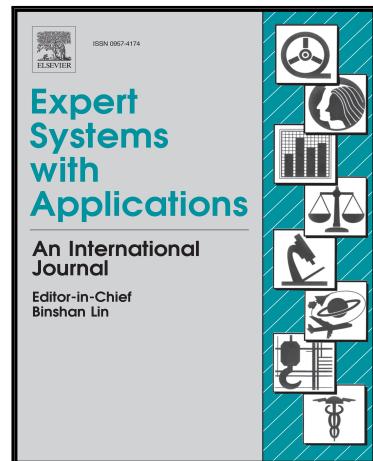


Accepted Manuscript

Stock Market One-Day Ahead Movement Prediction Using Disparate Data Sources

Bin Weng, Mohamed A. Ahmed, Fadel M. Megahed

PII: S0957-4174(17)30133-1
DOI: [10.1016/j.eswa.2017.02.041](https://doi.org/10.1016/j.eswa.2017.02.041)
Reference: ESWA 11149



To appear in: *Expert Systems With Applications*

Received date: 24 March 2016
Revised date: 25 February 2017
Accepted date: 26 February 2017

Please cite this article as: Bin Weng, Mohamed A. Ahmed, Fadel M. Megahed, Stock Market One-Day Ahead Movement Prediction Using Disparate Data Sources, *Expert Systems With Applications* (2017), doi: [10.1016/j.eswa.2017.02.041](https://doi.org/10.1016/j.eswa.2017.02.041)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Highlights

- A financial expert system for predicting the daily stock movements.
- “Knowledge base” captures both traditional and online data sources.
- The “inference engine” uses three artificial intelligence techniques.
- Prediction accuracy of 85% is higher than the reported results in the literature.
- The system is hosted online and freely available for investors and researchers.

Stock Market One-Day Ahead Movement Prediction Using Disparate Data Sources

Bin Weng^a, Mohamed A. Ahmed^b, Fadel M. Megahed^c

^a*Department of Industrial and Systems Engineering, Auburn University, AL 36849, USA — Email:
bzw0018@auburn.edu*

^b*Department of Industrial and Systems Engineering, Auburn University, AL 36849, USA — Email:
mza0068@auburn.edu*

^c*Corresponding Author. Farmer School of Business, Miami University, OH, 45056, USA — Email:
fmegahed@miamioh.edu*

Abstract

There are several commercial financial expert systems that can be used for trading on the stock exchange. However, their predictions are somewhat limited since they primarily rely on time-series analysis of the market. With the rise of the Internet, new forms of collective intelligence (e.g. Google and Wikipedia) have emerged, representing a new generation of “crowd-sourced” knowledge bases. They collate information on publicly traded companies, while capturing web traffic statistics that reflect the public’s collective interest. Google and Wikipedia have become important “knowledge bases” for investors. In this research, we hypothesize that combining disparate online data sources with traditional time-series and technical indicators for a stock can provide a more effective and intelligent daily trading expert system. Three machine learning models, decision trees, neural networks and support vector machines, serve as the basis for our “inference engine”. To evaluate the performance of our expert system, we present a case study based on the AAPL (Apple NASDAQ) stock. Our expert system had an 85% accuracy in predicting the next-day AAPL stock movement, which outperforms the reported rates in the literature. Our results suggest that: (a) the knowledge base of financial expert systems can benefit from data captured from nontraditional “experts” like Google and Wikipedia; (b) diversifying the knowledge base by combining data from disparate sources can help improve the performance of financial expert systems; and (c) the use of simple machine learning models for inference and rule generation is appropriate with our rich knowledge database. Finally, an intelligent decision making tool is provided to assist investors in making trading decisions on any stock, commodity or index.

Keywords: Artificial Intelligence, Feature Selection, Financial Expert System, Google News, R Programming, Wikipedia

¹ 1. INTRODUCTION

² Stock market prediction has attracted much attention from both academia and business. The
³ question remains: “To what extent can the past history of a common stock’s price be used to make
⁴ meaningful predictions concerning the future price of the stock?” Fama (1965). Early research on
⁵ stock market prediction was based on the Efficient Market Hypothesis (EMH) Fama (1965) and the
⁶ random walk theory Cootner (1964); Fama et al. (1969); Fama (1991, 1995). These early models
⁷ suggested that stock prices cannot be predicted since they are driven by new information (news)
⁸ rather than present/past prices. Thus, stock market prices will follow a random walk and their
⁹ prediction accuracy cannot exceed 50% (Bollen et al. (2011)).

¹⁰ There has been an increasing number of studies (see e.g.,Malkiel (2003); Smith (2003); Nofsinger
¹¹ (2005); Prechter Jr and Parker (2007); Bollen et al. (2011)) that provide evidence contrary to what is
¹² suggested by the EMH and random walk hypotheses. These studies show that the stock market can
¹³ be predicted to some degree and therefore, questioning the EMH’s underlying assumptions. Many
¹⁴ within the business community also view Warren Buffet’s ability to consistently beat the S&P index
¹⁵ (Kiersz (2015); Loomis (2012)) as a practical indicator that the market can be predicted.

¹⁶ The significant stock market movements (i.e. spikes) over short horizons cannot also be explained
¹⁷ by the EMH. These spikes are often driven by investor perceptions of a certain stock based on
¹⁸ information (news) collected from disparate data sources. As an illustration, on April 23, 2013,
¹⁹ at 1:07 p.m., Eastern Time, a tweet from the *Associated Press* (AP) account stated “Breaking:
²⁰ Two Explosions in the White House and Barack Obama is injured” (Megahed and Jones-Farmer
²¹ (2015)). The fraudulent tweet, which originated from the hacked AP Twitter account led to an
²² immediate drop in the *Dow Jones Industrial Average* (DJIA). Although the DJIA quickly recovered
²³ following an AP retraction and a White House press release, this example illustrates the immediate
²⁴ and dramatic effects of perception/news on stock prices.

²⁵ While the news may be unpredictable, some recent literature suggests that early indicators
²⁶ can be extracted from online sources (e.g., Google Trends and blogs) to predict changes in various
²⁷ economic indicators. For example, Google search queries have been shown to provide early indicators
²⁸ of disease infection rates and consumer spending (Choi and Varian (2012)). Schumaker and Chen
²⁹ (2009) showed that breaking financial news can be used to predict stock market movements. Bollen
³⁰ et al. (2011) used measurements of collective mood states derived from large-scale Twitter feeds to
³¹ predict the daily up and down changes in the DJIA. In addition, Moat et al. (2013) observed that

32 the frequency of views of Wikipedia’s financially-related pages can be an early indicator of stock
 33 market moves. The authors hypothesized that investors may be using such pages as a part of their
 34 decision making process. This work was extended in Preis et al. (2013) to include data from the
 35 number of relevant searches from Google Trends, and model the effect of search volume on trading
 36 behavior. Note that Mao et al. (2011) indicated that *search* and *usage* are more predictive than
 37 *survey sentiment indicators*.

38 From an expert systems perspective, the stock market prediction problem can be divided into
 39 two components: (1) what information and predictors need to be tracked as a part of our “knowledge
 40 base”; and (2) what artificial intelligence (AI) algorithms can be used for effective rule generation
 41 and predictions. The literature discussed in the previous paragraph indicate that online sources
 42 that capture the “collective intelligence” of investors should be an integral component of a financial
 43 expert system’s knowledge base. It is important to note that these online sources are not typically
 44 used in financial expert systems. Instead, the knowledge base of such systems typically rely on
 45 the historical prices of a stock and/or technical indicators extracted from a time-series analysis of
 46 stock prices (Kimoto et al. (1990); Lee and Jo (1999); Kim and Han (2000); Kim (2003); Hassan
 47 et al. (2007); Qian and Rasheed (2007); Lin et al. (2011); Chourmouziadis and Chatzoglou (2016)).
 48 We hypothesize that combining the expert’s knowledge from online sources with features extracted
 49 from the price and technical indicators will offer a more accurate representation of the dynamics
 50 that affect a stock’s price and its movement. Since these data sources were never combined in the
 51 context of financial expert systems, it is important to examine which AI algorithms are the most
 52 effective in translating the knowledge base into accurate predictions. Table 1 categorizes financial
 53 expert systems used for stock movement prediction based on their “knowledge base” and the AI
 54 approach used. From Table 1, it is clear that the all those papers relied on a single source for the
 55 knowledge base. The reader should note that there is a limited number of expert systems (e.g.,
 56 Bollen et al. (2011)) that combined traditional sources with crowd-sourced experts’ data; however,
 57 they are not included in our table since they predicted price (i.e. a continuous outcome instead
 58 of our binary outcome). The integration of diverse data sources can improve the knowledge base
 59 (see Alavi and Leidner (2001); Hendler (2014) for a detailed discussion) and thus, improving the
 60 performance of the expert system.

