

This is a pre-print version of the following article:

Evangelos Kalampokis, Efthimios Tambouris and Konstantinos Tarabanis: Understanding the Predictive Power of Social Media. Accepted for publication in Internet Research, special issue on the Predictive Power of Social Media. 2013.

Title: Understanding the Predictive Power of Social Media

Purpose: The purpose of this article is to consolidate existing knowledge and provide a deeper understanding of the use of Social Media (SM) data for predictions in various areas, such as disease outbreaks, product sales, stock market volatility, and elections outcome predictions.

Design/methodology/approach: The scientific literature was systematically reviewed to identify relevant empirical studies. These studies were analyzed and synthesized in the form of a proposed conceptual framework, which was thereafter applied to further analyze this literature, hence gaining new insights into the field.

Findings: The proposed framework reveals that all relevant studies can be decomposed into a small number of steps, and different approaches can be followed in each step. The application of the framework resulted in interesting findings. For example, most studies support SM predictive power, however more than one-third of these studies infer predictive power without employing predictive analytics. In addition, analysis suggests that there is a clear need for more advanced sentiment analysis methods as well as methods for identifying search terms for collection and filtering of raw SM data.

Value: The proposed framework enables researchers to classify and evaluate existing studies, to design scientifically rigorous new studies, and to identify the field's weaknesses, hence proposing future research directions.

Keywords: Social Networks; World Wide Web; Data Analysis; Open data.

1. Introduction

In the past years, the use of Social Media (SM) has dramatically increased with millions of users creating massive amounts of data every day. As of September 2012, the online social networking application Facebook reached one billion monthly active users, while the microblogging service Twitter reported more than 140 million active users. SM data is typically in the form of textual content (e.g. in blogs, reviews and status updates), rating scores in Likert scales or stars (e.g. review ratings), like or dislike indications (e.g. reviews helpful votes and Facebook's like or Google's '+1' buttons), Web search queries (e.g. Google trends), tags and profile information (e.g. social network graphs).

SM data incorporates personal opinions, thoughts and behaviours making it a vital component of the Web and a fertile ground for a variety of business and research endeavours. In this context, the predictive power of SM has been recently explored. For instance, empirical studies have analyzed the Yahoo! Finance message board to predict stock market volatility (Antweiler and Frank, 2004), weblog content to predict movies success (Mishne and Glance, 2006), Google search queries to track influenza-like illnesses (Ginsberg et al., 2009), Amazon reviews to predict product sales (Ghose and Ipeirotis, 2011) and Twitter posts (aka tweets) to infer levels of rainfall (Lampos and Cristianini, 2012).

These research efforts require cross-disciplinary skills as they involve both the transformation of noisy raw SM data into high quality data suitable for statistical analysis as well as the employment of predictive analytics, which comprise *'predictive models designed for predicting new/future observations or scenarios as well as methods for evaluating the predictive power of a model'* (Shmueli and Koppius, 2011: 555). In this setting, a number of researchers have recently challenged

This is a pre-print version of the following article:

Evangelos Kalampokis, Efthimios Tambouris and Konstantinos Tarabanis: Understanding the Predictive Power of Social Media. Accepted for publication in Internet Research, special issue on the Predictive Power of Social Media. 2013.

the methods employed and the results reported by empirical studies in the area. For instance, Jungherr et al. (2012) repeated the study conducted by Tumasjan et al. (2010) and reported controversial results. In addition, Gayo-Avello (2011) and Metaxas et al. (2011) conducted a number of experiments and criticized generalizations regarding the predictive power of SM.

This article aims at consolidating the knowledge created by empirical studies in recent years that exploit SM for predictions, thus enabling an in-depth understanding of SM predictive power. More specific objectives are: (a) to identify steps that characterize all relevant studies as well as approaches that can be followed in each step, and (b) to understand how different steps and approaches are related to SM predictive power.

We anticipate that the proposed framework will enable researchers to classify and evaluate existing studies, to design scientifically rigorous new studies, and to identify the field's weaknesses hence proposing future research directions.

The rest of the paper is structured as follows. Section 2 presents the research approach taken, while section 3 describes the proposed framework in detail. Section 4 presents the results of employing the framework to further analyze this literature, hence providing interesting results. Finally, section 5 draws conclusions.

2. Research Approach

In order to achieve the objectives of the paper we capitalize on the method proposed by Webster and Watson (2002) for conducting systematic literature reviews in the field of information systems. Initially, we performed a systematic search in order to accumulate a relatively complete body of relevant scientific literature. Towards this end, we started with Google Scholar using the key words *predict* OR *forecast* AND *social media* and we collected an initial pool of articles. Thereafter, we went

backward by reviewing citations in the identified articles and *forward* by using Google Scholar's functionality to identify articles citing the previously identified articles. We thereafter studied and filtered these initially identified articles in order to come up with the final set that was included in our research. For this purpose, we used the following inclusion and exclusion criteria:

- We excluded qualitative or purely theoretical articles (e.g. Louis and Zorlu, 2012).
- We included only studies aiming at making predictions. As a result, we have excluded empirical studies that aim at studying the relationship between SM data and phenomena outcome following an explanatory approach (e.g. Corley et al., 2010; Chen et al., 2011; Chevalier and Mayzlin, 2006; Chunara et al., 2012; Duan et al., 2008; Morales-Arroyo and Pandey, 2010; Reinstein and Snyder, 2005; Ye et al., 2006).
- We included only studies that attempt to predict real world outcomes. Thus, we excluded studies that predict online features such as tie strength (Gilbert and Karahalios, 2009), volume of comments on online news (Tsagkias et al., 2010) or movie rating on IMDB (Oghina et al., 2012).

This approach resulted in a set of 52 articles. For the sake of clarity, the list of these articles is presented at the end of this paper in the *Literature Review References* section.

In order to synthesize the accumulated knowledge we performed a *concept-centric* analysis. The main steps and most important aspects composing the whole prediction analysis process were extracted and combined in a conceptual SM data analysis framework for predictions that structures and depicts the area. Finally, the framework

was employed to further analyze the literature and to extract insights into the predictive power of SM.

3. The Social Media Data Analysis Framework

The proposed framework comprises two discrete phases, namely the *Data Conditioning Phase* and the *Predictive Analysis Phase*. The former refers to the transformation of noisy raw Social Media (SM) data into high quality data that is structured based on some *predictor variables*. The latter phase refers to the creation and evaluation of a *predictive model* that enables estimating outcome from a *new set of observations*.

