

# Yahoo! for Amazon: Sentiment Extraction from Small Talk on the Web

Sanjiv R. Das

Department of Finance, Leavey School of Business, Santa Clara University,  
Santa Clara, California 95053, srdas@scu.edu

Mike Y. Chen

Ludic Labs, San Mateo, California 94401, mike@ludic-lab.com

Extracting sentiment from text is a hard semantic problem. We develop a methodology for extracting small investor sentiment from stock message boards. The algorithm comprises different classifier algorithms coupled together by a voting scheme. Accuracy levels are similar to widely used Bayes classifiers, but false positives are lower and sentiment accuracy higher. Time series and cross-sectional aggregation of message information improves the quality of the resultant sentiment index, particularly in the presence of slang and ambiguity. Empirical applications evidence a relationship with stock values—tech-sector postings are related to stock index levels, and to volumes and volatility. The algorithms may be used to assess the impact on investor opinion of management announcements, press releases, third-party news, and regulatory changes.

*Key words:* text classification; index formation; computers-computer science; artificial intelligence; finance; investment

*History:* Accepted by David A. Hsieh, finance; received May 4, 2004. This paper was with the authors 1 year and 2 weeks for 1 revision. Published online in *Articles in Advance* July 20, 2007.

## 1. Introduction

Language is itself the collective art of expression, a summary of thousands upon thousands of individual intuitions. The individual gets lost in the collective creation, but his personal expression has left some trace in a certain give and flexibility that are inherent in all collective works of the human spirit—Edward Sapir, cited in *Society of Mind* by Minsky (1985, p. 270).

We develop hybrid methods for extracting opinions in an automated manner from discussions on stock message boards, and analyze the performance of various algorithms in this task, including that of a widely used classifier available in the public domain. The algorithms are used to generate a sentiment index and we analyze the relationship of this index to stock values. As we will see, this analysis is efficient, and useful relationships are detected.

The volume of information flow on the Web has accelerated. For example, in the case of Amazon Inc., there were cumulatively 70,000 messages by the end of 1998 on Yahoo's message board, and this had grown to about 900,000 messages by the end of 2005. There are almost 8,000 stocks for which message board activity exists, across a handful of message board providers. The message flow comprises valuable insights, market sentiment, manipulative behavior, and reactions to other sources of news. Message

boards have attracted the attention of investors, corporate management, and of course, regulators.<sup>1</sup>

In this paper, "sentiment" takes on a specific meaning, that is, the net of positive and negative opinion expressed about a stock on its message board. Hence, we specifically delineate our measure from other market conventions of sentiment such as deviations from the expected put-call ratio. Our measure is noisy because it comprises information, sentiment, noise, and estimation error.

Large institutions express their views on stocks via published analyst forecasts. The advent of stock chat and message boards enables small investors to express their views too, frequently and forcefully. We show that it is possible to capture this sentiment using *statistical language techniques*. Our algorithms are validated using revealed sentiment on message boards, and from the statistical relationship between sentiment and stock returns, which track each other.

<sup>1</sup>Das et al. (2005) present an empirical picture of the regularities found in messages posted to stock boards. The recent case of Emulex Corp. highlights the sensitivity of the Internet as a sentiment channel. Emulex's stock declined 62% when an anonymous, false news item on the Web claimed reduced earnings and the resignation of the CEO. The Securities Exchange Commission (SEC) promptly apprehended the perpetrator, a testimony to the commitment of the SEC to keeping this sentiment channel free and fair. In relation to this, see the fascinating article on the history of market manipulation by Leinweber and Madhavan (2001).

Posted messages offer opinions that are bullish, bearish, and many that are confused, vitriolic, rumor, and spam (null messages). Some are very clear in their bullishness, as is the following message on Amazon's board (Msg 195006):

The fact is . . . .

The value of the company increases because the leader (Bezos) is identified as a commodity with a vision for what the future may hold. He will now be a public figure until the day he dies. That is value.

In sharp contrast, this message was followed by one that was strongly bearish (Msg 195007):

Is it famous on infamous? A commodity dumped below cost without profit, I agree. Bezos had a chance to make a profit without sales tax and couldn't do it. The future looks grim here.

These (often ungrammatical) opinions provide a basis for extracting small investor sentiment from discussions on stock message boards.

While financial markets are just one case in point, the Web has been used as a medium for information extraction in fields such as voting behavior, consumer purchases, political views, quality of information equilibria, etc. (see Godes and Mayzlin 2004, Lam and Myers 2001, Wakefield 2001, Admati and Pfleiderer 2000 for examples). In contrast to older approaches such as investor questionnaires, sentiment extraction from Web postings is relatively new. It constitutes a real-time approach to sentiment polling, as opposed to traditional point-in-time methods.

We use statistical and natural language processing techniques to elicit emotive sentiment from a posted message; we implement five different algorithms, some language dependent, others not, using varied parsing and statistical approaches. The methodology used here has antecedents in the text classification literature (see Koller and Sahami 1997, Chakrabarti et al. 1998). These papers classify textual content into natural hierarchies, a popular approach employed by Web search engines.

Extracting the *emotive* content of text, rather than *factual* content, is a complex problem. Not all messages are unambiguously bullish or bearish. Some require context, which a human reader is more likely to have, making it even harder for a computer algorithm with limited background knowledge. For example, consider the following from Amazon's board (Msg 195016):

You're missing this Sonny, the same way the cynics pronounced that "Gone with the Wind" would be a total bust.

Simple, somewhat ambiguous messages like this also often lead to incorrect classification even by human

subjects. We analyze the performance of various algorithms in the presence of ambiguity, and explore approaches to minimizing its impact.

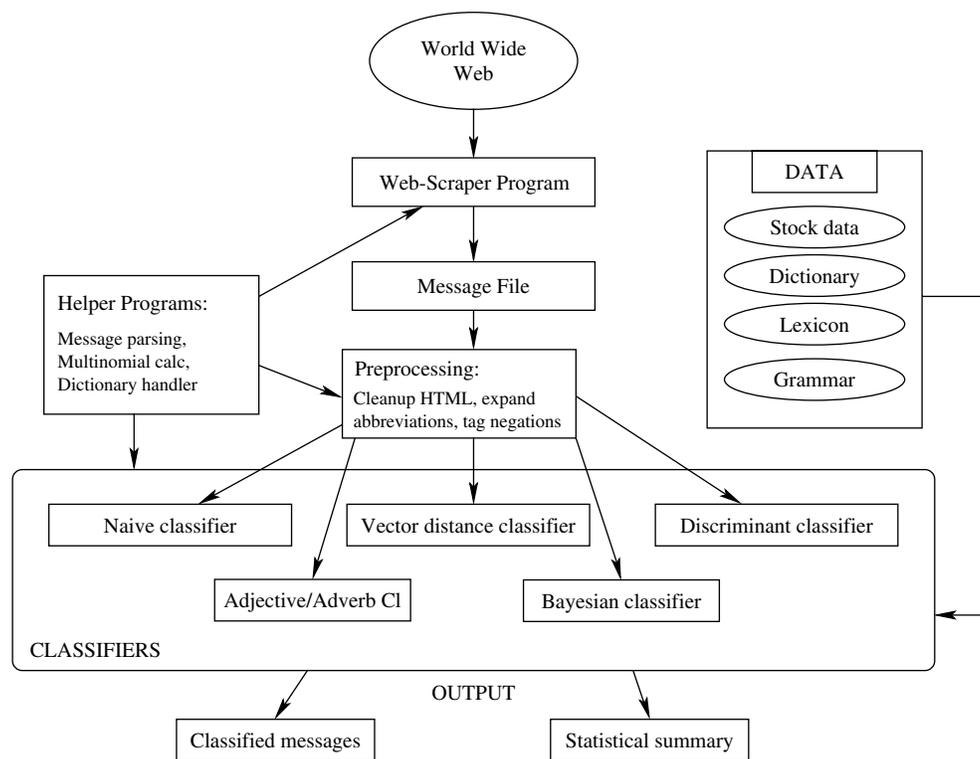
The technical contribution of this paper lies in the coupling of various classification algorithms into a system that compares favorably with standard Bayesian approaches, popularized by the phenomenal recent success of spam-filtering algorithms. We develop metrics to assess algorithm performance that are well suited to the finance focus of this work. There are unique contributions within the specific algorithms used as well as accuracy improvements overall, most noticeably in the reduction of false positives in sentiment classification. An approach for filtering ambiguity in known message types is also devised and shown to be useful in characterizing algorithm performance.

Recent evidence suggests a link between small investor behavior and stock market activity. Noticeably, day-trading volume has spurted.<sup>2</sup> Choi et al. (2002) analyze the impact of a Web-based trading channel on the trading activity in corporate 401(k) plans, and find that the "Web effect" is very large—trading frequency doubles, and portfolio turnover rises by over 50%, when investors are permitted to use the Web as an information and transaction channel. Wysocki (1998), using pure message counts, reports that variation in daily message posting volume is related to news and earnings announcements. Lavrenko et al. (2000) use computer algorithms to identify news stories that influence markets, and then trade successfully on this information. Bagnoli et al. (1999) examine the predictive validity of whisper forecasts, and find them to be superior to those of First Call (Wall Street) analysts.<sup>3</sup> Antweiler and Frank (2004) examine the bullishness of messages, and find that while Web talk does not predict stock movements, it is predictive of volatility. Tumarkin and Whitelaw (2001) also find similar results using self-reported sentiment (not message content) on the Raging Bull message board. Antweiler and Frank (2002) argue that message posting volume is a priced factor, and higher posting activity presages high volatility and poor returns. Tetlock (2005) and Tetlock et al. (2006) show that negative sentiment from these boards may be predictive of future downward moves in firm values.