61 Based on the insights from Table 1 and the discussion above, we outline a novel methodology to
 62 predict the future movements in the value of securities after tapping data from disparate sources,
 63 including: (a) the number of page visits to pertinent Wikipedia pages; (b) the amount of online

Table 1: A review of financial expert systems that are used in stock movement prediction. ANN, GA, SVM and DT correspond to artificial neural network, genetic algorithm, decision tree and support vector machine, respectively.

Paper	Sources for Knowledge Base			AI Approach
	Traditional	Crowd-sourcing	News	
Kimoto et al. (1990)	✓			ANN
Lee and Jo (1999)	✓			Time Series
Kim and Han (2000)	✓			ANN, GA
Kim (2003)	✓			SVM
Qian and Rasheed (2007)	✓			ANN, DT
Li and Kuo (2008)	✓			ANN
Schumaker and Chen (2009)			✓	SVM
Vu et al. (2012)		✓		DT
Chen et al. (2013)	✓			ANN
Adebiyi et al. (2014)	✓			ANN, ARIMA
Nguyen et al. (2015)		✓		SVM
Shynkevich et al. (2015)			✓	ANN, SVM
Chourmouziadis and Chatzoglou (2016)	✓			Fuzzy System
Our Financial Expert System	✓	✓	✓	ANN, SVM, DT

content produced on a particular day about a company, the stock of which is publicly traded; and
(c) commonly used technical indicators and company value indicators in stock value prediction.
In the AI component of our expert system, we compare the performance of ANN, SVM and DT for stock movement prediction. We have chosen these three specific approaches since: (i) neural networks have been widely deployed in intelligent trading systems (Kimoto et al. (1990); Li and Kuo (2008); Guresen et al. (2011); Bollen et al. (2011)); (ii) SVM was successfully used by Kim and Han (2000) and Schumaker and Chen (2009); and (iii) decision trees have been effectively used in crowd-sourced expert systems Vu et al. (2012). In those papers, the authors reported that these AI models outperformed the more traditional approaches. However, it is unclear whether: (1) such results will hold for our predictions since our knowledge base is more diverse, and (2) the results will hold when predicting different stocks and indices. Thus, our expert system will evaluate the performance of these models and select the best approach for a given prediction problem.

To demonstrate the utility of our system, we predict the one-day ahead movements in AAPL stocks over a three year period. Based on our case study, we show that the combination of online data sources with traditional technical indicators provide a higher predictive power than any of these sources alone. The remainder of the paper is organized as follows. In Section 2, we present a detailed description of the methodology we used to extract the data from the online sources, the variable selection techniques employed, and the corresponding predictive models. In Section 3, we highlight the main results and offer our perspective on their importance/interpretation. Our concluding remarks and recommendations for future work are provided in Section 4. In Appendices I-III, we explain how Google News data was captured, present the formulas for our generated

85 features, and define the predictors identified from our variable selection steps. We also present a
 86 copy of our full dataset, code and prediction tool at <https://github.com/binweng/ShinyStock>.

87 **2. METHODS**

88 To predict stock movements, we propose a data-driven approach that consists of three main
 89 phases, as shown in Figure 1. In Phase I, we scrape four sets of data from online resources. These
 90 datasets include: (a) publicly available market information on stocks, including opening/closing
 91 prices, trade volume, NASDAQ and the DJIA indices, etc.; (b) commonly used technical indicators
 92 that reflect price variation over time; (c) daily counts of Google News on the stocks of interest;
 93 and (d) the number of unique visitors for pertinent Wikipedia pages per day. We also populated
 94 additional features (i.e. summary statistics) in an attempt to uncover more significant predictors
 95 for stock movement. In Phase II, we use variable selection methods to select a subset of predictors
 96 that provide the most predictive power/accuracy. Then, in Phase III, we utilize three AI tech-
 97 niques to predict stock movement. These models are compared and evaluated based on a 10-fold
 98 cross validation sample using the area under the operating characteristics curve (AUC) and seven
 99 other metrics. Based on the evaluation, we select an appropriate model for real-time stock market
 100 prediction. We present the details for each of the phases in the subsections below.

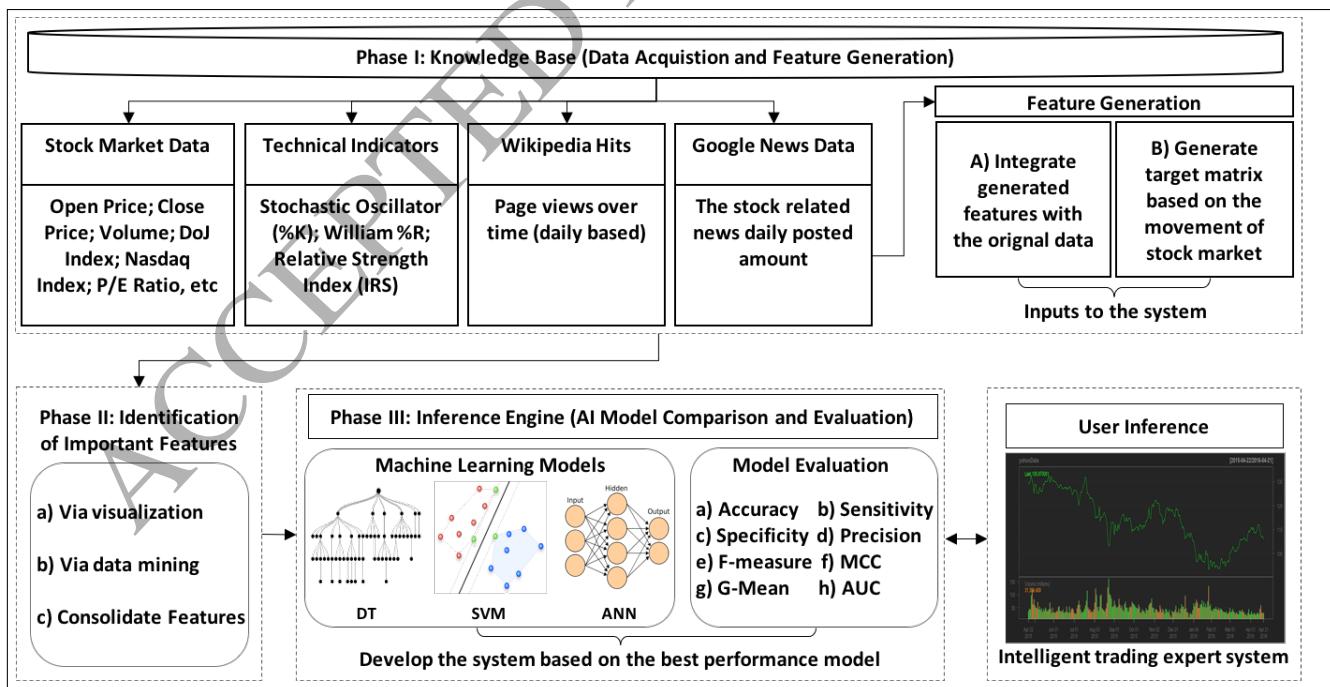


Figure 1: An Overview of the Proposed Method

101 *2.1. Data Acquisition and Feature Generation for Our “Knowledge Base”*

102 In this paper, we focus on predicting the AAPL, Apple NASDAQ, stock movement based on a
 103 37 month-period from May 1, 2012 to June 1, 2015. There are four datasets that were obtained,
 104 preprocessed and merged in Phase I. First, we obtain publicly available market data on AAPL
 105 using the *Yahoo Finance* website. We considered the following common predictors of stock prices
 106 (see e.g., Wang (2002); Lee and Jo (1999); Li and Kuo (2008); Jasemi et al. (2011)): the daily
 107 opening and closing prices, daily high/low, and volume of trades of the AAPL stock. In addition,
 108 we included the day-to-day movements in the DJIA and NASDAQ composite indices as indirect
 109 measures of risk that the AAPL stock is subject to due to the general market movements. We also
 110 used the price to earnings ratio (P/E) as an estimate for the fundamental health of the company
 111 (Gabrielsson and Johansson (2015)).