Each of these phases can be further divided into a sequence of stages and each stage into a number of steps. Finally, different approaches can be followed in each step. Figure 1 presents our framework with the two phases, the respective stages along with their steps.

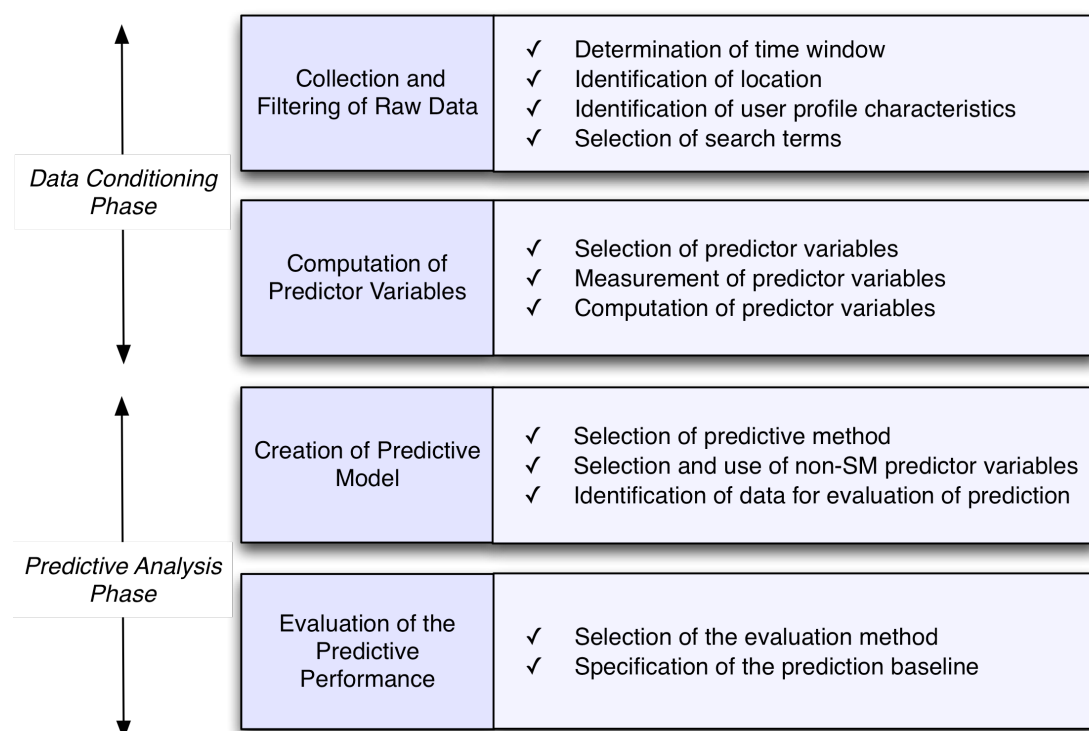


Figure 1 The two phases and the four stages of the Social Media data analysis framework for predictions along with the steps that compose each stage.

3.1 Phase I: Data Conditioning

The main purpose of the Data Conditioning phase is the transformation of noisy raw SM data into *high quality data* that will enable the computation of predictor variables. In order to define data quality we adopt and adapt a model proposed by Strong et al. (1997). In particular, we employ three data quality dimensions from Strong's model that we consider important in the SM data analysis realm, namely objectivity, completeness and amount of data.

Data objectivity is related to the accuracy of data production or the accuracy of the interpretation process, and specifies whether data is what it claims to be and measures what is supposed to measure. For instance, the data produced by interpreting text's sentiment could be of questionable objectivity in the case of non-rigorous sentiment analysis. The same holds when irrelevant data is interpreted as relevant. *Data completeness* deals with missing values from a data analysis perspective. It specifies whether or not collected data cover all aspects of a phenomenon in terms of e.g. entities characterizing it and/or predictor variables. Finally, *amount of data* (or sufficiency) specifies whether or not collected data is sufficient for predictive analysis.

The stages included in this phase along with the steps in each stage are described below.

3.1.1 Stage 1.1: Collection and Filtering of Raw Data

This stage deals with both raw SM data collection from various sources and filtering of data in order to determine those relevant. After its completion, the final data set that will be further analyzed during the next stage is produced. In order to determine the relevant raw data, the *when*, *where*, *who* and *what* questions should be answered.

For example, it can be inferred that a tweet mentioning the Conservative Party one week before the *UK elections of 2010* is related to these elections. The same holds for a tweet posted by David Cameron in the same period. The information used to determine relevance is extracted from the actual SM data or their metadata.

The effort required for this stage depends on both the SM and the application area. For example, data filtering in Twitter is challenging because of its noisy nature, while in Amazon it is straightforward as the reviews are aggregated in the product's Web page. Detailed steps that are involved in this stage are described below.

Determination of time window

The time window is related to the *when* question as it specifies the duration of the collection activity as well as its relation to the characteristic period of the phenomenon. The characteristic period for product sales could be related to the new-product lifecycle (Liu et al., 2010), while for a disease outbreak to duration of pandemic stages (Ritterman et al., 2009). Clearly, the time window affects both the completeness and the sufficiency of the data.

Identification of location

The identification of location characterizing data is related to the *where* question. It is crucial in some phenomena (e.g. determination of natural phenomena occurrence) and thus accurate extraction of location is very important. The location characterizing SM data can be extracted from metadata (e.g. Lampos and Cristianini, 2012; Achrekar et al., 2011) or inferred from actual data.

Identification of user profile characteristics

The information related to the online profile of a user answers the *who* question. In a number of empirical studies (e.g. Forman et al., 2008; Skoric et al., 2012) it is suggested that this information is very important. For instance, Achrekar et al. (2011)

filter tweets from the same user within a certain syndrome elapsed time in order to avoid duplication from multiple encounters associated with a single episode of the illness.

Selection of search terms

The search terms selection step deals with the *what* question. In complex phenomena the identification of both the complete and correct set of search terms can be challenging. For example, Da et al. (2011) measured the search volume for 3,606 stocks through Google trends based on both “stock ticker” and “company name” and they, interestingly, identified that their correlation was only 9%. The inadequate completion of this task could result in poor quality data regarding its completeness and objectivity.

The different approaches for this step can fall into two broad categories: (a) *manual approaches* where researchers set search terms (e.g. Polgreen et al., 2008) and (b) *dynamic approaches* where search terms are derived through a computational process (e.g. Ginsberg et al., 2009). We should note that we consider the use of Google Trends’ as a dynamic selection approach since the resulting categories are determined based on Google’s natural language classification engine.