<sup>2</sup> Stone (2001) cites a Bear Stearns report that reports a huge spurt in volume, and a total number of day-traders in excess of 50,000.

<sup>3</sup> The "whisper" number, an aggregate of informal earnings forecasts self-reported by individual investors, is now watched extensively by market participants, large and small. Whispers are forecasts of the quarterly earnings of a firm posted to the Web by individuals in a voluntary manner. The simple average of these forecasts is presented on the whisper Web page, along with the corresponding forecast from First Call, which is an aggregate of the sentiment of Wall Street analysts.

Figure 1 Schematic of the Algorithms and System Design Used for Sentiment Extraction



These results suggest the need for algorithms that can rapidly access and classify messages with a view to extracting sentiment—the goal of this paper.<sup>4</sup> The illustrative analyses presented in this paper confirm many of these prior empirical findings, and extend them as well.

Overall, this paper comprises two parts: (i) methodology and validation, in §2, which presents the algorithms used and their comparative performance, and (ii) the empirical relationship of market activity and sentiment, in §3. Section 4 contains discussion and conclusions.

## 2. Methodology

### 2.1. Overview

The first part of the paper is the extraction of opinions from message board postings to build a sentiment index. Messages are classified by our algorithms into one of three types: bullish (optimistic), bearish (pessimistic), and neutral (comprising either spam or messages that are neither bullish nor bearish). We use five algorithms, each with different conceptual underpinnings, to classify each message. These comprise a

blend of language features such as parts of speech tagging, and more traditional statistical methods.<sup>5</sup> Before initiating classification, the algorithms are tuned on a training corpus, i.e., a small subset of preclassified messages used for training the algorithms.<sup>6</sup> The algorithms “learn” sentiment classification rules from the preclassified data set, and then apply these rules out-of-sample. A simple majority across the five algorithms is required before a message is finally classified, or else it is discarded. This *voting approach* results in a better signal to noise ratio for extracting sentiment.

Figure 1 presents the flowchart for the methodology and online Appendix A (provided in the

<sup>4</sup> In contrast, Antweiler and Frank (2005) recently used computational linguistic algorithms to sort news stories into topics, instead of sentiment, and uncovered many interesting empirical regularities relating news stories and stock values.

<sup>5</sup> This paper complements techniques such as support vector machines (SVMs) that are optimization methods that classify content. See the papers by Vapnik (1995), Vapnik and Chervonenkis (1964), and Joachims (1999) for a review. A recent paper by Antweiler and Frank (2004) uses SVMs to carry out an exercise similar to the one in this paper. These approaches are computationally intensive and are often run on parallel processors. Moreover, they have been used for more than 30 years, and the technology is well developed. In this paper, we did not employ support vector machines, instead choosing to focus on purely analytic techniques that did not require optimization methods in the interests of computational efficiency.

<sup>6</sup> The training corpus is kept deliberately small to avoid overfitting, which is a common ailment of text classification algorithms.

e-companion)<sup>7</sup> contains technical details. The sequence of tasks is as follows. We use a “Web-scraper” program to download messages from the Internet, which are fed to the five classification algorithms to categorize them as buy, sell, or null types. Three supplementary databases support the classification algorithms.

- First, an electronic English “dictionary,” which provides base language data. This comes in handy when determining the nature of a word, i.e., noun, adjective, adverb, etc.

- Second, a “lexicon” which is a hand-picked collection of finance words (such as bull, bear, uptick, value, buy, pressure, etc.). These words form the variables for statistical inference undertaken by the algorithms. For example, when we count positive and negative words in a message, we will use only words that appear in the lexicon, where they have been pre-tagged for sign.

- Third, the “grammar” or the preclassified training corpus. It forms the base set of messages for further use in classification algorithms. These preclassified messages provide the in-sample statistical information for use on the out-of-sample messages.

These three databases (described in the online appendix) are used by five algorithms (denoted “classifiers”) to arrive at the three-way classification of each message.

## 2.2. Classifiers

Each of our five classifier algorithms relies on a different approach to message interpretation. Some of them are language independent, and some are not. Each approach is intuitive. They are all analytical, and do not require any lengthy optimization or convergence issues, hence they are computationally efficient, making feasible the processing of huge volumes of data in real time. We describe each one in turn.

**2.2.1. Naive Classifier.** This algorithm is based on a word count of positive and negative connotation words. It is the simplest and most intuitive of the classifiers. Recall that the lexicon is a list of hand-picked words that we found to be commonly used to describe stocks. Each word in the lexicon is identified as being positive, negative, or neutral. Each lexical entry is matched by a corresponding counterpart with a negation sign (see the online appendix for an example). Before processing any message, it is treated by a parsing algorithm that negates words if the context of the sentence requires it; for example, if the sentence reads “this stock is not good,” the word “good” is replaced by “good\_n,” signifying a negation. It

would then be counted as a sell word, rather than a buy word. This approach to negation is an innovation that is now used by computer scientists (see Pang et al. 2002).

Each word in a message is checked against the lexicon, and assigned a value (−1, 0, +1) based on the default value (sell, null, buy) in the lexicon. After this assignment, the net word count of all lexicon-matched words is taken, and if this value is greater than one, we sign the message as a buy. If this value is less than one, the message is taken to be a sell. All other messages are treated as neutral. The threshold value of one was chosen by experiment, and this may be adjusted for other applications. We did not subsequently attempt to improve the threshold value, so as to eliminate data-snooping bias. This classifier depends critically on the composition of the lexicon. By choosing words carefully in the lexicon, this approach may be adapted to other uses.

**2.2.2. Vector Distance Classifier.** If there are  $D$  words in the lexicon, and each word is assigned a dimension in vector space, then the lexicon represents a  $D$ -dimensional unit hypercube. Every message may be thought of as a word vector ( $m \in R^D$ ) in this space. The elements in the vector take values in the set  $\{0, 1, 2, \dots\}$  depending on how many times a word appears in the message. Suppose that the lexicon contains about 300 words. Then, each message may be characterized as a vector in 300-dimension space. As would be expected, each message contains a few lexical words only, and is therefore represented by a sparse vector.

A hand-tagged message (or grammar rule) in the training corpus (grammar) is converted into a vector  $G_j$ , and occupies a location in this  $D$ -dimensional Euclidian space. Each new message is classified by comparison to the cluster of pretrained vectors in this space. The angle  $\theta_j$  between the message vector ( $m$ ) and the vectors in the grammar ( $G_j$ ) provides a measure of closeness, i.e.,

$$\cos(\theta_j) = \frac{m \cdot G_j}{|m| \cdot |G_j|} \in [0, 1] \quad \forall j, \quad (1)$$

where  $|X|$  stands for the norm of vector  $X$ . Each message is assigned the classification of the grammar rule with which it has the smallest angle, i.e., that of  $\max_j[\cos(\theta_j)]$  (variations on this theme could use sets of top- $n$  closest vectors). Because  $\cos(\theta_j) \in [0, 1]$ , the vector distance classifier provides a measure of proximity in the form of percentage closeness—when the angle is small,  $\cos(\theta_j)$  is closer to one.

**2.2.3. Discriminant-Based Classifier.** The naive classifier (NC) weights lexical words equally. However, lexical words may have differential importance for classification. Some words, such as “buy,” may

<sup>7</sup> An electronic companion to this paper is available as part of the online version that can be found at <http://mansci.journal.informs.org/>.

be more indicative of sentiment than words such as “position.” Using the training corpus, we compute a measure of the discriminating ability of each word in our lexicon. We then replace the simple word count in the naive algorithm (NC) by a weighted word count.

The weights are based on a simple discriminant function for each word modified from the literature (see Chakrabarti et al. 1998, 2003—the latter paper demonstrates the usefulness of using the so-called Fisher discriminant statistic). Let the set  $i = \{\text{null, sell, buy}\}$  index the categories for our messages, and  $n_i$  be the number of messages in each category. Let the average number of times word  $w$  appears in a message of category  $i$  be denoted  $\mu_i$ . The number of times word  $w$  appears in a message  $j$  of category  $i$  is denoted  $m_{ij}$ . The discriminant formula for each word is

$$F(w) = \frac{(1/3) \sum_{i \neq k} (\mu_i - \mu_k)^2}{(\sum_i \sum_j (m_{ij} - \mu_i)^2) / (\sum_i n_i)} \quad \forall w. \quad (2)$$

This equation assigns a score  $F(w)$  to each word  $w$  in the lexicon, which is the ratio of the mean across-class squared variation to the average of within-class squared variation. The larger the ratio, the greater the discriminating power of word  $w$  in the lexicon. A good discriminant word maximizes across-class variation and minimizes within-class variation. Online Appendix B provides examples of the discriminant values of some of the words in the lexicon.

Each word in a message is checked against the lexicon, and assigned a signed value  $(-d, 0, +d)$ , based on the sign (sell = -1, null = 0, buy = +1) in the lexicon multiplied by the discriminant value  $d = F(w)$ . After this assignment, the net word count of all lexicon-matched words is taken, and if this value is greater than 0.01 (threshold), we sign the message as a buy. If this value is less than -0.01, the message is taken to be a sell. All other messages are treated as neutral. Again, the threshold was not improved subsequently so as to keep the empirical experiments in the sequel fair.

**2.2.4. Adjective-Adverb Phrase Classifier.** This classifier is based on the assumption that adjectives and adverbs emphasize sentiment and require greater weight in the classification process. This algorithm uses a word count, but restricts itself to words in specially chosen phrases containing adjectives and adverbs. Hence, the goal here is to focus only on the emphatic portions of the message.