112 The second set of predictors is comprised of three indicators that are used in technical analysis.
 113 Technical analysis is used to forecast future stock prices by studying historical prices and volumes
 114 (Chourmouziadis and Chatzoglou (2016)). Since all information is reflected in stock prices, it is
 115 sufficient to study specific technical indicators (created by mathematical formula) to predict price
 116 fluctuations and evaluate the strength of the prevailing trend (Bao and Yang (2008)). In this paper,
 117 we consider three technical indicators:

- 118 (A) Stochastic Oscillator (%K), developed by George C. Lane as a momentum indicator that can
 119 warn of the strength or weakness of the market. When the market is trending upwards, it
 120 tries to measure when the closing price would get close to the lowest price in a given period.
 121 On the other hand, when the market is trending downwards, it estimates when the closing
 122 price would get close to the highest price in the given period. For additional details on the
 123 %K and its calculation, the reader is referred to: Bao and Yang (2008) and Lin et al. (2011).
- 124 (B) The Larry William (LW) % R Indicator - It is a momentum indicator that facilitates the
 125 spotting of overbought and oversold levels. For its calculation, refer to Kim and Han (2000).
- 126 (C) The Relative Strength Index (RSI)- Similar to the LW %R, it compares the magnitude of
 127 recent gains to recent losses in an attempt to determine overbought and oversold conditions
 128 of an asset. RSI ranges from 0 to 100. In practice, investors sell if its value is ≥ 80 and buy
 129 if it is ≤ 20 . For more details, see Bao and Yang (2008) and Lin et al. (2011).

130 The reader should note that the values for these three technical indicators were calculated based
 131 on the market price data obtained from *Yahoo Finance*.

132 In the third data source, we scrape the amount of daily online content produced about a company,

133 and its products/services. In this paper, we obtain a count for aggregated news and blogs based on
 134 the daily count of content on Google News. We detail this step in Appendix I. The fourth and final
 135 data source is based on the Wikipedia page view counts of terms related to Apple stock (AAPL,
 136 Apple Inc., iPhone, iPad, Macbook, and Mac OS). We queried the daily visits for these pages
 137 from www.wikipediatrends.com. A graphical summary of the second and third set of predictors is
 138 provided in Figure 2.

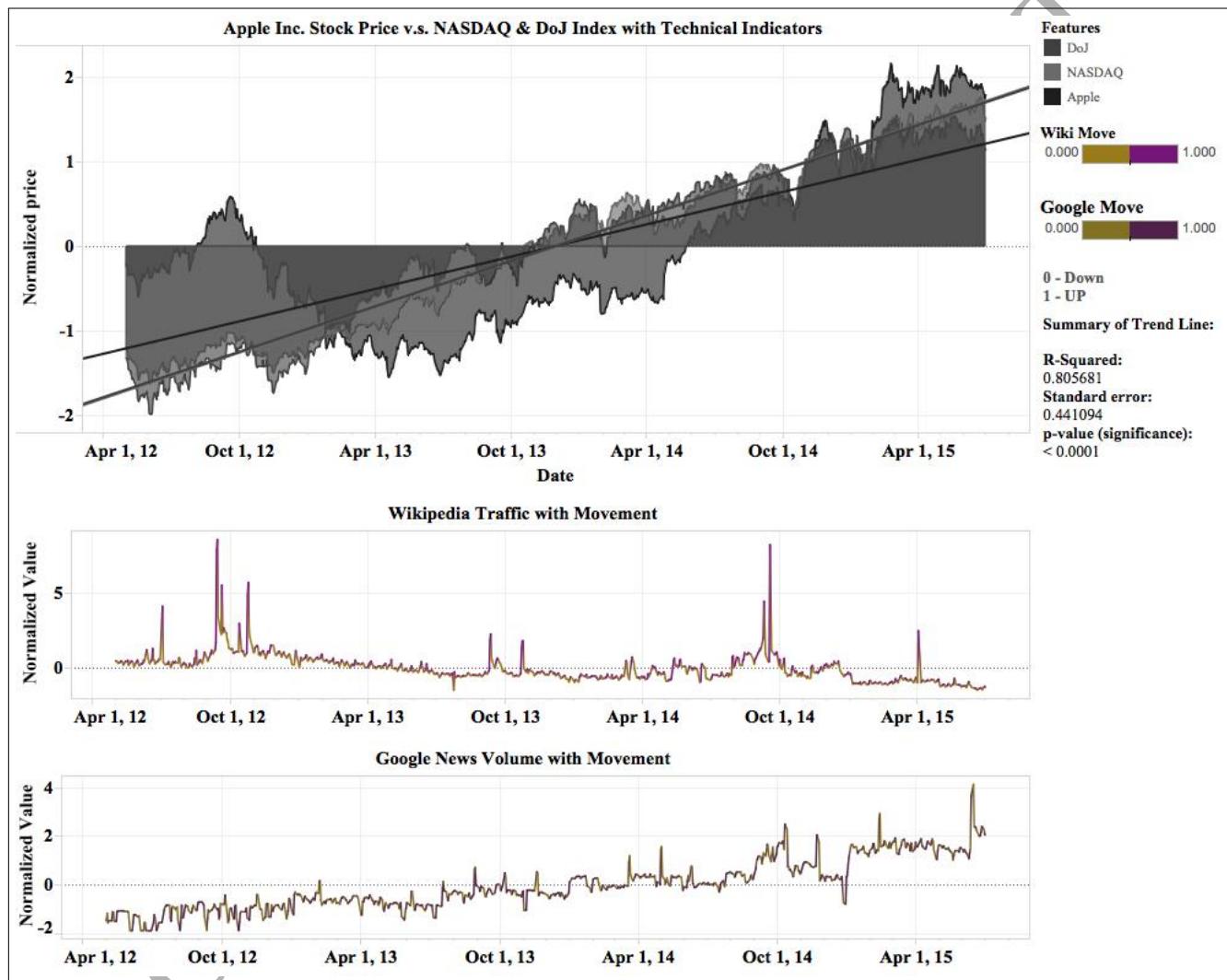


Figure 2: A Visual Summary of the Main Predictors from the Four Data Sources. An interactive version of this plot is available at: <https://goo.gl/fZSQEy>. Note that we rescaled the variables (by subtracting the mean and dividing by the standard deviation) to facilitate the visualization of the data.

139 To enhance the performance of the predictive models, we generate some additional features from
 140 the four predictor sets. We incorporate some of the underlying principles behind technical analysis
 141 (see e.g., Bao and Yang (2008)) to generate our feature set. Therefore, our generated features
 142 include: Wikipedia Momentum, Wikipedia Rate of Change, Google Momentum, Google Relative

¹⁴³ Strength Index, and three moving averages of stock prices (where $n = 3, 5$, and 10 , respectively).
¹⁴⁴ For the sake of completion, we explain how each of these features are calculated in Appendix II.

¹⁴⁵ *2.2. Variable/Feature Selection*

¹⁴⁶ The end goal of this phase is to have the data processed for the artificial intelligence models.
¹⁴⁷ This phase is comprised of two steps. First, we define different one-day-ahead outcomes (hereafter
¹⁴⁸ targets). Then, we use recursive feature elimination (RFE) to select the features/variables that
¹⁴⁹ offer the highest predictive power.

¹⁵⁰ There are several one-day-ahead outcomes that can be of interest to investors. We examine five
¹⁵¹ different targets. These targets are defined in Table 2. Target 1 compares the opening stock price
¹⁵² of day $i + 1$ with the closing price of the previous trading day. In Target 2, we compare the opening
¹⁵³ stock price of day $i + 1$ with the opening price of the previous trading day. Targets 3 and 4 follow a
¹⁵⁴ similar logic with the closing price used for day $i + 1$ instead of the opening price. In Target 5, we
¹⁵⁵ examine the differences in trade volume between day $i + 1$ and day i . It is important to note that we
¹⁵⁶ only calculate these targets for the AAPL stock as a case study. In addition, we have transformed
¹⁵⁷ all targets to a binary variable where $0 \rightarrow$ no increase in target, and $1 \rightarrow$ an increase in the target
¹⁵⁸ value from the previous day.