3.1.2 Stage 1.2: Computation of Predictor Variables

This stage deals with analysis of the raw data resulting from the previous stage in order to compute the values of predictor variables. In this stage, only variables related to SM are considered despite the fact that more variables (e.g. product price) can be finally employed in the predictive analysis stage. The steps composing this stage are the following:

Selection of predictor variables

Although a number of different variables have been used in the literature, we classify them into the following categories:

- *Volume-related variables*: these measure the amount of SM data in the form of number of tweets, number of reviews, number of queries etc.
- *Sentiment-related variables*: these measure the sentiment expressed through the data. The sentiment has been measured in the literature with the bullishness index (Oh and Sheng, 2011), review valence (Forman et al., 2008), review rating (Ghose and Ipeirotis, 2011), etc.
- *Profile characteristics of online users* such as Facebook friends (Franch, 2012), number of followers of users that posted a tweet (Rui and Whinston, 2011), total posts (Oh and Sheng, 2011), the location of the reviewer (Forman et al., 2008) and in-degree (Livne et al., 2011).

The proper selection of the variables that are employed in the analysis can influence the completeness of the data.

Measurement of predictor variables

The majority of variables are usually measured at successive time instants separated by uniform time intervals and are thus expressed as time series. The time intervals that have been used in the literature vary from hours to months. However, in some cases variables are measured just once hence resulting in one value per variable (e.g. Tumasjan et al., 2010).

Careful selection of measurement time intervals allows predictor variables to be comparable to the actual outcome. For instance, Forman et al. (2008) aggregated data by month because the evaluation data of the outcome was formed in monthly reports. However, in some cases the measurement of variables follows different time intervals than the actual outcome data (e.g. Tumasjan et al., 2010).

Computation of predictor variables

Although the computation of volume-related variables is straightforward and provides accurate results, the computation of sentiment expressed in text can be cumbersome and may provide poor results. Literature reveals that many research efforts have come up with poor sentiment analysis results (e.g. Gayo-Avello, 2011; Metaxas et al., 2011), mainly because of the informal and noisy nature of SM that creates problems to widely used NLP tools. The poor performance of sentiment analysis is a major source of weakness in the quality of data objectivity as the interpreted sentiment is different than that actually expressed.

In general the approaches used for sentiment computation can be categorized as follows: (a) *lexicon-based*, where sentiment is defined by the occurrence in the text of words included in a pre-defined lexicon (e.g. Metaxas et al., 2011; O'Connor et al., 2010) and (b) *machine learning*, where sentiment is computed by language model classifiers (e.g. Asur and Huberman, 2010).

3.2 Phase 2: Predictive Analysis

The aim of this phase is the creation and evaluation of a *predictive model* that will enable accurate prediction of phenomenon outcomes based on a new set of observations, where *new* can be interpreted as observations in future or observations that were not included in the original data sample.

Statisticians recognize that analyses aimed at prediction are different from those aimed at explanation (Konishi and Kitagawa, 2007). Predictive power refers to the ability of predicting new observations accurately, while explanatory power to the strength of association indicated by a statistical model. *'A statistically significant effect or relationship does not guarantee high predictive power, because the precision or magnitude of the causal effect might not be sufficient for obtaining levels of*

predictive accuracy that are practically meaningful' (Shmueli and Koppius, 2010: 561). Although statistically significant effects or relationships do not guarantee high predictive power, empirical studies that make predictive claims often infer predictive power from explanatory power without employing predictive analytics (Shmueli and Koppius, 2010).

3.2.1 Stage 2.1: Creation of Predictive Model

In this stage the actual model is created based on statistical or data mining methods. The steps that compose this stage are described below.

Selection of predictive method

The actual model of the predictive analysis is built based on different statistical or data mining methods. The most common method in literature is linear regression but many others have been also employed such as logistic regression (Livne et al., 2011), Markov models (Gruhl et al., 2005), neural networks (Bollen et al., 2011), support vector machine (Ritterman et al., 2009) and Granger causality (Gilbert and Karahalios, 2010).

Selection and use of non-SM predictor variables

Apart from the predictor variables computed through SM data, other predictor variables are also used in the predictive model. These usually express objective facts, such as past values of phenomenon outcomes and demographics. For instance, Forman et al. (2008) studied the relation between both the average valence of a review and the percentage of reviews disclosing real name or location, and product sales on Amazon. Towards this end, they also employed product price as a *control variable* in order to reduce the possibility that results reflect differences in average unobserved product quality rather than aspects of the reviews per se. In addition, Rui and Whinston (2011) employed non-SM predictor variables such as budget of a

movie or the fact that a movie is a sequel in order to enhance the accuracy of the model and Da et al. (2011) employed the number of news data from the Wall Street Journal in order to predict stock prices.

Identification of data for evaluation of prediction

The data referred to here represent the actual phenomenon outcome. This data is taken from official sources such as governmental documents and Web sites (e.g. Lamos and Cristianini, 2012; Sakaki et al., 2010; Ettredge et al., 2005), other trustworthy Web sites (e.g. Bollen et al., 2011), international organizations (e.g. O'Connor et al., 2010), etc. The accuracy and timely collection of this data is important for the creation of the predictive model.

3.2.2 Stage 2.2: Evaluation of the Predictive Performance

In this stage prediction accuracy is evaluated against the actual outcome. The steps that comprise this stage are described below.

Selection of the evaluation method

The evaluation of predictive performance is very important as it provides the actual result of the study as a whole. In the literature two different approaches are mainly employed: (a) *explanatory* analytics and (b) *predictive* analytics. The former assesses the statistical significance of the model using metrics such as p-values or R^2 (e.g. Asur and Huberman, 2010). The latter usually obtains out-of-sample data to be used for actual evaluation based on metrics such as out-of-sample error rate and statistics such as Predicted Residual Sums of Squares (e.g. Bordino et al., 2012), Root Mean Square Error (e.g. Achrekar et al., 2011), Mean Absolute Percentage Error (e.g. Bollen et al., 2011; Liu et al., 2007) and cross-validation summaries.

In general, the criteria that specify whether a study follows a predictive evaluation method or not are the following (Shmueli and Koppius, 2010):

- Was predictive accuracy based on out-of-sample assessment?
- Was predictive accuracy assessed with adequate predictive measures?