We wrote program logic for a parts of speech “tagger” which, in conjunction with the dictionary, searches for noun phrases containing adjectives or adverbs (i.e., in its simplest form, this would be an adjective-noun pair). Whenever this is detected, we form a “triplet,” which consists of the adjective or

adverb and the two words immediately following or preceding it in the message. This triplet usually contains meaningful interpretive information because it contains the adjective or adverb, both of which are parts of speech that add emphasis to the phrase in which they are embedded. This simple heuristic identifies significant phrases, and the lexicon is used to determine whether these connote positive or negative sentiment. If the net count in these phrases was greater than or equal to 1 (-1), a positive (negative) tag is assigned to the message, or else it is neutral.

**2.2.5. Bayesian Classifier.** The Bayesian classifier relies on a multivariate application of Bayes’ theorem (see Mitchell 1997, Neal 1996, Koller and Sahami 1997, Chakrabarti et al. 1998). Recently, it has been used for Web search algorithms, for detecting web communities, and in classifying pages on Internet portals.<sup>8</sup> These methods are now also widely used for spam filters, and the following description summarizes the ideas of prior work as specifically applied here.

The classifier comprises three components: (i) lexical words, (ii) message text, and (iii) classes or categories (bullish, bearish, or neutral), resulting in the literature standard word-message-class  $(w, m, c)$  model. The Bayesian classifier uses word-based probabilities, and is thus indifferent to the structure of the language. Because it is language independent, it has wide applicability, which enables investigation of message boards in other financial markets, where the underlying language may not be English.

The notation is as follows. The total number of categories or classes is  $C (=3)$ ,  $c_i$ ,  $i = 1, \dots, C$ . Each message is denoted  $m_j$ ,  $j = 1, \dots, M$ , where  $M$  is the total number of messages. We define  $M_i$  as the total number of messages per class  $i$ , and  $\sum_{i=1}^C M_i = M$ . Words  $(w)$  are indexed by  $k$ , and the total number of lexical words is  $D$ . The set of lexical words is  $F = \{w_k\}_{k=1}^D$ .

Let  $n(m_j, w_k)$  be the total number of times word  $w_k$  appears in message  $m_j$ . We maintain a count of the number of times each lexical item appears in every message in the training data set. This leads naturally to the variable  $n(m_j)$ , the total number of lexical words in message  $m_j$  including duplicates. This is a simple sum,  $n(m_j) = \sum_{k=1}^D n(m_j, w_k)$ .

<sup>8</sup> Koller and Sahami (1997) develop a hierarchical model, designed to mimic Yahoo’s indexing scheme. Hence, their model has many categories and is more complex. On the other hand, their classifier was not discriminating emotive content, but factual content, which is arguably more amenable to the use of statistical techniques. Our task is complicated because the messages contain opinions, not facts, which are usually harder to interpret. The reader may obtain details of the hierarchical scheme by referring to the technical descriptions in Koller and Sahami (1997) and Chakrabarti et al. (1998) for the document model approach of a naive Bayes classifier. The exposition here briefly summarizes these approaches. We modify but try to retain as closely as possible the notation of the naive Bayes classifier of Chakrabarti et al. (1998).

An important quantity is the frequency with which a word appears in a message class. Hence,  $n(c_i, w_k)$  is the number of times word  $w$  appears in all  $m_j \in c_i$ . This is  $n(c_i, w_k) = \sum_{m_j \in c_i} n(m_j, w_k)$ . This measure has a corresponding probability:  $p(c_i, w_k)$  is the probability with which word  $w_k$  appears in all messages  $m$  in class  $c$ :

$$p(c_i, w_k) = \frac{\sum_{m_j \in c_i} n(m_j, w_k)}{\sum_{m_j \in c_i} \sum_k n(m_j, w_k)} = \frac{n(c_i, w_k)}{n(c_i)}. \quad (3)$$

We require that  $p(c_i, w_k) \neq 0 \forall c_i, w_k$ . Hence, an adjustment is often made to Equation (3) via Laplace's formula, which is

$$p(c_i, w_k) = \frac{n(c_i, w_k) + 1}{n(c_i) + D}.$$

If  $n(c_i, w_k) = 0$  and  $n(c_i) = 0 \forall k$ , then every word is equiprobable, i.e.,  $1/D$ . We now have the required variables to compute the conditional probability of a message  $j$  in category  $i$ , i.e.,  $\Pr[m_j | c_i]$ :

$$\Pr[m_j | c_i] = \frac{n(m_j)!}{\prod_{k=1}^D n(m_j, w_k)!} \times \prod_{k=1}^D p(c_i, w_k)^{n(m_j, w_k)}.$$

We also compute  $\Pr[c_i]$ , the proportion of messages in the training set classified into class  $c_i$ .

The classification goal is to compute the most probable class  $c_i$  given any message  $m_j$ . Therefore, using the previously computed values of  $\Pr[m_j | c_i]$  and  $\Pr[c_i]$ , we obtain the following conditional probability (applying Bayes' theorem):

$$\Pr[c_i | m_j] = \frac{\Pr[m_j | c_i] \cdot \Pr[c_i]}{\sum_{i=1}^C \Pr[m_j | c_i] \cdot \Pr[c_i]}. \quad (4)$$

For each message, Equation (4) delivers three posterior probabilities,  $\Pr[c_i | m_j] \forall i$ , one for each message category. The message is classified as being from the category with the highest probability.

### 2.3. Voting Amongst Classifiers

All the classifier methods used here are analytical and do not require optimization or search algorithms. Hence, there are no issues of numerical convergence. Given the huge data sets involved, this is an important consideration in the overall algorithm design. The numerical speed of the algorithms is complemented by enhancing statistical reliability using an original voting scheme, based on the intuition that all available information is not exploited when classifiers are used in isolation, instead of in conjunction. We will see that the primary benefit of the voting scheme lies in reducing false positives.

Final classification is based on achieving a simple majority vote amongst the five classifiers, i.e., three of five classifiers should agree on the message type. If a majority is not obtained, the message is not classified. This approach reduces the number of messages classified, but enhances classification accuracy.

### 2.4. Training and Evaluation

The classification algorithms are initially trained using a portion of the data, which we designate as the "training set," which we typically wanted to restrict to a size of less than 1,000 messages. The number of messages is deliberately kept small so as to assess whether the classifiers are amenable to a minimal amount of training. Hence, our approach is biased against performing well in-sample versus commercial Bayes classifiers. But, the small training set also prevents overfitting of the data (leading to poor out-of-sample performance), a common ailment in text classification algorithms.

### 2.5. Metrics

Our goal is to develop a sentiment index formed from a time-series accumulation of the sentiment from individual messages. The quality of this index will depend on the performance of our classification algorithms. We developed various assessment metrics, and also compared our models to a widely used algorithm in the public domain.

First, a standard approach to measuring the performance of classifiers is to examine the "confusion matrix" for statistical significance. The confusion matrix is a tableau that presents a cross-classification of actual message type versus classified message type. The confusion matrix has three rows and three columns. Each of the rows signifies the sentiment (hold, sell, or buy) posted by the author of a message to the stock board. The columns detail how many of these messages were classified in each of three categories: null, sell, or buy. The greater the weight of the diagonal of the confusion matrix, the lesser the confusion experienced by the algorithm. The null hypothesis for our test postulates no classification ability of the algorithm, i.e., the rows and columns of the confusion matrix are independent. We checked this using a standard  $\chi^2$  test:

$$\chi^2(4) = \frac{1}{9} \sum_{i=1}^9 \frac{(O_i - E_i)^2}{E_i} \quad (df = 4), \quad (5)$$

where  $O_i$  are the elements of the observed confusion matrix, and  $E_i$  are the elements of the matrix when no classification ability is present.

As a first step in assessment of algorithm performance, we collected and hand-classified a few hundred messages to comprise a training data set. We tuned the classifiers on the training set, and then undertook classification on testing sets. The quality of text on message boards is very poor, resulting in a hard problem even for human classification. Our ultimate goal lies in developing a sentiment index formed from the cumulated classification of messages over time, where buys, holds, and sells are  $\{+1, 0, -1\}$ , respectively.

Second, in addition to the  $\chi^2$  metric described above, we gain a useful understanding of the algorithm by examining the percentage of messages correctly classified. Note that the low clarity of messages implies that quite often, people would be likely to disagree on the classification. To gauge the extent of ambiguity, a reclassification of the training corpus was undertaken by a second human subject. This subject had nothing to do with the design of the study, and is not one of the authors. We believe that no bias existed, even for this informal test. Of the 374 training messages and 64 test messages, the two human subjects agreed on the classification of only 72.46% of the messages. We may like to think of the mismatch percentage of 27.54% ( $100.00 - 72.46$ ) as the “ambiguity coefficient” of the message boards. A more stable version of this coefficient would be one obtained from many (say  $n$ ) human subjects, for reasonably large  $n$  (approximately  $n \sim 10$ ), where the agreement percentage is based on the consensus of all  $n$  people. This might well result in an ambiguity coefficient a little higher than from just a few subjects. It is also intuitive that as we increase  $n$ , the ambiguity coefficient will first rise rapidly and then taper off to an asymptote, as there will be a core set of messages on which there can be little disagreement. Hence, there are two benchmarks of algorithm performance. One is perfect performance, i.e., a comparison with 100% accuracy rates, and the second is the human benchmark, i.e., an “agreement” coefficient, equivalent to 100 minus the ambiguity coefficient.

A third useful metric is the extent to which false positives occur. These are buy messages classified as sells and vice versa. We report the percentage of false positives to total messages. Note that the value of the sentiment index is doubly impacted if a false positive error is made because sentiment is incremented by the wrong sign. Hence, such errors tend to be more costly than other types of misclassification.