Table 2: One-Day-Ahead Targets Used in this Paper

Target	Formula
Target 1	$Open(i + 1) - Close(i)$
Target 2	$Open(i + 1) - Open(i)$
Target 3	$Close(i + 1) - Close(i)$
Target 4	$Close(i + 1) - Open(i)$
Target 5	$Trade Volume(i + 1) - Trade Volume(i)$

¹⁵⁹ In Step 2, we selected the significant features using the SVM recursive feature elimination (RFE)
¹⁶⁰ algorithm. RFE is implemented through backwards selection of predictors based on predictor im-
¹⁶¹ portance ranking. The predictors are ranked and the less important ones are sequentially eliminated
¹⁶² prior to deploying our three predictive models. The goal of this step is to find a subset of predictors
¹⁶³ that can result in accurate predictions without overfitting. It should be noted that we used the
¹⁶⁴ SVM-RFE algorithm for each target. Thus, we have obtained different predictor sets for each tar-
¹⁶⁵ get. These predictors are presented in the Table 3 in Section 3. For more details on how RFE can
¹⁶⁶ be deployed using open-source programming languages, the readers are referred to the **R** Package

167 *Caret* (Kuhn (2008)) and to the *sklearn.feature_selection* module in *Python* (Scikit-Learn-Developers
 168 (2014)).

169 *2.3. The Inference Engine: AI Model Comparison and Evaluation*

170 In this phase, we compare the effectiveness of artificial neural networks (ANN), decision trees
 171 (DT), and support vector machines (SVM) for predicting movements in the AAPL stock based on
 172 the predictors identified in Subsection 2.2. In the paragraphs below, we first introduce how we used
 173 a 10-fold cross validation approach to minimize the sampling bias. Then, we provide a very short
 174 overview of the three classification approaches, and introduce the performance evaluation metrics
 175 used to identify the most suitable approach. The reader should note that, in this paper, we deploy
 176 the described methodology for each of the five targets. Hereafter, we use the term *dataset* to reflect
 177 each set of features/variables with its associated target for the AAPL stock over the 37 months of
 178 the study.

179 The k -fold cross-validation approach is used to minimize the bias associated with the random
 180 sampling of the training and test data samples (Kohavi (1995)). The entire dataset is randomly
 181 split into k mutually exclusive subsets of approximately equal size. The prediction model is tested k
 182 times by using the test sets. The estimation of the k -fold cross validation for the overall performance
 183 criteria is calculated as the average of the k individual performances as shown in Dag et al. (2016).
 184 In our analysis, we use the stratified 10-fold cross validation approach to estimate the performance
 185 of the different classification models. Our choice for $k = 10$ is based on literature results (see
 186 e.g., Kohavi (1995); Dag et al. (2016)) that show that 10-folds provide an ideal balance between
 187 performance and the time required to run the folds.

188 ANNs are widely employed in a wide variety of computational data analytics problems that
 189 include classification, regression and pattern recognition. In the context of stock market prediction,
 190 ANNs have been extensively applied in predicting stocks and indices at different markets (see Kim
 191 and Han (2000); Hassan et al. (2007); Atsalakis and Valavanis (2009); Zhang and Wu (2009); Dase
 192 and Pawar (2010); Bollen et al. (2011); Guresen et al. (2011), and the references within). We assume
 193 that the reader is familiar with ANNs and their construction (otherwise, refer to Hastie et al. (2011)).
 194 In this paper, we use the sigmoid function as the activation function for our ANN. We have also
 195 used the Multi-layer Perceptron (MLP) learning model with a back-propagation algorithm due to
 196 its superior performance to the radial basis function (RBF) in our preliminary analysis.

197 Decision trees are widely used in several data mining and stock market prediction problems

(Qian and Rasheed (2007); Lai et al. (2009); Atsalakis and Valavanis (2009)) since they are very easy to interpret. The modeling procedure starts with splitting the dataset into several subsets each of which consists of more or fewer homogeneous states of the target variable (Breiman et al. (1984)). Then the impacts of each independent variable on the target variable are measured. This procedure takes place successively until the decision tree reaches a stable state. Popular decision tree algorithms include Quinlan's ID3, C4.5, C5 (Quinlan (1986, 2014)) and C&RT (Breiman et al. (1984)). In our data analysis, the C5 algorithm was used since it: a) is computationally efficient; and b) has outperformed the other methods examined in our preliminary analysis.

Similar to the previous two other classification approaches, SVM is a popular approach for stock market prediction (Yang et al. (2002); Kim (2003); Schumaker and Chen (2009); Nassirtoussi et al. (2014)). More interestingly, SVMs are favored in applications where text mining is used for market prediction (Nassirtoussi et al. (2014)). SVMs can be used for both linearly and non-linearly separable datasets. When the data is linearly separable, SVMs construct a hyperplane on the feature space to distinguish the training tuples in the data such that the margin between the support vectors is maximized. For nonlinear cases, the data is typically mapped into a higher-dimensional space so that the new dataset in higher-dimension becomes linearly separable. This problem can be handled efficiently by using a Kernel function (see Han et al. (2011) for more details). Based on our preliminary analysis, we have used the Radial Basis Function (RBF) Kernel function in our SVM classification algorithm since it has resulted in the best performance.

To evaluate the performance of the three classification methods, we present eight commonly used metrics in the literature: a) *accuracy*, b) *area under the receiver operating characteristic curve (AUC)*, c) *F-measure*, d) *G-mean*, e) *MCC*, f) *precision*, g) *sensitivity*, and h) *specificity*. In addition, we provide our code and the confusion matrix for the sake of completion. Our selected measures are all suitable for our binary classification problem. For details on how any of the above metrics can be calculated, we refer the reader to Han et al. (2011), and Hastie et al. (2011). We use the AUC as our primary evaluation metric for the reasons explained in Dag et al. (2016).

3. RESULTS AND DISCUSSION

In this section, we first highlight the results from the variable/feature selection phase of our methodology. Then, we present the results from the prediction accuracy of our expert system with respect to the five potential targets. This is followed by some preliminary analysis to evaluate the impact of the information attained from the five different data sources on our prediction power.

229 For the sake of completion and to allow for the replication of our results, we present our code and
 230 a detailed tabular view of our results as supplementary documents to this manuscript.

231 *3.1. Variable/Feature Selection*

232 As mentioned in Section 2.2, the end goal of this phase is to prepare the data for the three
 233 machine learning models. Here, we employed the SVM RFE model five times (once for each target).
 234 This resulted in five different sets of twenty variables/features that offer the most predictive power
 235 for each of the five respective targets. We list these sets in Table 3. There are several additional
 236 observations to be made from Table 3:

- 237 (A) For any of the five targets, the selected variables/features span all predictor sets. This implies
 238 that there are non-redundant, useful information that can be captured from each data source.
- 239 (B) The previous day's closing, opening and high prices were significant predictors for all five
 240 targets. The previous day's low price is a strong predictor for four of the five targets (with
 241 the exception of Target 3).
- 242 (C) The *Price to Earnings (P/E) Ratio* is predictive for the four price targets, but not for Target
 243 5 (i.e., trade volume target). In our opinion, this makes sense since the P/E Ratio measures
 244 the current share price relative to the per-share earnings. Thus, it may not be suited for
 245 predicting trade volume since it does not capture any movements.
- 246 (D) Target 5 had the highest number of Google features of 10. This was a somewhat expected
 247 result since *Google Trends* should reflect interest more than price fluctuations. The number
 248 of Google features selected for any of the other targets varied between 4 and 6.
- 249 (E) Perhaps the most important observation has to do with the order of the variables/features
 250 selected. We have arranged the items in a descending order (left to right and then to the
 251 next row). For all targets, variables selected from the *first set of predictors* were the most
 252 significant predictors. They were followed by one or more *technical indicators*. Then, the list
 253 would include several *Wikipedia* features, which were followed by some *Google News* features.
 254 The final grouping included a mixture of *technical indicators* and *Google/Wikipedia features*.
- 255 Note that we provide the definition for each of the features listed in Table 3 in Appendix III.