Specification of the prediction baseline

The baseline for prediction is an important element in the literature as it provides an extra metric for evaluating predictive power. The predictive power of an SM data based model is often judged in relation to statistical models fit with traditional data sources (e.g. Goel et al., 2010; Rui and Whinston, 2011) or past values (e.g. Bollen et al., 2011; Ritterman et al., 2009; Wu and Brynjolfsson, 2009). In addition, the results of prediction are sometimes also evaluated against prior models and approaches (e.g. Ghose and Ipeirotis, 2011).

4. Understanding the Predictive Power of Social Media

We now employ our framework in order to gain insight into the predictive power of Social Media (SM). We initially categorize the identified articles based on the application area studied (Table 1) and the type of SM employed (Table 2).

Table 1 The application areas studied in the literature

<i>Disease outbreaks</i>	Achrekar et al. (2011); Althouse et al. (2011); Culotta (2010); Ginsberg et al. (2009); Hulth et al. (2009); Polgreen et al. (2008); Ritterman et al. (2009); Signorini et al. (2011); Wilson and Brownstein (2009)
<i>Elections</i>	Franch (2012); Gayo-Avello (2011); He et al. (2012); Jin et al. (2010); Jungherr et al. (2012); Livne et al. (2011); Lui et al. (2011); Metaxas et al. (2011); Skoric et al. (2012); Tjong et al. (2012); Tumasjan et al. (2010); Tumasjan et al. (2012)
<i>Macroeconomics</i>	Choi and Varian (2012); Ettredge et al. (2005); Guzman (2011); O'Connor et al. (2010); Vosen and Schmidt (2011); Vosen and Schmidt (2012); Wang et al. (2012); Wu and Brynjolfsson (2009)
<i>Movies</i>	Asur and Huberman (2010); Bothos et al. (2010); Goel et al. (2010); Krauss et al. (2008); Liu et al. (2007); Liu et al. (2010); Mishne and Glance (2006); Rui and Whinston (2011)

This is a pre-print version of the following article:

Evangelos Kalampokis, Efthimios Tambouris and Konstantinos Tarabanis: Understanding the Predictive Power of Social Media. Accepted for publication in Internet Research, special issue on the Predictive Power of Social Media. 2013.

<i>Natural phenomena</i>	Earle et al. (2011); Lampos and Cristianini (2012); Sakaki et al. (2010)
<i>Product sales</i>	Choi and Varian (2012); Forman et al. (2008); Ghose and Ipeirotis (2011); Goel et al. (2010); Gruhl et al. (2005); Jin et al. (2010)
<i>Stock Market</i>	Antweiler and Frank (2004); Bollen et al. (2011); Bordino et al. (2012); Da et al. (2011); De Choudhury et al. (2008); Gilbert and Karahalios (2010); Oh and Sheng (2011); Zhang et al. (2011a); Zhang et al. (2011b)

Table 2 The Social Media analyzed in the literature

<i>Blogs</i>	De Choudhury et al. (2008); Franch (2012); Gilbert and Karahalios (2010); Gruhl et al. (2005); Liu et al. (2007); Mishne and Glance (2006)
<i>Web search</i>	Althouse et al. (2011); Bordino et al. (2012); Choi and Varian (2012); Da et al. (2011); Ettredge et al. (2005); Ginsberg et al. (2009); Goel et al. (2010); Guzman (2011); Hulth et al. (2009); Lui et al. (2011); Polgreen et al. (2008); Vosen and Schmidt (2011); Vosen and Schmidt (2012); Wilson and Brownstein (2009); Wu and Brynjolfsson (2009)
<i>Message boards</i>	Antweiler and Frank (2004); Bothos et al. (2010); Krauss et al. (2008); Liu et al. (2010); Oh and Sheng (2011)
<i>Reviews</i>	Bothos et al. (2010); Forman et al. (2008); Ghose and Ipeirotis (2011)
<i>Microblogs (Twitter and Facebook updates)</i>	Achrekar et al. (2011); Asur and Huberman (2010); Bollen et al. (2011); Bothos et al. (2010); Culotta (2010); Earle et al. (2011); Franch (2012); Gayo-Avello (2011); He et al. (2012); Jungherr et al. (2012); Lampos and Cristianini (2012); Livne et al. (2011); Lui et al. (2011); Metaxas et al. (2011); O'Connor et al. (2010); Oh and Sheng (2011); Ritterman et al. (2009); Rui and Whinston (2011); Sakaki et al. (2010); Signorini et al. (2011); Skoric et al. (2012); Tjong et al. (2012); Tumasjan et al. (2010); Tumasjan et al. (2012); Wang et al. (2012); Zhang et al. (2011a); Zhang et al. (2011b)
<i>Social multimedia (YouTube, Flickr)</i>	Franch (2012); Jin et al. (2010)

Table 3 presents classification of the literature according to the approach employed for selecting search terms, which is vital in Stage 1.1 of the framework. The table suggests that the vast majority of the studies employs manual selection methods.

Table 3 Classification of literature according to the approach for search term selection

<i>Manual selection</i>	Achrekar et al. (2011); Althouse et al. (2011); Asur and Huberman (2010); Bollen et al. (2011); Bordino et al. (2012); Da et al. (2011); De Choudhury et al. (2008); Ettredge et al. (2005); Franch (2012); Gayo-Avello (2011); Gruhl et al. (2005); Guzman (2011); He et al. (2012); Jungherr et al. (2012); Liu et al. (2007); Lui et al. (2011); Metaxas et al. (2011); Mishne and Glance (2006); O'Connor et al. (2010); Oh and Sheng (2011); Polgreen et al. (2008); Rui and Whinston (2011); Signorini et al. (2011); Skoric et al. (2012); Tjong et al. (2012); Tumasjan et al. (2010); Wilson and Brownstein (2009); Wu and Brynjolfsson (2009); Zhang et al. (2011a); Zhang et al. (2011b)
<i>Dynamic selection</i>	Choi and Varian (2012); Culotta et al. (2010); Ginsberg et al. (2009); Goel et al. (2010); Hulth et al. (2009); Lamos and Cristianini (2012); Ritterman et al. (2009); Sakaki et al. (2010); Vosen and Schmidt (2011); Wang et al. (2012)

In Table 4 the studies that involve sentiment analysis are aggregated and categorized according to the method they have employed. Selecting such a method is important in Stage 1.2 of the proposed framework. In this table we do not include studies that express the sentiment as review ratings since its measurement is straightforward.