Fourth, in addition to false positives, we compare the value of the aggregate sentiment given no classification error versus that obtained based on our classifier. Note that if the classification error has no bias, then it is likely to have a smaller effect on the index than if there is bias.

To summarize, we examine four different metrics of classification performance: percentage classification accuracy, percentage of false positives, percentage error in aggregate sentiment, and a  $\chi^2$  test of no classification ability. We report the results for all our five individual algorithms, and for the majority voting algorithm. This voting scheme is applied in two ways—one, where messages with no majority are treated as “hold” messages, reported as type “Vote” in the subsequent tables, and two, where such messages are discarded, reported as “Vote- $d$ ” in the

tables. Finally, we also applied a popular, well-known software tool, the Rainbow algorithm of McCallum (1996). This is a highly optimized and widely used algorithm, and provides a good benchmark of performance.

First, in Table 1, Panel A, we run the classifier in-sample, i.e., setting the training data set to be the testing one. As is to be expected, the algorithms perform quite well. Note that, unlike our algorithms that use a fixed lexicon, the Rainbow algorithm first under-

**Table 1 Tests of Various Algorithms**

Algorithm	Accuracy	False positives	Sentiment error	$\chi^2$	Number of messages
Panel A: In-sample					
NvWtd	92.2460	0.2674	0.5348	64.2581	374
vecDist	49.1979	12.2995	23.5294	6.1179	374
DiscWtd	45.7219	16.0428	35.0267	4.4195	374
AdjAdv	63.3690	7.7540	20.3209	17.4351	374
BayesCI	60.6952	8.2888	14.4385	13.4670	374
Vote	58.2888	7.7540	17.3797	11.1799	374
Vote-d	62.8743	8.6826	17.0659	14.7571	334
Rainbow	97.0430	0.8065	3.2258	75.2281	372
Panel B: Test sample (out-of-sample)					
NvWtd	25.7393	18.2913	6.4622	2.0010	913
vecDist	35.7065	25.8488	5.8050	1.4679	913
DiscWtd	39.1019	25.1917	4.1621	2.6509	913
AdjAdv	31.3253	29.7919	51.3691	0.5711	913
BayesCI	31.5444	19.8248	12.2673	2.0873	913
Vote	30.0110	17.8532	8.3242	2.1955	913
Vote-d	33.1242	20.4517	7.7792	2.3544	797
Rainbow	33.1140	33.0044	38.7061	0.5458	912
Panel C: Test sample (out-of-sample)					
NvWtd	27.6024	20.8000	9.5031	33.8485	50,952
vecDist	38.4224	25.7281	8.5728	103.1038	50,952
DiscWtd	40.6049	25.0530	5.8172	131.8502	50,952
AdjAdv	32.6366	30.2186	54.4434	35.2341	50,952
BayesCI	33.2254	19.9835	14.1525	128.9712	50,952
Vote	31.8614	18.0268	10.2665	136.8215	50,952
Vote-d	35.8050	20.8978	11.0348	138.4210	43,952
Rainbow	35.2335	33.0893	40.0484	24.3460	49,575

*Notes.* In this table, we present statistics of the performance of various approaches to message classification. Five basic algorithms as described in the text are used: (i) Naive, (ii) Vector distance, (iii) Discriminant weighted, (iv) Adjective-adverb, (v) Bayes. The results of a majority vote are given, as are the same when no-majority votes are discarded (Vote- $d$ ,  $d =$  discard). The results of the Rainbow algorithm are also presented. We note that this algorithm is a Bayes classifier that develops a word set from the training set; the one in our algorithm uses an independent word set that is hand generated from the training set and other messages. Panel A has in-sample (test sample equals training sample) results from a training set of 374 messages taken from 1999. Panel B uses this training set on a small testing sample of 913 messages randomly taken from 25 stocks from the period July–August, 2001, which is a data set described in §3. Panel C uses a larger test sample of about 50,000 randomly chosen messages, taken from the same data set of §3. The first three measures are expressed in percentage terms. Sentiment error is expressed without sign. Tickers used in Panel C for July–August 2001: AMAT, EDS, INTU, MU, PSFT, BRCM, EMC, JNPR, NT, SCMR, CA, ERTS, LU, ORCL, SUNW, CSCO, IBM, MOT, PALM, TLAB, DELL, INTC, MSFT, PMTC, TXN.

takes a scan of words in the training set to determine which words are the best discriminants and then does the classification. Therefore, it performs very well in-sample in terms of classification accuracy. As can be seen, it behaves much like our naive algorithm, which delivers similar performance. The other algorithms do less well, as they are based on a pre-selected smaller set of words, all of which may not be ideal for this specific data set. However, the somewhat weaker in-sample performance is likely to be offset by a reduction in over fitting when working out-of-sample.

In Panels B and C of Table 1, we see that the out-of-sample performance is much lower than in-sample, as is to be expected. Looking at the false positives, the Rainbow model now performs worse than the others. It also has a high error in aggregate sentiment. We note that the voting algorithm produces the lowest false positive rate. The overall performance across all algorithms highlights the fact that classifying very dirty emotive text is indeed a hard problem. There are two routes to improving this baseline performance. First, we increase the size of the training set without making it too big to result in overfitting. Second, we screen messages for ambiguity before classification, discarding messages we find ambiguous. We take this up in the next section.

## 2.6. Ambiguity

Messages posted to stock boards are highly ambiguous. There is little adherence to correct grammar, and many words in the messages do not, and probably will never, appear in standard dictionaries. This makes the task of classification algorithms exceedingly hard, as opposed to say, spam filters, where the characteristics of spam versus nonspam e-mails are quite distinct. There is often only a subtle difference between a buy and a sell message on a stock board, resulting in many false positives. Ambiguity is related to the absence of “aboutness,” in that classification based on word counts may not always correctly capture what a message is about (see Morville 2005).

Stock postings are highly ambiguous and we would like to examine whether algorithm performance is a function of ambiguity or not. To reduce ambiguity, we developed a method to exclude possibly ambiguous messages. We did this by using the General Inquirer, a computer-assisted approach for content analyses of textual data from Harvard University.<sup>9</sup> The algorithms in this approach return a count of optimistic and pessimistic words in text. We filtered messages

for ambiguity using a simple metric, i.e., the difference between the number of optimistic and pessimistic words as a percentage of the total words in a body of text. We denote this as the optimism score. Because we use a completely different dictionary and approach for the optimism score, we ensure that none of the results are biased in favor of our algorithm relative to others. As a first pass, we examined over 50,000 messages for which we had buy, sell, and hold classifications provided by posters on Yahoo!, and looked at the average optimism score in each category. These are as follows:

Message type	Optimism score	
	Mean	Std. dev.
Buy	0.032	0.075
Hold	0.026	0.069
Sell	0.016	0.071

Therefore, on average, the optimism score does rank order the categories well. However, note that the standard deviation of these scores is quite large. Hence, for example, if we would like to filter out ambiguous buy messages, we may do so by filtering in messages that were posted as buys and had optimism scores greater than one standard deviation away from the mean ( $>0.032 + 1 \times 0.07$ ). Likewise, to filter in sell messages with low-ambiguity scores, we would want sell messages with ambiguity scores lower than a standard deviation from the mean ( $<0.016 - 1 \times 0.07$ ). For hold messages, we reduce ambiguity by taking ever smaller intervals around the mean of 0.026. To decrease ambiguity levels, we filter in buy (sell) messages only when they are increasing sigma levels away from the mean score. A high-ambiguity (low noise reduction) message is filtered in when it is at least one-sigma from the mean score in the right direction. Medium-ambiguity messages are two-sigma and low-ambiguity messages are three-sigma from the mean score.

In Tables 2 and 3, we reduce ambiguity as we proceed from Panel A down to Panel C and examine the changes in classification metrics. In Panel A, the cutoff for filtering in a message is based on one standard deviation as above. In Panel B, we raise the cutoff to two standard deviations, and in Panel C, the cutoff is three standard deviations. Naturally, we sample fewer messages as we proceed to the least ambiguous case. For hold messages, Panel A uses a spread of 0.01 around the mean, Panel B uses 0.005, and Panel C uses 0.0025.

Table 2 uses a small training set (size 374), and the one used in Table 3 is larger (size 913). The testing sets are the same in both tables. Accuracy levels increase as the training set increases in size. However,

<sup>9</sup>See <http://www.wjh.harvard.edu/~inquirer/> for more details. This system is also used in recent work by Tetlock (2005) and Tetlock et al. (2006).