256 *3.2. Predictive Modeling Outcomes*

257 As explained in Section 2.3, we use the AUC as the primary evaluation criterion to evaluate the
 258 performance of the ANN, DT, and SVM models in predicting the five different day-ahead outcomes

Table 3: The Twenty Most Predictive Variables/Features for Each Target

Target	Variables/Features Selected				
Target 1	Close	Open	High	Low	P/E Ratio
	Wiki_3.day.disparity	Wiki_5.day.disparity	Wiki_10.day.disparity	Wiki_Momentum_1	Wiki_ROC
	Google_MA_5	Google_EMA_3	Google_3.Day.disparity	Google_5.day.disparity	RSI
	Stochastic Oscillator (%K)	Wiki_RSI	Google_MA_4	William %R	Google_MA_3
Target 2	Close	Open	High	Low	P/E Ratio
	Wiki_5.day.disparity	Wiki_Move	Wiki_MA3.Move	Wiki_EMA5.Move	Wiki_5day_disparity_Move
	Google_EMA5.Move	Google_3day_disparity_Move	Google_ROC.Move	Google_RSI.Move	Wiki_3.day.disparity
	Stochastic Oscillator (%K)	RSI.Move	Wiki_RSI.Move	Google_MA_6	Google_Move
Target 3	Close	Open	High	P/E Ratio	Stochastic_Move
	Wiki_Momentum_1	Wiki_Move	Wiki_MA3.Move	Wiki_EMA5.Move	Wiki_ROC.Move
	Google_EMA5.Move	Google_3day_disparity_Move	Google_ROC.Move	Google_RSI.Move	Wiki_10.day.disparity
	RSI.Move	Wiki_RSI.Move	Wiki_3.day.disparity	Google_Move	Google_MA5.Move
Target 4	Close	Open	High	Low	P/E Ratio
	RSI.Move	Wiki_10.day_Disparity	Wiki_Move	Wiki_MA3.Move	Wiki_EMA5.Move
	Google_Move	Google_3day_disparity_Move	Google_ROC.Move	Google_RSI.Move	William %R
	Stochastic Oscillator (%K)	Stochastic_Move	Wiki_3day_disparity.Move	Wiki_ROC.Move	Wiki_RSI.Move
Target 5	Close	Open	High	Low	William %R
	Wiki_Momentum_1	Wiki_RSI	Google_MA_2	Google_MA_3	Google_MA_4
	Google_MA_9	Google_3.day.disparity	Google_5.day.disparity	Google_10.day.disparity	Wiki_10.day.disparity
	Wiki_3.day.disparity	Wiki_5.day.disparity	Google_MA_6	Google_MA_7	Google_MA_8

(while presenting the 7 other metrics for completion). In Figure 3, we present the best-case, worst-case and the mean performance of the three machine learning models for each of the five targets. Note that the best-case, worst-case, and mean performances are determined based on the 10-fold cross validation step of our approach. The reader is encouraged to visit the interactive version of this plot at <https://goo.gl/L06FSA>. Based on Figure 3, there are several interesting observations that can be made:

- (A) Based on the AUC metric, SVM outperforms the ANN model for all five targets, DT outperforms the ANN model for Targets 2-4, DT outperforms the SVM model in predicting Targets 2-3 (while having a similar performance in Target 4), and the DT model failed to predict one of the classes for both Targets 1 and 5.
- (B) For Targets 2-4, the recommended models have an AUC value greater than 0.89. The AUC is the probability that the model will rank a randomly chosen positive instance (i.e. increase in price) higher than a randomly chosen negative one (i.e., decrease in stock price).
- (C) The acquired data may not be capturing the underlying factor's for changes in trade volume (i.e., Target 5). The DT model could not predict decreases in trade volume (i.e., all its predictions were "1"s), and the ANN has a similar prediction to that of a random predictor (i.e., flipping a coin). The SVM had a somewhat reasonable mean AUC value of 0.632.
- (D) Based on the eight evaluation metrics' values, our disparate data sources and machine learning models can best predict Target 2. Recall that Target 2 compares *next day's opening price* with *today's opening price*. This is a somewhat surprising result since we expected Target 1

279 to have the best results.

- 280 (E) Perhaps more importantly, our results (especially for Target 2) are more accurate than those
 281 typically reported in the literature. Our model resulted in $\approx 85\%$ accuracy/hit ratio with
 282 an average AUC of > 0.874 for *SVM and DT for Target 2*. In the literature, the previous
 283 predictions did not exceed an accuracy of 83% (see Table 5 in Nassirtoussi et al. (2014), which
 284 summarizes the outcomes of 24 text-mining-based financial expert systems).
- 285 (F) Building on the previous result, it is also clear that the addition of data from disparate data
 286 sources have resulted in improved accuracy. For example, Kim (2003) used the SVM model
 287 with only technical indicators as inputs, and obtained an accuracy rate of 65% accuracy for
 288 their best performance model. Our $> 20\%$ accuracy improvement (when SVM or DT are
 289 used) is significant and justifies the effort needed to include new data sources.

290 From the above discussion, we have established that we can predict Targets 1-4 reasonably well
 291 through the deployment of an adequate machine learning model with inputs identified in Table 3.
 292 To formally understand the usefulness of the four disparate data sources and our generated features,
 293 we consider several scenarios that are summarized in Table 4. Note that *Scenarios 1-2* involve the
 294 data sources most commonly used in traditional stock market prediction. *Scenarios 3-4* build on
 295 *Scenario 2* with the additional of one online data source. In *Scenarios 5-6*, we add the generated
 296 features to *Scenarios 3* and *4*, respectively. Scenario 7 include all five data sources.

Table 4: Examining the Impact of the Non-Traditional Data Sources

Scenario #	Description
1	Market data
2	Market data, technical indicators
3	Market data, technical indicators, Wikipedia Traffic
4	Market data, technical indicators, Google news counts
5	Market data, technical indicators, Wikipedia Traffic, generated features
6	Market data, technical indicators, Google news counts, generated features
7	Market data, technical indicators, Wikipedia traffic, Google news counts, and generated features

297 As an example, consider the SVM model for Target 2. Let us examine how the inclusion of data
 298 sources, according to the seven scenarios presented in Table 4, impact the eight evaluation metrics.
 299 We present the results in Table 5. From the results, it is clear that the best performance is obtained
 300 when all data sources are included. In addition, by comparing S5 and S6 (or alternatively S3 and