Table 4 Classification of literature according to the text's sentiment analysis approach

<i>Lexicon-based</i>	Bollen et al. (2011); Gayo-Avello (2011); Liu et al. (2010); Metaxas et al. (2011); O'Connor et al. (2010); Zhang et al. (2011a); Zhang et al. (2012b)
<i>Machine Learning</i>	Antweiler and Frank (2004); Asur and Huberman (2010); Bothos et al. (2010); Gayo-Avello (2011); Gilbert and Karahalios (2010); He et al. (2012); Krauss et al. (2008); Liu et al. (2007); Mishne and Glance

	(2006); Oh and Sheng (2011); Rui and Whinston (2011)
--	--

Based on the criteria employed by Shmueli and Koppius (2010) we also classify (Table 5) literature according to the approach used to infer SM predictive power.

Table 5 Classification of literature according to the evaluation approach

<i>Explanatory evaluation</i>	Antweiler and Frank (2004); Asur and Huberman (2010); Bordino et al. (2012); Da et al. (2011); Ettredge et al. (2005); Forman et al. (2008); Gayo-Avello (2011); He et al. (2012); Jin et al. (2010); Jungherr et al. (2012); Krauss et al. (2008); Livne et al. (2011); Liu et al. (2010); Lui et al. (2011); Metaxas et al. (2011); Mishne and Glance (2006); Polgreen et al. (2008); Skoric et al. (2012); Tjong et al. (2012); Tumasjan et al. (2010); Wilson and Brownstein (2009); Zhang et al. (2011a); Zhang et al. (2011b)
<i>Predictive evaluation</i>	Achrekar et al. (2011); Althouse et al. (2011); Bollen et al. (2011); Bothos et al. (2010); Choi and Varian (2012); Culotta (2010); De Choudhury et al. (2008); Franch (2012); Ghose and Ipeirotis (2011); Gilbert and Karahalios (2010); Ginsberg et al. (2009); Goel et al. (2010); Gruhl et al. (2005); Guzman (2011); Hulth et al. (2009); Lampos and Cristianini (2012); Liu et al. (2007); O'Connor et al. (2010); Oh and Sheng (2011); Ritterman et al. (2009); Rui and Whinston (2011); Sakaki et al. (2010); Signorini et al. (2011); Vosen and Schmidt (2011); Vosen and Schmidt (2012); Wang et al. (2012); Wu and Brynjolfsson (2009)

Finally, Table 6 categorizes literature according to their final outcome with regard to the predictive power of SM. Some studies provide evidence for both outcomes. These are included in both categories.

Table 6 Classification of literature based on main outcome

<i>Support SM predictive power</i>	Achrekar et al. (2011); Althouse et al. (2011); Antweiler and Frank (2004); Asur and Huberman (2010); Bollen et al. (2011); Bordino et al. (2012); Bothos et al. (2010); Choi and Varian (2012); Culotta (2010); Da et al. (2011); De Choudhury et al. (2008); Ettredge et al. (2005); Forman et al. (2008); Franch (2012); Ghose and Ipeirotis (2011);
------------------------------------	---

	Gilbert and Karahalios (2010); Ginsberg et al. (2009); Goel et al. (2010); Gruhl et al. (2005); Guzman (2011); Hulth et al. (2009); Jin et al. (2010); Krauss et al. (2008); Lampos and Cristianini (2012); Liu et al. (2007); Liu et al. (2010); Livne et al. (2011); Oh and Sheng (2011); Polgreen et al. (2008); Ritterman et al. (2009); Rui and Whinston (2011); Sakaki et al. (2010); Signorini et al. (2011); Tjong et al. (2012); Tumasjan et al. (2010); Vosen and Schmidt (2011); Vosen and Schmidt (2012); Wang et al. (2012); Wu and Brynjolfsson (2009); Zhang et al. (2011a); Zhang et al. (2011b)
<i>Challenge SM predictive power</i>	Bollen et al. (2011); Forman et al. (2008); Gayo-Avello (2011); Goel et al. (2010); He et al. (2012); Jungherr et al. (2012); Liu et al. (2010); Lui et al. (2011); Metaxas et al. (2011); Mishne and Glance (2006); O'Connor et al. (2009); Skoric et al. (2012); Tjong et al. (2012); Wilson and Brownstein (2009)

By synthesizing Tables 1-6 we can further analyze the empirical studies in the literature and make some interesting observations.

Search term selection

Table 3 suggests that although dynamic search term selection is used in most application areas (Table 1), it only appears in studies that employ *Web search* and *microblog data* (Table 2). Furthermore, all these studies support SM predictive power (Table 6) based on predictive analytics (Table 5). In the case of manual search term selection when considering the same two SM categories, the percentage of studies that support SM predictive power falls off to fifty percent (50%). Hence, we can conclude that search term selection is of vital importance in microblog and Web search data, and thus these SM categories call for sophisticated search terms selection methods. For instance, Lampos and Cristianini (2012) successfully estimated daily rainfall rates for five UK cities by identifying relevant tweets through the application of Bolasso

(i.e. the bootstrapped version of Least Absolute Shrinkage and Selection Operator) for search term selection.

Sentiment analysis

Table 4 suggests that the majority of studies that employ sentiment analysis investigate stock market and movies (Table 1). Although sentiment seems to be important in application areas such as elections, product sales and macroeconomics, only six (6) out of twenty four (24) studies include a sentiment-related independent variable. Disease outbreaks and natural phenomena related studies do not employ sentiment, as one might have expected. Interestingly however, forty per cent (40%) of studies that have used sentiment-related variables challenge SM predictive power. This number increases to sixty five percent (65%) in the case of *lexicon-based* approaches, while it falls off to twenty percent (20%) in those of *machine learning*. Hence, it seems that sentiment analysis in SM requires innovative approaches that could address the noisy and informal nature of SM.

Evaluation method

In general, half of the studies do not use predictive analytics to draw conclusions on the predictive performance of SM. These studies span equally across all SM categories (Table 2). With regard to application areas (Table 1), the vast majority of election-related cases do not follow a predictive analytics evaluation, while most studies related to macroeconomic indices, natural phenomena and product sales application areas evaluate predictive power based on prediction analytics. The evaluation of a predictive model with out-of-sample data is sometimes challenging. For instance, in the case of election-related studies the outcome is produced once every four or five years. In order to overcome this limitation Franch (2012) used poll data.