**Table 2** Classification as a Function of Ambiguity: Training Set Size of 374

Algorithm	Accuracy	False positives	Sentiment error	$\chi^2$	Number of messages
Panel A: High ambiguity					
NvWtd	25.4777	20.3556	9.5011	35.7433	7,536
vecDist	49.7479	15.9501	3.1714	138.5787	7,536
DiscWtd	54.8832	13.1900	10.0318	204.3926	7,536
AdjAdv	29.8301	32.8822	53.3439	4.7059	7,536
BayesCl	43.0467	10.6290	9.0366	159.1631	7,536
Vote3	41.7596	8.8110	5.6661	169.2776	7,536
Vote3-d	47.7541	10.1794	6.8067	168.9292	6,523
Rainbow	34.8511	34.7309	39.5246	3.0299	7,489
Panel B: Medium ambiguity					
NvWtd	23.9664	17.5711	14.0181	15.2671	1,548
vecDist	53.2946	9.5607	2.7778	47.9415	1,548
DiscWtd	58.1395	8.5917	9.4961	58.2234	1,548
AdjAdv	26.8734	32.8165	52.9716	2.0661	1,548
BayesCl	45.2196	6.3307	9.0439	45.9316	1,548
Vote3	44.1860	4.7804	5.9432	49.0391	1,548
Vote3-d	49.5549	5.4896	4.8961	49.1041	1,348
Rainbow	34.1952	32.2560	43.5036	1.3772	1,547
Panel C: Low ambiguity					
NvWtd	20.0000	14.8276	10.0000	5.5572	290
vecDist	55.8621	3.7931	0.3448	12.6265	290
DiscWtd	57.2414	5.5172	8.2759	11.7723	290
AdjAdv	24.8276	34.1379	55.5172	0.5376	290
BayesCl	48.9655	2.0690	6.8966	11.6353	290
Vote3	49.3103	1.3793	5.1724	12.2487	290
Vote3-d	52.7778	1.5873	1.9841	12.1598	252
Rainbow	37.5862	24.1379	33.4483	1.0625	290

*Notes.* In this table, we present statistics of the performance of various approaches to message classification. Five basic algorithms as described in the text are used: (i) Naive, (ii) Vector distance, (iii) Discriminant weighted, (iv) Adjective-adverb, (v) Bayes. The results of a majority vote are given, as are the same when no-majority votes are discarded (Vote-*d*, *d* = discard). The results of the Rainbow algorithm are also presented. We note that this algorithm is a Bayes classifier that develops a word set from the training set; the one in our algorithm uses an independent word set that is hand generated from the training set and other messages. The extent of ambiguity in the test data set declines as we proceed from Panel A to Panel C (as described in §2.6). A high-ambiguity message is filtered in when it is at least one-sigma from the mean score in the right direction. Medium-ambiguity messages are two-sigma and low-ambiguity messages are three-sigma from the mean score. Hence, test sample size falls with decreasing ambiguity. Accuracy, false positives, and sentiment error are expressed in percentage terms. Sentiment error is expressed without sign.

this does not impact the percentage of false positives. The false positives decline dramatically with a reduction in ambiguity as expected. The algorithms that persistently return a higher number of false positives than the others may be overfitting the data. Overall, the algorithms increase in performance as we tune down the ambiguity of the messages. We achieve an accuracy range of between 60%–70%, which is good for classification of sentiment (see also Pang et al. 2002 for an application to sentiment parsing for movies with higher-accuracy levels).

**Table 3** Classification as a Function of Ambiguity: Training Set Size of 913

Algorithm	Accuracy	False positives	Sentiment error	$\chi^2$	Number of messages
Panel A: High ambiguity					
NvWtd	45.5016	34.2224	1.8843	16.7918	7,536
vecDist	61.0934	14.2118	10.2309	236.2058	7,536
DiscWtd	54.8832	13.1900	10.0318	204.3926	7,536
AdjAdv	64.1056	33.5987	41.2818	74.0030	7,536
BayesCl	57.4443	12.3275	7.8025	238.8515	7,536
Vote3	53.3838	10.1778	14.3577	242.0414	7,536
Vote3-d	63.2705	12.1534	12.4703	227.2089	6,311
Rainbow	64.8818	32.6479	13.0191	86.8046	7,489
Panel B: Medium ambiguity					
NvWtd	46.9638	30.7494	1.0982	6.2881	1,548
vecDist	64.5995	8.6563	8.6563	69.9800	1,548
DiscWtd	58.1395	8.5917	9.4961	58.2234	1,548
AdjAdv	65.7623	28.4884	41.5375	23.2180	1,548
BayesCl	61.4341	7.7519	8.0749	67.8975	1,548
Vote3	58.3979	6.3307	13.5659	68.8180	1,548
Vote3-d	66.7671	7.3851	11.6805	65.8436	1,327
Rainbow	65.4816	28.9593	10.7304	27.4832	1,547
Panel C: Low ambiguity					
NvWtd	46.5517	25.5172	9.3103	1.9822	290
vecDist	66.8966	3.7931	7.9310	16.9034	290
DiscWtd	57.2414	5.5172	8.2759	11.7723	290
AdjAdv	61.3793	24.1379	40.0000	4.7444	290
BayesCl	64.4828	4.4828	8.2759	15.2331	290
Vote3	63.4483	4.1379	12.4138	15.3550	290
Vote3-d	66.7939	4.5802	11.0687	14.5289	262
Rainbow	67.5862	18.2759	12.0690	9.0446	290

*Notes.* Same as in Table 2.

In the rest of the paper, we explore the sentiment index generated by the voting classifier. We note that this model has fairly low error in sentiment index generation, especially out-of-sample. It also has lower incidence of false positives than most of the remaining classifiers. The rest of our analyses explore whether the sentiment index offers return predictability or not, and its relation to various stock market variables.

### 3. Experiments: Sentiment, and Stock Markets

In the previous section, we assessed the ability of the index to match human classification. In this section, we analyze the relationship of sentiment to stock market data so as to understand the connection of message board discussion to market economics.

#### 3.1. Data

Our data comprises 24 tech-sector stocks, present in the Morgan Stanley High-Tech Index (MSH). These stocks were chosen so as to focus on the tech sector, and also because their message boards showed a wide range of activity. For a period of two months, July and August 2001, we downloaded every message posted to these boards. This resulted in a total of

145,110 messages. These messages were farmed by a lengthy process of Web-scraping.

We then employed our voting algorithm with the larger training set (found to provide the best performance in the previous section) to assess each message and determine its sentiment. The messages up to 4 P.M. for each day were used to create the aggregate sentiment for the day up to close of trading for each of the stocks. Buy messages increment the index by one, and sell messages decrement the index by one. We then aggregated sentiment daily (equally weighted) across all sample stocks so as to obtain an aggregate sentiment index for our tech portfolio.

We downloaded the MSH index for the same period. The following measures are constructed for further analysis:

1. *Normalized indexes.* We normalized both the MSH index and the aggregate sentiment index by subtracting the mean value and dividing by the standard deviation of the data series. We also did this for the sentiment of each individual stock series. This makes the scale of the values the same across all stocks, so that we may combine them in our analysis. All the analyses presented are based on the normalized series.

2. *Disagreement.* Following the work of Das et al. (2005), we constructed their disagreement measure which is as follows:

$$\text{DISAG} = \left| 1 - \frac{B - S}{B + S} \right|,$$

where  $B$  is the number of buy messages and  $S$  is the number of sell messages. This measure lies between zero (no disagreement) and one (high disagreement). It may be computed for any time period; in our analyses, we computed it daily.

3. *Volatility.* We measure intraday volatility as the difference between the high and low stock prices for the day divided by the average of the open and closing price.

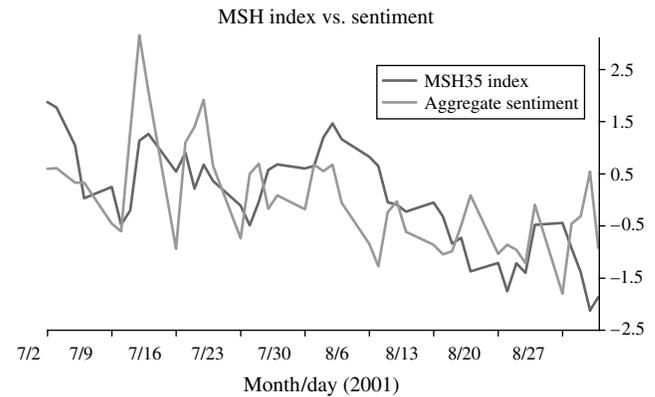
4. *Volume.* This is the trading volume in number of shares traded in the day. We also keep track of message volume for each stock.

Weekends are eliminated as message posting is not matched with corresponding trading. We also normalized the data on disagreement, volume, and volatility, so as to be able to combine the data across stocks. This is expedient given the short time series. We first use the data to look at the relationship of aggregate sentiment to the stock index. The analysis is undertaken using normalized time series.

### 3.2. Sentiment and the Stock Index

We aggregated sentiment across all the stocks, and examined how closely this sentiment related to the MSH index. This is represented in Figure 2, which

**Figure 2** Normalized MSH Index and Aggregate Sentiment, Daily, July–August 2001



indicates that these two series do track each other closely, implying that the extracted sentiment is not excessively noisy. The correlation between the stock index and the sentiment time series is 0.48. The sentiment index is highly autocorrelated out to two trading weeks (evidenced by Box-Ljung tests, not reported). The autocorrelation ranges from 0.8 at one lag to 0.15 at a lag of 10 trading days. Therefore, we need to be careful in interpreting the results of small sample regressions relating sentiment to other market variables. Individual normalized stock sentiment autocorrelations are lower (one lag correlation of normalized sentiment across all stocks is 0.27). Therefore, the sentiment index is more autocorrelated than individual stock sentiment. This is consistent with positive cross-autocorrelations of sentiment, and similar results are known for cross-autocorrelations of stock returns, as in Lo and MacKinlay (1990).