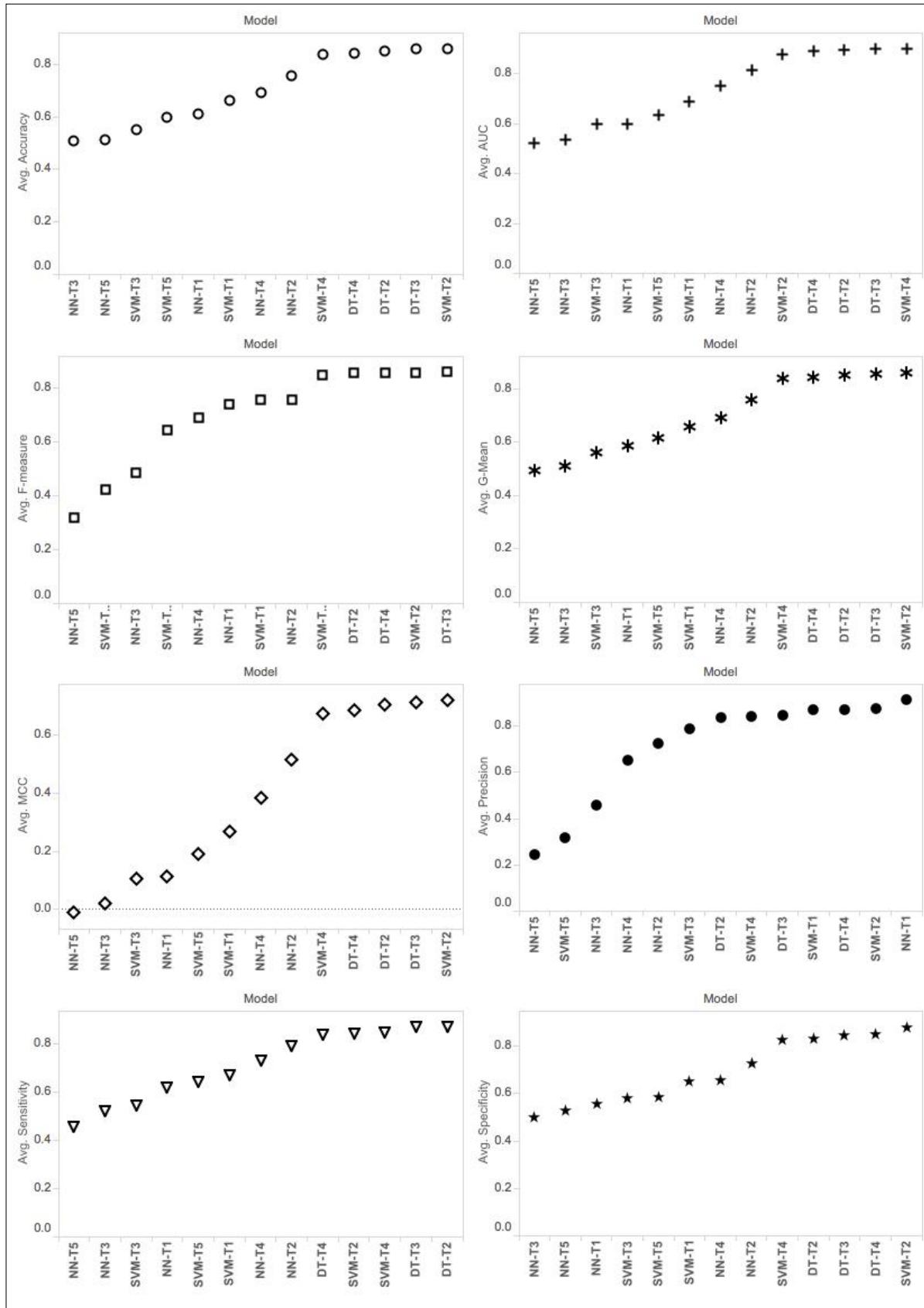


Figure 3: A Visual Summary of the Performance of the 3 Data Mining Models for Each of the Five Targets. An interactive version of this plot can be found at: <http://goo.gl/L06FSA>

301 S4), one can see that the Wikipedia data is more informative than the Google Data for Target 2.
 302 Note that we consider the results presented in Table 5 as a formal way to evaluate our observation
 303 in Point (F) above.

Table 5: Comparison of seven scenarios using eight evaluation metrics

Scenario	Accuracy	Sensitivity	Specificity	Precision	F-measure	MCC	G-Mean	AUC
S1	0.565	0.577	0.551	0.601	0.589	0.127	0.564	0.634
S2	0.616	0.618	0.614	0.648	0.633	0.232	0.616	0.711
S3	0.618	0.634	0.601	0.629	0.632	0.235	0.617	0.703
S4	0.618	0.639	0.595	0.639	0.639	0.233	0.616	0.708
S5	0.822	0.835	0.807	0.824	0.830	0.642	0.821	0.800
S6	0.813	0.821	0.804	0.805	0.813	0.625	0.813	0.856
S7	0.858	0.838	0.879	0.873	0.854	0.719	0.858	0.874

304 4. CONCLUSIONS AND FUTURE WORK

305 4.1. An Overview of the Impacts and Contributions of our Proposed Expert System

306 In this paper, we developed a financial expert system to predict movements in the one-day
 307 ahead stock price/volume. To construct our “knowledge base”, we scrapped four different data
 308 sets: (i) historical stock market data, (ii) commonly used *technical indicators*, (iii) *Wikipedia* traffic
 309 statistics pertaining to the company’s pages (i.e. general company profile, stock page, and pages
 310 pertaining to the company’s main products), and (d) *Google News*. Since these data sources were
 311 never used in combination in an expert system, we generated features from the *two online sources*
 312 to further improve our knowledge base. Our AI framework consisted of two major phases: (1)
 313 variable/ feature selection, which helps improve the performance of our AI algorithms by reducing
 314 the dimensions of the data without the loss of information; and (2)the incorporation of ANN, SVM
 315 and DT for prediction, which allows us to select the “best” model for a given target and stock. We
 316 provide a web-based user interface (see <https://github.com/binweng/ShinyStock>) to promote the
 317 adoption of our expert system by investors and financial planners.

318 From an *Expert and Intelligent Systems* research perspective, our system is innovative and
 319 novel. Specifically, the related literature on stock movement prediction (shown in Table 1) primarily
 320 considered the use of traditional data sources (i.e. market data and technical indicators) and none,
 321 to our knowledge, combined multiple data sources. Our system utilizes disparate data sources in
 322 an attempt to have a more holistic representation of the factors and conditions that precede stock
 323 movement. The proposed expert system is tested using a large and feature-rich *Apple Inc.* dataset

324 collected for a period of 37 months (May 1, 2012 to June 1, 2015), providing a hit ratio of 85%
 325 (which exceeds the reported results in the literature). Perhaps more importantly, we have addressed
 326 the following theoretical questions that relate to the design of expert and intelligent systems:

- 327 (a) What is the value of using online sources (specifically Wikipedia and Google News) when
 328 predicting the one-day ahead stock movement? In contrast with the majority of literature,
 329 we analyze this question through combining variables/features from these online sources with
 330 more traditional predictors. This allows us to quantify the value added rather than just
 331 obtaining a predictive model.
- 332 (b) Does the added value of these online sources differ with different targets? We chose five
 333 different *one-day-ahead* targets to examine if the value obtained from these sources changes
 334 according to different prediction questions.
- 335 (c) Which targets are most suitable for prediction based on the aforementioned five data sources?
- 336 (d) Which AI models provide the best predictive performance for each of the five targets?

337 From our case study, we have learned that the addition of these online sources are useful (especially
 338 for Targets 1-4). In addition, based on the Apple stock, it seems that Wikipedia has more predictive
 339 power than Google News. That being said, the addition of Google News indicators improve the
 340 predictive accuracy the AI models utilized by our expert system (see Figure 3 and Table 5). From
 341 our seven scenarios of data aggregation, it is clear that the addition of online data sources and our
 342 generated features can significantly improve the prediction accuracy. This can imply that there
 343 are *news* hidden in these sources according to the followers of the *Efficient Market Hypothesis*.
 344 Alternatively, one can say that changes in these data sources precede changes in the stock market.
 345 Our analysis also indicates that all five targets can be predicted (using the best model) better than
 346 a coin-flip. Our intelligent system’s prediction performance is better than the results reported in
 347 the literature (see Section 3.2).

348 *4.2. Implementing our Expert System in Practice*

349 From an *Expert and Intelligent Systems* practical implementation perspective, our proposed
 350 system can be used in a number of different ways. First, on a basic level and through our interface,
 351 an investor who does not have a strong programming background can use our “knowledge base”
 352 to capture the total number of “Google News” articles and visitors of relevant Wikipedia pages.
 353 Through our plotting tools, that investor can visualize the crowd’s perception of a given stock or
 354 index. We have shown that these perceptions can be predictive of stock movement. It is important

355 to note that this information is not available by current commercial products. Second, on a more
 356 advisory level, our expert system can be used to provide a data-driven recommendation for investors;
 357 an informed short term buy, or sell strategy of stocks can be made relative to whether the investors
 358 portfolio carries the stock. From that viewpoint as well, investors can use our system to construct an
 359 ensemble of predictions (with at least two models - their current approach and our expert system's
 360 recommendation). In the case of a two-model scenario, our expert system can indirectly help
 361 with quantifying risk/uncertainty (i.e., if both models agree, the likelihood of a correct outcome
 362 increases). If the investor already had access to multiple forecasting systems, then our expert
 363 system will present a new perspective on a stock since our model combines both traditional and
 364 nontraditional sources. In such a case, the investor can make his/her decision through a simple
 365 voting procedure. Third, our code, which provide through a link in this article, can be deployed in
 366 an existing fully automated short term trading system to make its decision-making process more
 367 comprehensive.

368 *4.3. Limitations and Future Research*

369 There are several limitations and opportunities for future work that arise from this study. First,
 370 we have only examined Apple stock over a certain time-period. Thus, it is not clear if our results
 371 and/or conclusions can extend prior or past this period. More generally, it would be interesting
 372 to examine if our conclusion would differ if a different type of commodity stock is chosen and/or
 373 if a stock index is desired. Second, we did not attempt to include other online data sources. It
 374 is not clear if the relevance of our sources would change if, for example, Facebook data is used.
 375 Therefore, there are several opportunities to extend our work by the inclusion of additional data
 376 sources. We expect a diminishing return with the inclusion of new data sources, since we expect
 377 some redundancy in the information captured from online data sources. That being said, it would
 378 be interesting to rank the value obtained from the different online data sources (for different stocks
 379 and indices). A third direction can be to consider the stochastic nature of the prediction. Our
 380 “inference engine” presented a deterministic prediction; however in practice, it might be interesting
 381 to have a level of certainty that is associated with the prediction. This can be accomplished through
 382 the incorporation of fuzzy systems, Bayesian Belief Networks (BBN), and ensemble approaches. The
 383 fourth, and perhaps the largest improvement on this financial expert system, is to attempt to predict
 384 the actual price rather than the movement. From an investor's point of view, a 20% increase in
 385 stock price is very different than a 1% increase. In our analysis, these two scenarios are identical

386 since they are both coded as an increase in stock price.