Tables 5 and 6 suggest that ten (10) out of fourteen (14) studies that challenge SM predictive power have used explanatory evaluation methods. This fact does not imply that these studies do not contribute to the understanding of SM predictive power as lack of a statistically significant relationship indicates low predictive power. In addition, fourteen (14) out of forty (40) studies that support SM predictive power infer predictive power without employing predictive analytics. Here we should also note that if these studies had used predictive evaluation methods, they could have presented high predictive power. However, based on the reported results we cannot assess their predictive power because a statistically significant relationship does not always ensure high predictive power. For example, *'low predictive power can result from over-fitting, where an empirical model fits the training data so well that it underperforms in predicting new data'* (Breiman, 2001: 204).

Application areas

The application area of a study seems to be related to the accuracy of the prediction that the study presents. Some application areas, such as disease outbreak and natural phenomena, do not involve the expression of any kind of opinion or sentiment. The signal that the researcher has to decode in these cases has to do with the occurrence or not of the event. As a result, these studies are expected to provide more accurate predictions than studies requiring extracting opinions or sentiment out of raw data. Moreover, some application areas, such as elections or macroeconomics, can be characterized as complex because they involve multiple and interrelated real-world entities such as political parties and politicians or complex concepts such as consumer confidence or inflation rate. The identification of the complete set of relevant raw SM data in these cases is challenging and hence call for sophisticated methods.

This becomes evident if we elaborate on two of the identified applications areas, namely elections and disease outbreak. The former involves opinion expression and is characterized by multiple and interrelated real-world entities (i.e. political parties, candidates, election constituencies), while the latter does not require opinion extraction. Table 1 suggests that all eleven (11) election-related studies selected their search terms manually (Table 3) and only three of them employed sentiment-related variables (Table 4). These facts could provide an explanation of the unfavourable and controversial results reported in the literature regarding predictability of election results through SM. In addition, half of the disease outbreak related studies employed sophisticated search unit selection approaches, eighty percent (80%) used predictive analytics evaluation and ninety percent (90%) supported SM predictive power.

5. Conclusions

Social media (SM) are a vital component of the contemporary Web as they enable the production of data that reflects personal opinions, thoughts and behaviour. Since the emergence of blogs and forums, several research efforts have explored the potential of SM data for *predictions of outcomes* such as disease outbreaks, product sales, stock market volatility, and elections. As the field is immature, some studies produce controversial results and doubtful outcomes.

In this article, we aim at consolidating knowledge created in the past eight years by empirical studies that aim at predicting real world outcomes through SM, thus enabling an in-depth understanding of SM predictive power. Towards this end, we identify and synthesize the literature and we create a SM data analysis conceptual framework for predictions. Using this framework we further analyze the literature and classify studies according to the approaches they follow and the results they report.

The proposed framework suggests that all relevant studies can be decomposed into a small number of *steps* and that different choices can be made in each step. The application of the framework enabled us to make some interesting observations. The majority of the empirical studies support SM predictive power, however more than one-third of these studies infer predictive power without employing predictive analytics. Sophisticated search term selection is crucial in Web search and microblog data. In addition, the use of sentiment-related variables resulted often in controversial outcomes proving that SM data call for sophisticated sentiment analysis approaches.

We anticipate that both the framework and analysis results will enable researchers to design scientifically rigorous new studies and to more easily identify the field's weaknesses, hence proposing new future research directions.

References

1. Breiman, L. (2001), "Statistical Modeling: The Two Cultures", *Statistical Science* Vol. 16, No. 3, pp. 199-215.
2. Corley, C.D., Cook, D.J., Mikler, A.R. and Singh, K.P. (2010), "Text and Structural Data Mining of Influenza Mentions in Web and Social Media", *International Journal of Environmental Research and Public Health*, Vol. 7, pp. 596-615.
3. Chevalier, J.A. and Mayzlin, D. (2006), "The effect of Word of Mouth on Sales: Online Book Reviews", *Journal of Marketing Research*, Vol. 43, No. 3, pp. 345-354.
4. Chen, Y., Wang, Q. and Xie, J. (2011), "Online Social Interactions: A Natural Experiment on Word of Mouth Versus Observational Learning", *Journal of Marketing Research*, Vol. 48, No. 2, pp. 238-254.

5. Chunara, R., Andrews, J. R. and Brownstein, J. S. (2012), "Social and News Media Enable Estimation of Epidemiological Patterns Early in the 2010 Haitian Cholera Outbreak", *The American Journal of Tropical Medicine and Hygiene*, Vol. 86, No. 1, pp. 39-45.
6. Duan W., Gu, B. and Whinston, A.B. (2008), "Do online reviews matter? - An empirical investigation of panel data", *Decision Support Systems*, Vol. 45, No. 4, pp. 1007-1016.
7. Gilbert, E. and Karahalios, K. (2009), "Predicting Tie Strength With Social Media", in *27th International Conference on Human Factors in Computing Systems - CHI 2009*, pp. 211-220.
8. Konishi, S. and Kitagawa, G. (2007), *Information Criteria and Statistical Modeling*, Springer, New York.
9. Louis, C.St. and Zorlu, G. (2012), "Can Twitter predict disease outbreaks?" *British Medical Journal*, 344:e2353
10. Morales-Arroyo, M. and Pandey, T. (2010), "Identification of Critical eWOM Dimensions for Music Albums", in *IEEE International Conference on Management of Innovation and Technology*, IEEE, pp. 1230-1235.
11. Oghina, A., Breuss, M., Tsagkias, M. and de Rijke, M. (2012), "Predicting IMDB Movie Rating Using Social Media", in *34th European Conference on Information Retrieval (ECIR 2012)*, Springer, pp. 503-507.
12. Reinstein, D. and Snyder, C.M. (2005), "The Influence of Expert Reviews on Consumer Demand for Experience Goods: A Case Study of Movie Critics", *Journal of Industrial Economics*, Vol. 53, No. 1, pp. 27-51.
13. Shmueli, G. (2010), "To Explain or to Predict?", *Statistical Science*, Vol. 25, No.3, pp. 289-310.

14. Shmueli, G. and Koppius, O.R. (2010), "Predictive Analytics in Information Systems Research", *MIS Quarterly*, Vol. 35, No. 3, pp. 553-572.
15. Strong, D.M., Lee, Y.W. and Wang, R.Y. (1997), "Data Quality in Context", *Communications of the ACM*, Vol. 40, No.5, pp. 103-110.
16. Tsagkias, E., Weerkamp, W. and de Rijke, M. (2010), "News comments: Exploring, modeling, and online predicting", in *ECIR 2010*, Springer, pp. 191–203.
17. Webster, J. and Watson, R.T. (2002), "Analyzing the Past to Prepare for the Future: Writing a Literature Review", *MIS Quarterly*, Vol. 26, No2, pp. xiii-xxiii.
18. Ye, Q., Law, R. and Gu, B. (2009), "The impact of online user reviews on hotel room sales", *International Journal of Hospitality Management*, Vol. 28, pp. 180-182.

