by taking a series of steps informed by first-order gradient information. In their paper, they claim HMCs performance to be highly sensitive to two user-specified parameters: a step size, and a desired number of steps. NUTS is an improvement from HMC as it eliminates the need to set a number of steps, and works by building a set of likely candidate points that spans a wide swath of the target distribution, stopping automatically when it starts to double back and retrace its steps. They achieved all of this by making use of a recursive algorithm.

## 4. Conclusion

Three financial forecasting methods were presented in this report, two of which showed little to no potential of ever producing any statistically significant result when the correct methodology was applied. The third method, machine learning, showed some potential in the tests carried out, which is why this method was built into an automated algorithmic strategy to trade with. The algorithm proved to be successful in forecasting future prices, using both classification and regression methods. However, the backtesting proved this method to fail in forecasting price falls. Once this factor was removed from the equation, the algorithms were very successful and reported a profit by the end of the test. This is however not always ideal as stocks which could fall in price could be catastrophic to the strategy. A stop loss would be ideal in insuring that no positions are held in downward falling stocks. It was also evident that regression methods were more successful in forecasting future price movements when compared to classification methods.

If there is anything that this report shows, is that profitable stock market prediction is an extremely tough problem. Even though the strategies reported a profit by the end of the backtest, they still did not beat the market. Whether it is at all possible to use such methods to outperform the market returns, ultimately remains an open question. These findings support the Efficient Market Hypothesis, proving that casual investors are better off investing in passive buy and

hold strategies consisting of index funds and ETFs. However, there was some evidence found showing that the Random Walk Hypothesis does not hold true for all cases, as some stocks did show signs of repeating trends.

## Acknowledgements

The author would like to express his special thanks of gratitude to my supervisor, Alan Gatt, for the patient guidance, encouragement, and advice he has provided throughout my time as his student. He would also like to thank Luke Vella Critien, for guiding him towards the right path in the early stages of my research and for recommending Alan Gatt as my tutor. The author's gratitude is also extended to Emma Galea and Miguel Attard for their valuable input while carrying out his research. Completing this work would have been all the more difficult were it not for the support and friendship provided by the other members of the Malta College of Arts, Sciences, and Technology, and the institute of Information and Technology. The author is indebted to them for their help. Finally, the author wishes to thank his family who have supported me all throughout the final year of his bachelors.

## References

- [1] K. R. Wilson and V. V. Yakovlev, „Ultrafast rainbow: tunable ultrashort pulses from a solid-state kilohertz system”, *J. Opt. Soc. Am. B*, vol. 14, pp. 444–448, 1997.
- [2] J. Comly and E. Garmire, „Second harmonic generation from short pulses”, *Appl. Phys. Lett.*, vol. 12, no. 7-9, 1968.
- [3] O. E. Martinez, „Achromatic phase matching for second harmonic generation of femtosecond pulses”, *IEEE J. Quant. Electron.*, vol. QE-25, pp. 2464–2468, 1989.
- [4] G. Szabo and Z. Bor, „Broadband frequency doubler for femtosecond pulses”, *Appl. Phys. B*, vol. 50, pp. 51–54, 1990.
- [5] J.-Y. Zhang, J. Y. Huang, H. Wang, K. S. Wong, and G.K. Wong, „Second-harmonic generation from regeneratively amplified femtosecond laser pulses in BBO and LBO crystals”, *J. Opt. Soc. Am. B*, vol. 15, pp. 200–209, 1998.
- [6] K. Hayata and M. Koshihara, „Group-velocity-matched second-harmonic generation: an efficient scheme for femtosecond ultraviolet pulse generation in periodically domain-inverted  $\beta$ -BaB<sub>2</sub>O<sub>4</sub>”, *Appl. Phys. Lett.*, vol. 62, pp. 2188–2190, 1993.
- [7] G. Y. Wang and E. M. Garmire, „High-efficiency generation of ultrashort second-harmonic pulses based on the Cherenkov geometry”, *Opt. Lett.*, vol. 19, pp. 254–256, 1994.
- [8] C. Radzewicz, Y. B. Band, G. W. Pearson, and J. S. Krasinski, „Short pulse nonlinear frequency conversion without group-velocity-mismatch broadening”, *Opt. Commun.*, vol. 117, pp. 295–303, 1995.
- [9] P. Di Trapani, A. Andreoni, G. P. Banfi, C. Solcia, R. Danielius, A. Piskarskas, P. Foggi, M. Monguzzi, and C. Sozzi, „Group-velocity self-matching of femtosecond pulses in noncollinear parametric generation”, *Phys. Rev. A*, vol. 51, pp. 3164–3168, 1995.
- [10] V. Krylov, A. Kalintsev, A. Rebane, D. Erni, and U. P. Wild, „Non-collinear parametric generation in LiIO<sub>3</sub> and  $\beta$ -barium borate by frequency-doubled femtosecond Ti:sapphire laser pulses”, *Opt. Lett.*, vol. 20, pp. 151–153, 1995.
- [11] P. Di Trapani, A. Andreoni, C. Solcia, P. Foggi, R. Danielius, A. Dubietis, and A. Piskarskas, „Matching of group velocities in three-wave parametric interaction with fs pulses and application to travelling-wave generators”, *J. Opt. Soc. Am. B*, vol. 12, pp. 2237–2244, 1995.