387 In conclusion, this paper presented a financial expert system to predict the movement of a stock
 388 on a daily basis. We have shown that taking into consideration predictive factors from multiple
 389 sources can improve its predictive performance. We have also shown that the performance of the
 390 AI models can change significantly depending on the target used. To encourage future research, we
 391 provide our code and data in <https://github.com/binweng/ShinyStock>.

392 **5. APPENDICES**

393 *5.1. Appendix I - Process to Acquire Data from Google News*

394 *Google News* data is acquired from *Alphabet Inc.'s Google Search Engine*. The use of *Google*
 395 *News* allows us to gather all sources of news produced over a particular time period based on some
 396 search keywords. From a stock market perspective, this allows an end user to search for a publicly
 397 traded company's stock, and obtain a number for the amount of news produced for that stock. The
 398 steps to obtain the amount of news produced are shown below:

- 399 1. Go to www.google.com.
- 400 2. Input the search keyword, such as "AAPL, Apple Stock".
- 401 3. Click "News" and then "Search tools".
- 402 4. Custom the date range to the date you want to search.
- 403 5. Click "Search tools" again to show the result.

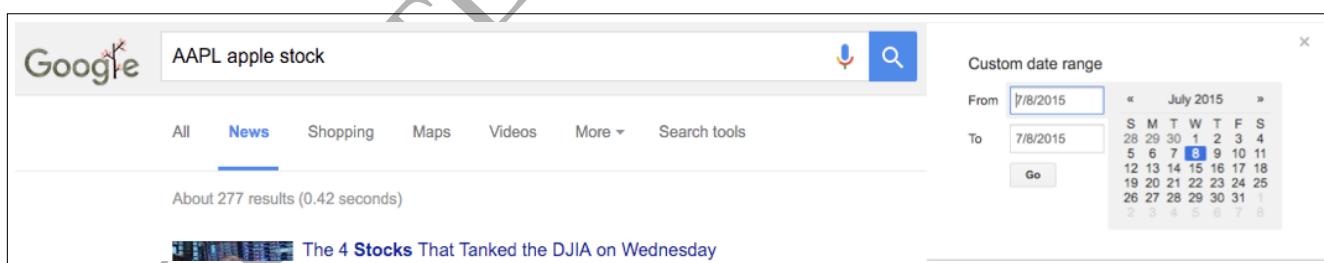


Figure 4: Screen-shot of An Illustration of Acquiring Data using Google's Search Engine

404 5.2. Appendix II - Formulas for the Generated Features

405 In this study, we generated seven different types of features for Wikipedia traffic data and Google
 406 news data. The formulas are shown below. In the formula, n means the time periods, V_t means the
 407 data point at period t.

408 1. *Moving Average*:

$$MA(n)_t = \frac{V_t}{n} + \frac{V_{t-1}}{n} + \dots + \frac{V_{t-n+1}}{n} \quad (1)$$

409 2. *Exponential Moving Average*:

$$EMA(n)_t = (V_t - MA(n)_{t-1}) \times \left(\frac{2}{n+1}\right) + MA(n)_{t-1} \quad (2)$$

410 3. *Disparity*:

$$Disparity(n)_t = \frac{V_t}{MA(n)_t} \times 100 \quad (3)$$

411 4. *Momentum1*:

$$Momentum1_t = \frac{V_t}{V_{t-5}} \times 100 \quad (4)$$

412 5. *Momentum2*:

$$Momentum2_t = (V_t - V_{t-5}) \times 100 \quad (5)$$

413 6. *Rate Of Change*:

$$ROC_t = \frac{V_t}{Momentum2_t} \times 100 \quad (6)$$

414 7. *Relative Strength Index*:

$$RSI(n) = 100 - \frac{100}{1 + \frac{AverageGain(n)}{AverageLoss(n)}} \quad (7)$$

⁴¹⁵ 5.3. Appendix III - Definition of Variables/Features in Table 3

We define the variables/features used in our model in the table below.

Table 6: Definition of The Most Predictive Variables/Features

Variable	Definition
Close	Closing price of the day
Google_x_day_disparity	Ratio of Google news volume to its x day moving average
Google_x_day_disparity_Move	Movement of Google_x_day_disparity as previous day
Google_EMA_x	x day exponential moving average of Google news volume
Google_EMA_x_Move	Movement of Google_EMA_x as previous day
Google_MA_x	x day moving average of Google news volume
Google_MA_x_Move	Movement of Google_MA_x
Google_Move	Movement of Google news volumes as previous day
Google_ROC_Move	Movement of the rate of change for Google news volume as previous day
Google_RSI_Move	Movement of relative strength index for Google news volume as previous day
High	Highest price of the day
Low	Lowest price of the day
Open	Opening price of the day
P/E Ratio	Price-Earning ratio
RSI	Relative strength index of the stock price
RSI_Move	Movement of RSI
Stochastic Oscillator	Compares a security's closing price to its price range over a given time period
Stochastic_Move	Movement of Stochastic Oscillator
Wiki_x_day_disparity	Ratio of Wikipedia traffic to its x day moving average
Wiki_x_day_disparity_Move	Movement of Wiki_x_day_disparity
Wiki_EMA_x_Move	Movement of x day exponential moving average for Wikipedia traffic
Wiki_MA_x_Move	Movement of x day moving average for Wikipedia traffic
Wiki_Momentum_1	Ratio of current close price to the price three day's ago
Wiki_Move	Movement of Wikipedia as previous day
Wiki_ROC	Rate of change (ROC) of Wikipedia traffic
Wiki_ROC_Move	Movement of Wiki_ROC
Wiki_RSI	Relative strength index of Wikipedia traffic
Wiki_RSI_Move	Movement of Wiki_RSI
William %R	The level of the close price relative to the highest high

⁴¹⁶

⁴¹⁷ **References**

- ⁴¹⁸ Adebiyi, A. A., Adewumi, A. O., Ayo, C. K., 2014. Comparison of arima and artificial neural
⁴¹⁹ networks models for stock price prediction. Journal of Applied Mathematics 2014.

- ⁴²⁰ Alavi, M., Leidner, D. E., 2001. Review: Knowledge management and knowledge management
⁴²¹ systems: Conceptual foundations and research issues. *MIS quarterly*, 107–136.
- ⁴²² Atsalakis, G. S., Valavanis, K. P., 2009. Surveying stock market forecasting techniques—part ii: Soft
⁴²³ computing methods. *Expert Systems with Applications* 36 (3), 5932–5941.
- ⁴²⁴ Bao, D., Yang, Z., 2008. Intelligent stock trading system by turning point confirming and proba-
⁴²⁵ bilistic reasoning. *Expert Systems with Applications* 34 (1), 620–627.
- ⁴²⁶ Bollen, J., Mao, H., Zeng, X., 2011. Twitter mood predicts the stock market. *Journal of Computa-
⁴²⁷ tional Science* 2 (1), 1–8.
- ⁴²⁸ Breiman, L., Friedman, J., Stone, C. J., Olshen, R. A., 1984. Classification and regression trees.
⁴²⁹ CRC press.
- ⁴³⁰ Chen, M.-Y., Chen, D.-R., Fan, M.-H., Huang, T.-Y., 2013. International transmission of stock
⁴³¹ market movements: an adaptive neuro-fuzzy inference system for analysis of taiex forecasting.
⁴³² *Neural Computing and Applications* 23 (1), 369–378.
- ⁴³³ Choi, H., Varian, H., 2012. Predicting the present with google trends. *Economic Record* 88, 2 – 9.
- ⁴³⁴ Chourmouziadis, K., Chatzoglou, P. D., 2016. An intelligent short term stock trading fuzzy system
⁴³⁵ for assisting investors in portfolio management. *Expert Systems with Applications* 43, 298–311.
- ⁴³⁶ Cootner, P. H., 1964. The random character of stock market prices.
- ⁴³⁷ Dag, A., Topuz, K., Oztekin, A., Bulur, S., Megahed, F. M., 2016. A preoperative recipient-donor
⁴³⁸ heart transplant survival score. *Decision Support Systems*, –.
⁴³⁹ URL <http://dx.doi.org/10.1016/j.dss.2016.02.007>
- ⁴⁴⁰ Dase, R., Pawar, D., 2010. Application of artificial neural network for stock market predictions: A
⁴⁴¹ review of literature. *International Journal of Machine Intelligence* 2 (2), 14–17.
- ⁴⁴² Fama, E. F., 1965. The behavior of stock-market prices. *The Journal of Business* 38 (1), 34–105.
⁴⁴³ URL <http://www.jstor.org/stable/2350752>
- ⁴⁴⁴ Fama, E. F., 1991. Efficient capital markets: Ii. *The journal of finance* 46 (5), 1575–1617.
- ⁴⁴⁵ Fama, E. F., 1995. Random walks in stock market prices. *Financial analysts journal* 51 (1), 75–80.
- ⁴⁴⁶ Fama, E. F., Fisher, L., Jensen, M. C., Roll, R., 1969. The adjustment of stock prices to new
⁴⁴⁷ information. *International economic review* 10 (1), 1–21.
- ⁴⁴⁸ Gabrielsson, P., Johansson, U., 2015. High-frequency equity index futures trading using recurrent
⁴⁴⁹ reinforcement learning with candlesticks. In: *Computational Intelligence, 2015 IEEE Symposium*
⁴⁵⁰ Series on. IEEE, pp. 734–741.
- ⁴⁵¹ Guresen, E., Kayakutlu, G., Daim, T. U., 2011. Using artificial neural network models in stock
⁴⁵² market index prediction. *Expert Systems with Applications* 38 (8), 10389–10397.
- ⁴⁵³ Han, J., Kamber, M., Pei, J., 2011. Data mining: concepts and techniques. Elsevier.