We examine the statistical relationship of the stock series (MSH) to the sentiment series (SENTY). We regress the aggregate sentiment level on lagged values of the stock index and sentiment. We also regress the stock index level on the same lagged values. The same pair of regressions is also undertaken in changes (first differences) in addition to the ones in levels. These results are presented in Table 4. The regressions in levels show that tech index is strongly related to its value on the previous day, and weakly related to the sentiment index value from the previous day at the 10% significance level. From the second regression, we see that the sentiment index on a given day is significantly related to its prior day's value, but not to that of the stock index. The interpretation of this pair of equations is that sentiment does offer some explanatory power for the level of the index. However, as noted earlier, the high autocorrelation of these series is a matter of concern, and therefore, we also present the regressions in changes. Whereas the stock index does appear to be related to the lagged sentiment value, the relationship is now weak, and the

**Table 4** Regressions Relating the Stock Index (MSH) to the Sentiment Index

Dependent variable	Intercept	MSH <sub>t</sub>	SENTY <sub>t</sub>	R <sup>2</sup>
Regressions in levels				
MSH <sub>t+1</sub>	-0.081	0.793***	0.154*	0.77
<i>t</i> -stat	-1.12	9.28	1.86	
SENTY <sub>t+1</sub>	-0.028	0.100	0.451***	0.25
<i>t</i> -stat	-0.21	0.62	2.90	
Dependent variable	Intercept	ΔMSH <sub>t</sub>	ΔSENTY <sub>t</sub>	R <sup>2</sup>
Regressions in changes				
ΔMSH <sub>t+1</sub>	-0.094	-0.082	0.138*	0.07
<i>t</i> -stat	-1.19	-0.51	1.71	
ΔSENTY <sub>t+1</sub>	-0.083	-0.485	-0.117	0.08
<i>t</i> -stat	-0.53	-1.51	-0.72	

*Notes.* We aggregated the daily sentiment equally weighted across stocks in the sample to get a sentiment index value (SENTY) for each day, ignoring weekends, when there is no trading. Regressions are presented both in levels and in first differences on normalized values. The number of asterisks determines the level of significance: \* = 10% level, \*\* = 5% level, \*\*\* = 1% level.

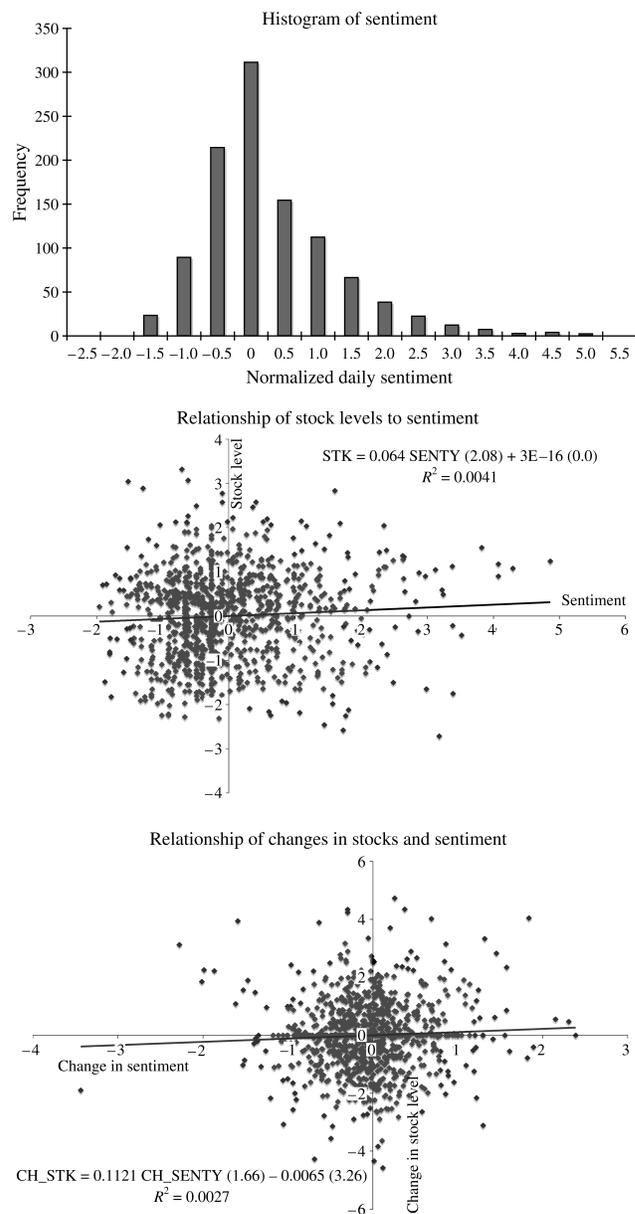
regression has low explanatory power as may be seen in the diminished *R*-square values. A more detailed empirical paper is required to explore a longer time series, and will also enable aggregation of sentiment over longer periods to eliminate the autocorrelation.

**3.3. Individual Stocks**

Given that there appears to be a link from sentiment to the index at the aggregate level, we now drill down to the individual stock level to see if there is a relationship. Our analysis uses the normalized stock price and normalized sentiment index for each stock. The normalization allows us to stack up all stocks together and then conduct the analysis using pooled data.

Figure 3 presents a look at the data on individual stock levels and sentiment for the pooled sample of all stocks. The top graph shows the histogram of normalized sentiment for all stock days. There is a noticeable positive skew. This suggests that there are days when the message boards are very bullish on a stock and tend to display aggressive optimism (see Antweiler and Frank 2004 for a bullishness index); clearly, the same does not occur with bearish views, which seem to be espoused less strongly, although Tetlock (2005) finds that negative media coverage presages downward moves in stock prices. The middle graph is a scatter plot of all stock-day pairs of normalized stock level and sentiment. The relationship of stock level to sentiment is significant, with a *t*-statistic over two (shown on the graph), but the overall fit of the model is understandably poor as the regression lacks several other variables that explain stock levels. However, as mentioned before, it is important to also examine this relationship in first differences as well, and the bottom graph shows that this significance is markedly

**Figure 3** Individual Stocks and Sentiment



*Notes.* The top graph shows the distribution of normalized sentiment for all stock days. There is a noticeable positive skew. The middle graph is a scatter plot of all stock-day pairs of normalized stock level and sentiment. The relationship of stock level to sentiment is significant, but the bottom graph shows that this significance is lost when the same data is presented in first differences. The list of tickers is: AMAT, BRCM, CA, CSCO, DELL, EDS, EMC, ERTS, IBM, INTC, INTU, JNPR, LU, MOT, MSFT, MU, NT, ORCL, PALM, PMTC, SCMR, SUNW, TLAB, TXN.

reduced (*t*-statistic of 1.66) when this is accounted for, and hence, it is hard to infer a strong predictive ability for sentiment in forecasting daily movements for individual stocks.

We can see from Figure 3 that there is no strong relationship from sentiment to stock prices on average across the individual stocks. Neither regression (levels or changes) shows statistical strength. On the

other hand, in §3.2 for the aggregate index, we found a statistical relation from sentiment to stock prices. An implication of these results is that aggregation of individual stock sentiment may be resulting in a reduction of idiosyncratic error in sentiment measurement, giving significant results at the index level. Further, this evidence is also consistent with the presence of cross-autocorrelations of sentiment amongst stocks.

**3.4. Sentiment, Disagreement, Volumes, and Volatility**

We now examine the relationship of our sentiment measure to the disagreement measure, message volume, trading volume, and volatility, defined previously in §3.1. Previous work (Das et al. 2005, Antweiler and Frank 2004) suggests strong relationships between these variables, and Figure 4 confirms this.

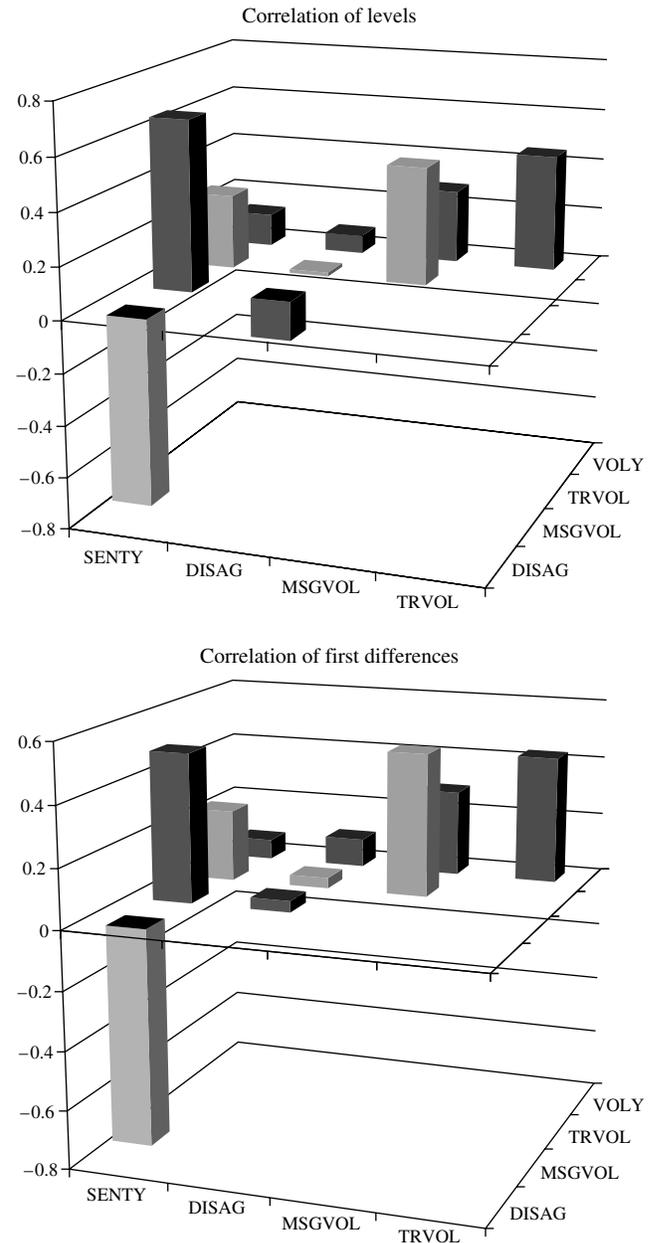
The following relationships are evident. First, sentiment is inversely related to disagreement. Hence, when disagreement increases, sentiment drops. Alternatively, there is greater disagreement when sentiment is falling rather than when it is rising. Second, sentiment is correlated to high posting volume, suggesting that increased discussion indicates optimism. It hints at the fact that people prefer making posts when they are bullish on a stock (Antweiler and Frank 2004 therefore name their extracted sentiment a “bullishness” index). Third, there is a strong relationship between message volume and volatility, consistent with the findings of Antweiler and Frank (2002), who find that sentiment extracted from message boards is not predictive of stock movements, but activity on these boards may presage increases in volatility. Finally, trading volume and volatility are strongly related to each other.