Literature Review References

19. Achrekar, H., Gandhe, A., Lazarus, R., Yu, S-H. and Liu, B. (2011), "Predicting Flu Trends using Twitter Data" in 2011 *IEEE Conference on Computer Communications Workshops*, IEEE, pp. 702-707.
20. Althouse, B.M., Ng, Y.Y. and Cummings, D.A.T. (2011), "Prediction of Dengue Incidence Using Search Query Surveillance", *Public Library of Science*, Vol. 5, No. 8, pp. 1-7.
21. Antweiler, W. and Frank, M.Z. (2004), "Is all that talk just noise? the information content of internet stock message boards" *Journal of Finance*, Vol. 59, No. 3, pp. 1259-1294.
22. Asur, S. and Huberman, B.A. (2010), "Predicting the Future With Social Media", in *2010 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology*, IEEE Press, pp. 492-499.

23. Bollen, J., Mao, H. and Zeng, X.J. (2011), "Twitter mood predicts the stock market", *Journal of Computational Science*, Vol. 2, No. 1, pp. 1-8.
24. Bordino, I., Battiston, S., Caldarelli, G., Cristelli, M., Ukkonen, A. and Weber, I. (2012), "Web search queries can predict stock market volumes", *PLoS ONE*, Vol. 7, No. 7, pp. e40014.
25. Bothos, E., Apostolou, D. and Mentzas, G. (2010), "Using Social Media to Predict Future Events with Agent-Based Markets", *IEEE Intelligent Systems*, Vol. 25, No. 6, pp. 50-58.
26. Choi, H. and Varian, H. (2012), "Predicting the Present with Google Trends", *The Economic Record*, Vol. 88, pp. 2-9.
27. Culotta, A. (2010), "Towards detecting influenza epidemics by analyzing Twitter messages", in *First Workshop on Social Media Analytics*, ACM Press, pp. 115-122.
28. Da, Z., Engelberg, J., and Gao, P. (2011), "In Search of Attention", *The Journal of Finance*, Vol. 66, No. 5, pp. 1461-1499.
29. De Choudhury, M., Sundaram, H., John, A., and Seligmann, D.D. (2008), "Can blog communication dynamics be correlated with stock market activity?" in *19th ACM Conference on Hypertext and Hypermedia*, ACM Press, pp. 55-60.
30. Earle, P., Bowden, D.C. and Guy, M. (2011), "Twitter earthquake detection: earthquake monitoring in a social world", *Annals of Geophysics*, Vol. 54, No. 6, pp. 708-715.
31. Ettredge, M., Gerdes, J. and Karuga, G. (2005), "Using web-based search data to predict macro-economic statistics", *Communications of the ACM*, Vol. 48, pp. 87-92.

32. Forman, C., Ghose, A. and Wiesenfeld, B. (2008), "Examining the Relationship Between Reviews and Sales: The Role of the Reviewer Identity Disclosure in Electronic Markets", *Information Systems Research*, Vol. 19, No. 3, pp. 291-313.
33. Franch, F. (2012), "(Wisdom of the Crowds)²: 2010 UK Election Prediction with Social Media", *Journal of Information Technology & Politics*, DOI: 10.1080/19331681.2012.705080
34. Gayo-Avello D. (2011), "Don't Turn Social Media Into Another 'Literary Digest' Poll" *Communications of the ACM*, Vol. 54, No. 10, pp. 121-128.
35. Ghose, A. and Ipeirotis, P.G. (2011), "Estimating the Helpfulness and Economic Impact of Product Reviews: Mining Text and Reviewer Characteristics", *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, Vol. 23, No. 10, pp. 1498-1512.
36. Gilbert, E. and Karahalios, K. (2010), "Widespread worry and the stock market" in *Fourth International Conference on Weblogs and Social Media*, AAAI Press, pp. 58-65.
37. Ginsberg, J., Mohebbi, M. H., Patel, R. S., Brammer, L., Smolinski, M. S. and Brilliant, L. (2009), "Detecting influenza epidemics using search engine query data", *Nature*, Vol. 457, No. 7232, pp. 1012-4.
38. Goel, S., Hofman, J.M., Lahaie, S., Pennock, D.M. and Watts, D.J. (2010), "Predicting consumer behaviour with Web search", *Proceedings of the National Academy of Sciences (PNAS)*, Vol. 107, No. 41, pp. 17486-17490.
39. Gruhl, D., Guha, R., Kumar, R., Novak, J. and Tomkins, A. (2005), "The Predictive Power of Online Chatter", in *Eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, ACM Press, pp. 78-87.

40. Guzman, G. (2011), "Internet search behavior as an economic forecasting tool: The case of inflation expectations", *Journal of Economic and Social Measurement*, Vol. 36, No. 3, pp. 119-167.
41. He, Y., Saif, H., Wei, Z. and Wong, K. (2012), "Quantising Opinions for Political Tweets Analysis", in *Eight International Conference on Language Resources and Evaluation*, European Language Resources Association, pp. 3901-3906.
42. Hulth, A., Rydevik, G. and Linde, A. (2009), "Web Queries as a Source for Syndromic Surveillance" *PLoS ONE*, Vol. 4, No. 2, pp. e4378.
43. Jin, X., Gallagher, A., Cao, L., Luo, J. and Han, J. (2010), "The Wisdom of Social Multimedia: Using Flickr For Prediction and Forecast" in *ACM Multimedia 2010*, ACM Press, pp. 1235-1244.
44. Jungherr, A., Jurgens, P., and Schoen, H. (2012), "Why the Pirate Party Won the German Election of 2009 or The Trouble With Predictions: A Response to Tumasjan, A., Sprenger, T.O., Sander, P.G., & Welpe, I.M. 'Predicting Elections With Twitter: What 140 Characters Reveal About Political Sentiment'", *Social Science Computer Review*, Vol. 30, No. 2, pp. 229-234.
45. Krauss, J., Nann, S., Simon, D., Fischbach, K. and Gloor, P. (2008), "Predicting Movie Success and Academy Awards through Sentiment and Social Network Analysis" in *16th European Conference on Information Systems*, pp. 2026-2037.
46. Lampos, V. and Cristianini, N. (2012), "Nowcasting Events from the Social Web with Statistical Learning", *ACM Transactions on Intelligent Systems and Technology*, Vol. 3, No. 4.
47. Liu, Y., Chen, Y., Lusch, R. F., Chen, H., Zimbra, D. and Zeng, S. (2010), "User-Generated Content on Social Media: Predicting Market Success with Online Word-of-Mouth", *IEEE Intelligent Systems*, Vol. 25, No. 1, pp. 75-78.