- 454 Hassan, M. R., Nath, B., Kirley, M., 2007. A fusion model of hmm, ann and ga for stock market
 455 forecasting. *Expert Systems with Applications* 33 (1), 171–180.
- 456 Hastie, T. J., Tibshirani, R. J., Friedman, J. H., 2011. *The elements of statistical learning: data
 457 mining, inference, and prediction*. Springer.
- 458 Handler, J., 2014. Data integration for heterogenous datasets. *Big data* 2 (4), 205–215.
- 459 Jasemi, M., Kimiagari, A. M., Memariani, A., 2011. A modern neural network model to do stock
 460 market timing on the basis of the ancient investment technique of japanese candlestick. *Expert
 461 Systems with Applications* 38 (4), 3884–3890.
- 462 Kiersz, A., 03 2015. Here's how badly warren buffett beat the market.
 463 <http://www.businessinsider.com/warren-buffett-berkshire-hathaway-vs-sp-500-2015-3>.
- 464 Kim, K.-j., 2003. Financial time series forecasting using support vector machines. *Neurocomputing*
 465 55 (1), 307–319.
- 466 Kim, K.-j., Han, I., 2000. Genetic algorithms approach to feature discretization in artificial neural
 467 networks for the prediction of stock price index. *Expert systems with Applications* 19 (2), 125–132.
- 468 Kimoto, T., Asakawa, K., Yoda, M., Takeoka, M., 1990. Stock market prediction system with
 469 modular neural networks. In: *Neural Networks, 1990., 1990 IJCNN International Joint Conference
 470 on. IEEE*, pp. 1–6.
- 471 Kohavi, R., 1995. A study of cross-validation and bootstrap for accuracy estimation and model
 472 selection. In: *International joint conference on artificial intelligence*, 1995. pp. 1137–1143.
- 473 Kuhn, M., 2008. Caret package. *Journal of Statistical Software* 28 (5).
- 474 Lai, R. K., Fan, C.-Y., Huang, W.-H., Chang, P.-C., 2009. Evolving and clustering fuzzy decision
 475 tree for financial time series data forecasting. *Expert Systems with Applications* 36 (2), 3761–
 476 3773.
- 477 Lee, K., Jo, G., 1999. Expert system for predicting stock market timing using a candlestick chart.
 478 *Expert Systems with Applications* 16 (4), 357–364.
- 479 Li, S.-T., Kuo, S.-C., 2008. Knowledge discovery in financial investment for forecasting and trading
 480 strategy through wavelet-based {SOM} networks. *Expert Systems with Applications* 34 (2), 935
 481 – 951.
- 482 Lin, X., Yang, Z., Song, Y., 2011. Intelligent stock trading system based on improved technical
 483 analysis and echo state network. *Expert systems with Applications* 38 (9), 11347–11354.
- 484 Loomis, C. J., 02 2012. Buffett beats the sp for the 39th year.
 485 <http://fortune.com/2012/02/25/buffett-beats-the-sp-for-the-39th-year/>.
- 486 Malkiel, B. G., 2003. The efficient market hypothesis and its critics. *The Journal of Economic
 487 Perspectives* 17 (1), 59–82.
- 488 Mao, H., Counts, S., Bollen, J., 2011. Predicting financial markets: Comparing survey, news, twitter
 489 and search engine data. arXiv preprint arXiv:1112.1051.

- 490 Megahed, F. M., Jones-Farmer, L. A., 2015. Frontiers in Statistical Quality Control 11. Springer
491 International Publishing, Cham, Ch. Statistical Perspectives on “Big Data”, pp. 29–47.
- 492 Moat, H. S., Curme, C., Avakian, A., Kenett, D. Y., Stanley, H. E., Preis, T., 2013. Quantifying
493 wikipedia usage patterns before stock market moves. *Scientific reports* 3.
- 494 Nassirtoussi, A. K., Aghabozorgi, S., Wah, T. Y., Ngo, D. C. L., 2014. Text mining for market
495 prediction: A systematic review. *Expert Systems with Applications* 41 (16), 7653–7670.
- 496 Nguyen, T. H., Shirai, K., Velcin, J., 2015. Sentiment analysis on social media for stock movement
497 prediction. *Expert Systems with Applications* 42 (24), 9603–9611.
- 498 Nofsinger, J. R., 2005. Social mood and financial economics. *The Journal of Behavioral Finance*
499 6 (3), 144–160.
- 500 Prechter Jr, R. R., Parker, W. D., 2007. The financial/economic dichotomy in social behavioral
501 dynamics: the socioeconomic perspective. *The Journal of Behavioral Finance* 8 (2), 84–108.
- 502 Preis, T., Moat, H. S., Stanley, H. E., 2013. Quantifying trading behavior in financial markets using
503 google trends. *Scientific reports* 3.
- 504 Qian, B., Rasheed, K., 2007. Stock market prediction with multiple classifiers. *Applied Intelligence*
505 26 (1), 25–33.
- 506 Quinlan, J. R., 1986. Induction of decision trees. *Machine learning* 1 (1), 81–106.
- 507 Quinlan, J. R., 2014. C4. 5: programs for machine learning. Elsevier.
- 508 Schumaker, R. P., Chen, H., 2009. Textual analysis of stock market prediction using breaking
509 financial news: The azfin text system. *ACM Transactions on Information Systems (TOIS)* 27 (2),
510 12.
- 511 Scikit-Learn-Developers, 2014. 1.13 feature selection - scikit-learn documentation.
512 <http://goo.gl/GDedwn>.
- 513 Shynkevich, Y., McGinnity, T. M., Coleman, S., Belatreche, A., 2015. Stock price prediction based
514 on stock-specific and sub-industry-specific news articles. In: *Neural Networks (IJCNN), 2015
515 International Joint Conference on*. IEEE, pp. 1–8.
- 516 Smith, V. L., 2003. Constructivist and ecological rationality in economics. *The American Economic
517 Review* 93 (3), 465–508.
- 518 Vu, T.-T., Chang, S., Ha, Q. T., Collier, N., 2012. An experiment in integrating sentiment features
519 for tech stock prediction in twitter.
- 520 Wang, Y.-F., 2002. Predicting stock price using fuzzy grey prediction system. *Expert Systems with
521 Applications* 22 (1), 33 – 38.
- 522 Yang, H., Chan, L., King, I., 2002. Support vector machine regression for volatile stock market
523 prediction. In: *Intelligent Data Engineering and Automated LearningIDEAL 2002*. Springer, pp.
524 391–396.
- 525 Zhang, Y., Wu, L., 2009. Stock market prediction of s&p 500 via combination of improved bco
526 approach and bp neural network. *Expert systems with applications* 36 (5), 8849–8854.