- [12] P. Di Trapani, A. Andreoni, P. Foggi, C. Solcia, R. Danielius, and A. Piskarskas, „Efficient conversion of femtosecond blue pulses by travelling-wave parametric generation in non-collinear phase matching”, *Opt. Commun.*, vol. 119, pp. 327–332, 1995.
- [13] T. R. Zhang, H. R. Choo, and M. C. Downer, „Phase and group velocity matching for second harmonic generation of femtosecond pulses”, *Appl. Opt.*, vol. 29, pp. 3927–3933, 1990.
- [14] A. Andreoni and M. Bondani, „Group-velocity control in the mixing of three non-collinear phase-matched waves”, *Appl. Opt.*, vol. 37, pp. 2414–2423, 1998.
- [15] P. Di Trapani, A. Andreoni, C. Solcia, G. P. Banfi, R. Danielius, A. Piskarskas, and P. Foggi, „Powerful sub-100-fs pulses broadly tunable in the visible from a blue-pumped parametric generator and amplifier”, *J. Opt. Soc. Am. B*, vol. 14, pp. 1245–1248, 1997.
- [16] R. Danielius, A. Piskarskas, P. Di Trapani, A. Andreoni, C. Solcia, and P. Foggi, „Matching of group velocities by spatial walk-off in collinear three-wave interaction with tilted pulses”, *Opt. Lett.*, vol. 21, pp. 973–975, 1996.
- [17] R. Danielius, A. Piskarskas, A. Stabinis, G. P. Banfi, P. Di Trapani, and R. Righini, „Travelling-wave parametric generation of widely tunable, highly coherent femtosecond light pulses”, *J. Opt. Soc. Am. B*, vol. 10, pp. 2222–2232, 1993.
- [18] S. A. Akhmanov, A. S. Chirkin, K. N. Drabovich, A. I. Kovrigin, R. V. Khokhlov, and A. P. Sukhorukov, „Nonstationary nonlinear optical effects and ultrashort light pulse formation”, *IEEE J. Quant. Electron.*, vol. QE-4, pp. 598–605, 1968.
- [19] R. Danielius, A. Piskarskas, P. Di Trapani, A. Andreoni, C. Solcia, and P. Foggi, „A collinearly phase-matched parametric generator/amplifier of visible femtosecond pulses”, *IEEE J. Quant. Electron.*, vol. 34, pp. 459–464, 1998.
- [20] S. Sartania, Z. Cheng, M. Lenzner, G. Tempea, Ch. Spielmann, F. Krausz, and K. Ferencz, „Generation of 0.1-TW 5-fs optical pulses at a 1-kHz repetition rate”, *Opt. Lett.*, vol. 22, pp. 1562–1564, 1997.
- [21] M. Nisoli, S. De Silvestri, O. Svelto, R. Szipoecs, K. Ferencz, Ch. Spielmann, S. Sartania, and F. Krausz, „Compression of high-energy laser pulses below 5 fs”, *Opt. Lett.*, vol. 22, pp. 522–524, 1997.
- [22] S. A. Akhmanov, V. A. Vysloukh, and A. S. Chirkin, *Optics of femtosecond laser pulses*. New York: American Institute of Physics, 1992.
- [23] A. Andreoni, M. Bondani, and M. A. C. Potenza, „Ultra-broadband and chirp-free frequency doubling in  $\beta$ -barium borate”, *Opt. Commun.*, vol. 154, pp. 376–382, 1998.
- [24] H. Wang, K. S. Wong, D. Deng, Z. Xu, G. K. L. Wong, and J. Zhang, „Kilohertz femtosecond UV-pumped visible  $\beta$ -barium borate and lithium triborate optical parametric generator and amplifier”, *Appl. Opt.*, vol. 36, pp. 1889–1893, 1997.