48. Liu, Y., Huang, X., An, A. and Yu, X. (2007), "ARSA: A sentiment-aware model for predicting sales performance using blogs" in *30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, ACM Press, pp. 607-614.
49. Livne, A., Simmons, P.S., Adar, E. and Adamic, L.A. (2011), "The Party is Over Here: Structure and Content in the 2010 Election" in *Fifth International AAAI Conference on Weblogs and Social Media*, AAAI Press, pp. 201-208.
50. Lui, C., Metaxas, P.T. and Mustafaraj, E. (2011), "On the predictability of the U.S. Elections through search volume activity" in *IADIS International Conference e-Society 2011*, pp. 165-172.
51. Metaxas P. T., Mustafaraj, E. and Gayo-Avello, D. (2011), "How (Not) To Predict Election" in *2011 IEEE Third International Conference on Social Computing*, IEEE, pp.165-171.
52. Mishne, G. and Glance, N. (2006), "Predicting Movie Sales from Blogger Sentiment" in *American Association for Artificial Intelligence 2006 Spring Symposium on Computational Approaches to Analysing Weblogs*
53. O'Connor, B., Balasubramanyan R., Routledge, B.R., and Smith, N.A. (2010), "From Tweets to Polls: Linking Text Sentiment to Public Opinion Time Series", in *Proceedings of the International AAAI Conference on Weblogs and Social Media*, AAAI Press, pp. 122-129.
54. Oh, C. and Sheng, O. (2011), "Investigating Predictive Power of Stock Micro Blog Sentiment in Forecasting Future Stock Price Directional Movement", in *32nd International Conference on Information Systems*, AIS, p.17.

55. Polgreen, P.M., Chen, Y., Pennock, D.M. and Nelson, F.D. (2008), "Using Internet Searches for Influenza Surveillance", *Clinical Infectious Diseases*, Vol.47, No.11, pp. 1443-1448.
56. Ritterman, J., Osborne, M. and Klein, E. (2009), "Using Prediction Markets and Twitter to Predict a Swine Flu Pandemic" in *First International Workshop on Mining Social Media*, pp. 9-17.
57. Rui, H. and Whinston, A. (2011), "Designing a Social-Broadcasting-Based Business Intelligence System", *ACM Transactions on Management Information Systems*, Vol. 2, No. 4, pp.22.
58. Sakaki, T., Okazaki, M. and Matsuo, Y. (2010), "Earthquake shakes twitter users: real-time event detection by social sensors" in 19th International Conference on World Wide Web (WWW'10), ACM Press, pp. 851-860.
59. Signorini, A., Segre, A.M. and Polgreen, P.M. (2011), "The Use of Twitter to Track Levels of Disease Activity and Public Concern in the U.S. during the Influenza A H1N1 Pandemic", *PLoS ONE*, Vol. 6, No. 5, pp. e19467.
60. Skoric, M., Poor, N., Achananuparp, P., Lim, E. and Jiang J. (2012), "Tweets and Votes: A Study of the 2011 Singapore General Election", in 45th Hawaii International Conference on System Sciences, IEEE, pp. 2583-2591.
61. Tjong, E., Sang, K. and Bos, J. (2012), "Predicting the 2011 Dutch Senate Election Results with Twitter", in 13th Conference of the European Chapter of the Association for Computational Linguistics, Association for Computational Linguistics, pp. 53-60.
62. Tumasjan, A., Sprenger, T.O., Sandner, P.G., and Welpe, I.M. (2010), "Predicting Elections with Twitter: What 140 Characters Reveal about Political Sentiment", in

Fourth International AAAI Conference on Weblogs and Social Media, AAAI Press, pp. 178-185.

63. Tumasjan, A., Sprenger, T.O., Sandner, P.G. and Welpe, I.M. (2012), "Where There is a Sea There are Pirates: Response to Jungherr, Jurgens, and Schoen", *Social Science Computer Review*, Vol. 30, No. 2, pp. 235-239.
64. Vosen, S. and Schmidt, T. (2011), "Forecasting Private Consumption: Survey-Based Indicators vs. Google Trends", *Journal of Forecasting*, Vol. 30, No. 6, pp. 565-578.
65. Vosen, S. and Schmidt, T. (2012), "A monthly consumption indicator for Germany based on Internet search query data", *Applied Economics Letters*, Vol. 19, No. 7, pp. 683-687.
66. Wang, X., Gerber, M.S. and Brown, D., E. (2012), "Automatic Crime Prediction Using Events Extracted from Twitter Posts" in S.J. Yang, A.M. Greengerg and M. Endsley (Eds.): SBP 2012, LNCS 7227, pp. 231-238, Springer-Verlag Berlin Heidelberg.
67. Wilson, K. and Brownstein, J.S. (2009), "Early detection of disease outbreaks using the Internet", *Canadian Medical Association Journal*, Vol. 180, No.8, pp. 829-831.
68. Wu, L. and Brynjolfsson, E. (2009), "The Future of Prediction: How Google Searches Foreshadow Housing Prices and Quantities" in *30th International Conference on Information Systems*, AISE, <http://aisel.aisnet.org/icis2009/147>
69. Zhang, X., Fuehres, H. and Gloor, P.A. (2011a), "Predicting Stock Market Indicators Through Twitter 'I hope it is not as bad as I fear'", *Procedia - Social and Behavioral Sciences*, Vol. 26, pp. 55-62.

This is a pre-print version of the following article:

Evangelos Kalampokis, Efthimios Tambouris and Konstantinos Tarabanis: Understanding the Predictive Power of Social Media. Accepted for publication in Internet Research, special issue on the Predictive Power of Social Media. 2013.

70. Zhang, X., Fuehres, H. and Gloor, P.A. (2011b), "Predicting Asset Value through Twitter Buzz", J. Altmann et al. (Eds): Advances in Collective Intelligence 2011, AISC 113, pp. 23-34.