Antweiler and Frank (2004) find that message volume explains volatility but not returns, and in another paper, Antweiler and Frank (2002) show that high message posting appears to be more likely in stocks with poorer returns. Our data allows us to revisit both these findings, and we find confirmation of these results using approaches different from theirs, albeit with a small sample. We regress changes in volatility on changes in sentiment, disagreement, message volume, and trading volume. Results are shown in the first panel in Table 5. Message volume significantly explains changes in stock levels as well as volatility. Changes in sentiment are positively related to changes in stock levels as well. Interestingly, disagreement does not explain volatility as one might expect.

**4. Conclusion**

We developed a methodology for extracting small investor sentiment from stock message boards. Five distinct classifier algorithms coupled by a voting scheme are evaluated using a range of metrics.

**Figure 4** Correlations between Sentiment, Disagreement, Message Volume, Trading Volume, and Volatility



*Notes.* The correlations are computed for each pair of variables using all stock days across all tickers in our sample. The top plot shows the correlations of levels and the one on the bottom is in first differences.

Time series and cross-sectional aggregation of message sentiment improves the quality of the sentiment index. Sentiment aggregated across stocks tracks index returns; this effect is weak for individual stocks. Whereas our methodology results in classification accuracy levels similar to that of widely used Bayes classifiers, the noise-reduction approaches we employ substantially reduces the number of false positives generated in the classification, as well as improves the accuracy of the index value itself. Our suite of

**Table 5** Relationship of Changes in Volatility and Stock Levels to Changes in Sentiment, Disagreement, Message Volume, and Trading Volume

Independent variables	$\Delta$ VOLY	$\Delta$ STK
Intercept	−0.000 −0.01	−0.056** −3.29
$\Delta$ SENTY	−0.106 −1.50	0.059* 1.83
$\Delta$ DISAG	0.008 0.14	0.001 0.03
$\Delta$ MSGVOL	0.197*** 3.34	−0.080*** −2.99
$\Delta$ TRVOL	0.447*** 11.84	0.000 0.01
$R^2$	0.20	0.02

*Notes.* The results relate to pooled regressions for all stocks in the sample. *t*-statistics are presented below the estimates. The normalized variables are: VOLY is volatility, STK is stock price, SENTY is the sentiment for a stock for the day, DISAG is disagreement, MSGVOL is the number of messages per day, TRVOL is the number of shares traded per day. The number of asterisks determines the level of significance: \* = 10% level, \*\* = 5% level, \*\*\* = 1% level.

techniques presents an effective approach to classify noisy stock message board postings to develop an index of sentiment. Given the quality of text in these messages, and the nonuniform emotive content therein, we believe that more research using these ideas is worth pursuing.

As an application example, we created a tech-sector sentiment index from a representative set of 24 stocks. This index is related to stock levels, but only weakly relates when assessed in first differences. However, the relationship is weaker at the individual stock level; therefore, aggregation of sentiment reduces some of the noise from individual stock board postings. Whereas our sentiment index has expected contemporaneous relationships with various market variables, the disagreement measure we create evidences very little correlation to other variables. The overall evidence suggests that market activity is related to small investor sentiment and message board activity. Thus, the algorithms developed in this paper may be used to assess the impact on investor opinion of management announcements, press releases, third-party news, and regulatory changes.

The techniques here may be tailor-made for other applications by changing the lexicon and grammar appropriately, and for other languages too, because the methods are general and are not domain specific. Some other areas in which application is possible are as follows. First, there is a limited understanding of the microstructure of tech stocks. Because these stocks have the most active message boards, the sentiment classifier may support empirical work in this

domain. Second, the algorithms may be used to investigate the mechanics of herding. A third application is that of monitoring market activity. Regulators are concerned about market manipulation that goes undetected amongst the millions of messages posted to message boards every day. Fourth, firms may use the classifier to monitor their message boards for investor reaction to management actions. Finally, the sentiment index may be applied to testing theories in the domain of behavioral finance.

## 5. Electronic Companion

An electronic companion to this paper is available as part of the online version that can be found at <http://mansci.journal.informs.org/>.

## Acknowledgments

The authors owe a special debt to the creative environment at University of California Berkeley's Computer Science Division, where this work was begun. The comments of the department editor, David Hsieh, an associate editor, and two referees contributed immensely to many improvements. Thanks to David Levine for many comments and for the title. The authors are also grateful to Vikas Agarwal, Chris Brooks, Yuk-Shee Chan, David Gibson, Geoffrey Friesen, David Leinweber, Asis Martinez-Jerez, Patrick Questembert, Priya Raghuram, Sridhar Rajagopalan, Ajit Ranade, Mark Rubinstein, Peter Tufano, Raman Uppal, Shiv Vaithyanathan, Robert Wilensky, and seminar participants at Northwestern University, University of California Berkeley-EECS; London Business School; University of Wisconsin, Madison; the Multinational Finance Conference, Italy; the Asia Pacific Finance Association Meetings, Bangkok; European Finance Association Meetings, Barcelona; Stanford University; and Barclays Global Investors for helpful discussions and insights. Danny Tom and Jason Waddle were instrumental in delivering insights into this paper through other joint work on alternative techniques via support vector machines. The first author gratefully acknowledges support from a Breetwor Fellowship, the Dean Witter Foundation, and a research grant from Santa Clara University.

## References

- Admati, A., P. Pfleiderer. 2000. Noisytalk.com: Broadcasting opinions in a noisy environment. Working Paper 1670R, Stanford University, Stanford, CA.
- Antweiler, W., M. Frank. 2002. Internet stock message boards and stock returns. Working paper, University of British Columbia, Vancouver, BC, Canada.
- Antweiler, W., M. Frank. 2004. Is all that talk just noise? The information content of Internet stock message boards. *J. Finance* 59(3) 1259–1295.
- Antweiler, W., M. Frank. 2005. The market impact of news stories. Working paper, University of British Columbia, Vancouver, BC, Canada.
- Bagnoli, M., M. D. Beneish, S. G. Watts. 1999. Whisper forecasts of quarterly earnings per share. *J. Accounting Econom.* 28(1) 27–50.

- Chakrabarti, S., B. Dom, P. Indyk. 1998. Enhanced hypertext categorization using hyperlinks. *Proc. 1998 ACM SIGMOD Internat. Conf. Management Data*, ACM Press, New York, 307–318.
- Chakrabarti, S., S. Roy, M. V. Soundalgekar. 2003. Fast and accurate text classification via multiple linear discriminant projections. *VLDB J.* **12**(2) 170–185.
- Chakrabarti, S., B. Dom, R. Agrawal, P. Raghavan. 1998. Scalable feature selection, classification and signature generation for organizing large text databases into hierarchical topic taxonomies. *VLDB J.* **7**(3) 163–178.
- Charniak, E. 1993. *Statistical Language Learning*. MIT Press, Cambridge, MA.
- Choi, J., D. Laibson, A. Metrick. 2002. Does the Internet increase trading? Evidence from investor behavior in 401(k) plans. *J. Financial Econom.* **64** 397–421.
- Das, S., A. Martinez-Jerez, P. Tufano. 2005. e-Information. *Financial Management* **34**(5) 103–137.
- Godes, D., D. Mayzlin. 2004. Using online conversations to study word of mouth communication. *Marketing Sci.* **23**(4) 545–560.
- Joachims, T. 1999. Making large-scale SVM learning practical. B. Scholkopf, C. Burges, A. Smola, eds. *Advances in Kernel Methods—Support Vector Learning*. MIT Press, Cambridge, MA, 243–253.
- Koller, D., M. Sahami. 1997. Hierarchically classifying documents using very few words. *Internat. Conf. on Machine Learning*, Vol. 14. Morgan-Kaufmann, San Mateo, CA.
- Lam, S. L., J. Myers. 2001. Dimensions of website personas. Working paper, University of California, Berkeley, CA.
- Lavrenko, V., M. Schmill, D. Lawrie, P. Ogilvie, D. Jensen, J. Allen. 2000. Mining of concurrent text and time series. *Proc. Knowledge Discovery Data Mining, 2000 Conf. Text Mining Workshop*, 37–44.
- Leinweber, D., A. Madhavan. 2001. Three hundred years of stock market manipulation. *J. Investing* **10** 7–16.
- Lo, A. W., A. C. MacKinlay. 1990. When are contrarian profits due to stock market overreaction? *Rev. Financial Stud.* **3**(2) 175–205.
- McCallum, A. 1996. *Bow: A Toolkit for Statistical Language Modeling, Text Retrieval, Classification and Clustering*. School of Computer Science, Carnegie-Mellon University, Pittsburgh, PA, <http://www.cs.umass.edu/~mccallum/code-data.html>.
- Minsky, M. 1985. *Society of Mind*. Simon & Schuster, New York.
- Mitchell, T. 1997. *Machine Learning*. McGraw-Hill, New York.
- Morville, P. 2005. *Ambient Findability*. O'Reilly, Sebastopol, CA.
- Neal, R. 1996. Bayesian learning for neural-networks. *Lecture Notes in Statistics*, Vol. 118. Springer-Verlag, New York.
- Pang, B., L. Lee, S. Vaithyanathan. 2002. Thumbs up? Sentiment classification using machine learning techniques. *Proc. Conf. Empirical Methods Nat. Language Processing*, 79–86.
- Stone, A. 2001. The Darwinism of day trading. *Business Week* (May 23), online edition.
- Tetlock, P. 2005. Giving content to investor sentiment: The role of media in the stock market. *J. Finance*. Forthcoming.
- Tetlock, P., M. Saar-Tsechansky, S. Macskassy. 2006. More than words: Quantifying language to measure firms' fundamentals. Working paper, University of Texas, Austin, TX.
- Tumarkin, R., R. Whitelaw. 2001. News or noise? Internet postings and stock prices. *Financial Analysts J.* **57**(3) 41–51.
- Vapnik, V. 1995. *The Nature of Statistical Learning Theory*. Springer-Verlag, New York.
- Vapnik, V., A. Chervonenkis. 1964. On the uniform convergence of relative frequencies of events to their probabilities. *Theory Probab. and Its Appl.* **16**(2) 264–280.
- Wakefield, J. 2001. Catching a buzz. *Scientific Amer.* **285**(5) 30–32.
- Wysocki, P. 1998. Cheap talk on the web: The determinants of postings on stock message boards. Working Paper 98025, University of Michigan Business School, Ann Arbor, MI.

**e - companion**

ONLY AVAILABLE IN ELECTRONIC FORM

## Electronic Companion—“Yahoo! for Amazon: Sentiment Extraction from Small Talk on the Web” by Sanjiv R. Das and Mike Y. Chen, *Management Science*, 10.1287/mnsc.1070.0704.

---

### Online Appendices

#### A. Overview of the Methodology Flow

This is a brief overview of the model components, which complements the model schematic presented earlier in Figure 1. The programs were coded in Java.

##### A.1. The Dictionary

Our data includes auxiliary information on the English language. To exploit parts-of-speech usage in messages, a dictionary was used to detect adjectives and adverbs for the classifier algorithms. This dictionary is called CUVOALD (Computer Usable Version of the Oxford Advanced Learner’s Dictionary).<sup>EC1</sup> It contains parts-of-speech tagging information, and we wrote appropriate program logic to use this dictionary while analyzing messages for grammatical information.

##### A.2. The Lexicon

Words are the heart of any language inference system, and in a specialized domain, this is even more so. The sentiment classification model relies on a lexicon of “discriminant” words, which comprise the lexicon. The lexicon is designed using domain knowledge and statistical methods. A discriminant function is used to statistically detect which words in the training corpus are good candidates for classifier usage (the details of the discriminant function are provided in §2.2.3). Therefore, the lexicon is essentially a collection of words relevant to the classification problem, which will be used by the classifier algorithms to discriminate buy messages from sell messages. Hence, we exercised care in creating the lexicon, so as to include many useful words that would enable the algorithms to discriminate positive from negative sentiment. Clearly, a different lexicon will result in different classifications; this injects flexibility and the ability to tune the algorithm, but also requires domain expertise. We had to read thousands of messages to cull the set of words that now comprise the lexicon. The user’s goal is to populate the lexicon with words of high discriminant value, and this is where the application of domain expertise is valuable. Over time, more words may be added to the lexicon, which improves in this evolutionary manner. More details on the lexicon are presented in Appendix B.

##### A.3. The Grammar

A grammar may be defined as a set of functions or rules applied in conjunction with the lexicon to extract sentiment from text. Correspondences between word sets, language features, and classification types comprise the grammar. In our setting, the training corpus is the grammar. This set of messages, once hand-tagged, may be thought of as a set of rules that govern the classification of other messages. One way to approach classification of any message is to search the grammar for a rule that may be applied to the message. For example, a distance function under a carefully chosen metric may be used to identify the applicable rule. Suppose we wish to analyze message  $M$ . We compare, using some metric, the relationship of this message  $M$  to a set of other messages  $G$ , and find the one that is its closest look-alike. We then equate the properties of message  $M$  to those of the proxy. The set

<sup>EC1</sup> The dictionary was downloaded from Birkbeck College, University of London. It is the creation of Roger Mitton of the Computer Science Department. It contains about 70,000 words, and covers most of the commonly used words in the English language. Informal tests of the dictionary showed that about 80–90 percent of the words in a message were found in the dictionary.

of preclassified messages  $G$  is denoted the grammar, and the rule that finds the proxy message or a proxy set of messages is codified in a classification algorithm. The classification algorithm implements a rule that finds closest messages in a grammar, using the words in the lexicon as variables. Some of the algorithms use only the grammar, or the lexicon, and some use both.<sup>EC2</sup>

#### A.4. Message Pre-Processing

Before applying the lexicon-grammar based algorithms, each message is preprocessed to enable cleaner interpretation. First, we carry out “HTML Cleanup,” which removes all HTML tags from the body of the message as these often occur concatenated to lexical items of interest. Examples of some of these tags are: <BR>, <p>, &quot;, etc. Second, we expand abbreviations to their full form, making the representation of phrases with abbreviated words common across the message. For example, the word “ain’t” is replaced with “are not,” “it’s” is replaced with “it is,” etc. Finally, we handle negation words. Whenever a negation word appears in a sentence, it usually causes the meaning of the sentence to be the opposite of that without the negation. For example, the sentence “It is not a bullish market” actually means the opposite of a bull market. Words such as “not,” “never,” “no,” etc., serve to reverse meaning. We handle negation by detecting these words and then tagging the rest of the words in the sentence after the negation word with markers, so as to reverse inference. These three parsers deliver a clean set of messages for classification.

## B. Construction of the Lexicon

The features of the lexicon are as follows:

1. These words are hand-selected based on a reading of several thousand messages.
2. The lexicon may be completely user-specified, allowing the methodology to be tailored to individual preference. For example, if the user is only interested in messages that relate to IPOs, a lexicon containing mostly IPO-related words may be designed. (The grammar, i.e. the training set would also be correspondingly tagged.)
3. For each word in the lexicon, we tag it with a “base” value, i.e. the category in which it usually appears. For example, the word “sell” would be naturally likely to appear in messages of type SELL, and we tag “sell” with base value 1. If the word is of BUY type, we tag it with value 3, and NULL words are tagged 0.<sup>EC3</sup> Every time a new word is added to the lexicon, the user is required to make a judgment on the base type.
4. Each word is also “expanded,” i.e. appears in the lexicon in all its forms, so that across forms, the word is treated as one word. This process is analogous to stemming words, except that we exhaustively enumerate all forms of the word rather than stem them.<sup>EC4</sup>
5. Each word is also entered with its “negation” counterpart, i.e. the sense in which the word would appear if it were negated. Negation is detected during preprocessing (described later) and is used to flag portions of sentences that would be reversed in meaning.

An example of a lexical entry along with its base value, expansion and negation is provided below:

```
3 favorable favorite favorites favoring favored
1 favorable_n favorite_n favorites_n favoring_n favored_n
```

All forms of the word appear in the same line of the lexicon. As can be seen, a tag is attached to each negated word in the second line above. The default classification value (the “base” value) is specified at the beginning of the line for each lexical item (i.e. a 0, 1, or 3).

<sup>EC2</sup> We may think of the grammar as Schank (1975) did, i.e., it is a “conceptual processor.” With stock market messages, the language is cryptic, and the grammar rules must work together so as to make sense of the “thought bullets” posted to the web. Schank states this particularly well: “People do not usually state all the parts of a given thought that they are trying to communicate because the speaker tries to be brief and leaves out assumed or inessential information. The conceptual processor searches for a given type of information in a sentence or a larger unit of discourse that will fill the needed slot.” Our algorithms combine grammar rules and lexical items to achieve automated classification. (See Schank, R. 1975. *Conceptual dependency theory. Conceptual Information Processing*, Chap. 3. North-Holland, Amsterdam, The Netherlands, 22–67.)

<sup>EC3</sup> These tag values seem odd, but are used in the algorithms; the numbers are an implementation detail, and may vary across algorithms. There is no special reason for the choice of the numbers used.

<sup>EC4</sup> Stemming is the process of mapping a word to its root word. For example, the root of “buying” is “buy.”

The current size of the lexicon is approximately 300 distinct words. Ongoing, incremental analysis results in additions to the word set.

Based on the training corpus, we can compute the *discriminant value* of each item in the lexicon. This value describes the power of the lexical item in differentiating message types. For example, the word “buy” is likely to be a strong discriminator, since it would be suggestive of positive sentiment. The goal is to populate the lexicon with words that are good discriminators.

### C. Discriminant Values

Example values for some words from the discriminant function are shown here (we report a selection of words only, not the entire lexicon). The last three words appear with their negation tags.

```
SAMPLE DISCRIMINANT VALUES
bad 0.040507943664639216
hot 0.016124148231134897
hype 0.008943543938332603
improve 0.012395140059803732
joke 0.02689751948279659
jump 0.010691670826157351
killing 0.010691670826157329
killed 0.016037506239236058
lead 0.003745650480005731
leader 0.0031710056164216908
like 0.003745470397428718
long 0.01625037430824596
lose 0.12114219092843743
loss 0.007681269362162742
money 0.15378504322023162
oversell 0.0
overvalue 0.016037506239236197
own 0.0030845538644182426
gold_n 0.0
good_n 0.04846852990132937
grow_n 0.016037506239236058
```

These values make for interesting study. For example, the word “lose” understandably has a high discriminant value. The word “oversell” is not used at all. One of the higher values comes from the negated word “good–n” which means that there is plenty of negation in the language used in the message boards. Compare this with its antonym “bad,” which actually has a lower discriminant value! The word “joke” is a good discriminator, which is somewhat surprising, though not totally nonintuitive. The highest valued discriminant is the word “money.”