# Estimating Sentiment Orientation in Social Media for Business Informatics

**Kristin Glass**[1]         **Richard Colbaugh**[1,2]

[1] New Mexico Institute of Mining and Technology, Socorro, NM USA
[2] Sandia National Laboratories, Albuquerque, NM USA

## Abstract

Inferring the sentiment of social media content, for instance blog postings or online product reviews, is both of great interest to businesses and technically challenging to accomplish. This paper presents two computational methods for estimating social media sentiment which address the challenges associated with Web-based analysis. Each method formulates the task as one of text classification, models the data as a bipartite graph of documents and words, and assumes that only limited prior information is available regarding the sentiment orientation of any of the documents or words of interest. The first algorithm is a semi-supervised sentiment classifier which combines knowledge of the sentiment labels for a few documents and words with information present in unlabeled data, which is abundant online. The second algorithm assumes existence of a set of labeled documents in a domain related to the domain of interest, and leverages these data to estimate sentiment in the target domain. We demonstrate the utility of the proposed methods by showing they outperform several standard methods for the task of inferring the sentiment of online reviews of movies, electronics products, and kitchen appliances. Additionally, we illustrate the potential of the methods for multilingual business informatics through a case study involving estimation of Indonesian public opinion regarding the July 2009 Jakarta hotel bombings.

## 1. Introduction

The enormous popularity of "social media", such as blogs, forums, and social networking sites, represents a significant challenge to standard business models and practices, as these media move the control of information from companies to consumers [e.g. 1-5]. However, social media also offer unprecedented opportunities to increase business responsiveness and agility. For example, recent surveys reveal that 32% of the nearly 250 million bloggers worldwide regularly give opinions on products and brands, 71% of active Internet users read blogs, and 70% of consumers trust opinions posted online by other consumers [6,7]. Thus social media is a vast source of business-relevant opinions. Moreover, this information has a reach that rivals any traditional media and an influence which substantially exceeds standard advertising channels.

Businesses are therefore strongly motivated to pay attention to social media and other online information sources. For instance, it is crucially important for companies to be able to rapidly discover and characterize both negative and positive sentiment expressed by current and potential customers. Complaints and other negative views are easier to address if detected quickly, while early positive "buzz" can be reinforced and amplified. Identifying nascent consumer concern or enthusiasm about topics which are relevant to company business can be of great strategic advantage. Indeed, the relevance and timeliness of the information available in social media has the potential to revolutionize the way business is conducted in many sectors.

While monitoring social media is of considerable interest to businesses, performing such analysis is technically challenging. The opinions of consumers are typically expressed as informal communications and are buried in the vast, and largely irrelevant, output of millions of bloggers and other online content producers. Consequently, effectively exploiting these data requires the development of new, automated methods of analysis [1-5]. Although powerful computational analytics have been derived for traditional forms of content, less has been done to develop techniques that are well-suited to the particular characteristics of the content found in social media.

This paper considers one of the central problems in the new domain of social media analytics: deciding whether a given document, such as a blog post or forum thread, expresses positive or negative opinion toward a given topic. The informal nature of social media content poses a challenge for language-based sentiment analysis. While statistical learning-based methods often provide good performance in unstructured settings like this [e.g., 8-15], obtaining the required labeled instances of data, such as a lexicon of sentiment-laden words for a given domain or a collection of "exemplar" blog posts of known polarity, is expensive and time-consuming for Web applications.

We present two new computational methods for inferring sentiment orientation of social media content which address these challenges. Each method formulates the task as

one of text classification, models the data as a bipartite graph of documents and words, and assumes that only limited prior information is available regarding the sentiment orientation of any of the documents or words of interest. The first algorithm adopts a semi-supervised approach to sentiment classification, combining knowledge of the sentiment polarity for a few documents and a small lexicon of words with information present in a corpus of unlabeled documents; note that such unlabeled data are readily obtainable in online applications. The second algorithm assumes existence of a set of labeled documents in a domain related to the domain of interest, and provides a procedure for transferring the sentiment knowledge contained in these data to the target domain. We demonstrate the utility of the proposed algorithms by showing they outperform several standard methods for the task of inferring the sentiment polarity of online reviews of movies, electronics products, and kitchen appliances. Additionally, we illustrate the potential of the methods for multilingual business informatics through a case study involving estimation of Indonesian public opinion regarding the July 2009 Jakarta hotel bombings.

## 2. Preliminaries

We approach the task of estimating the sentiment orientations of a collection of documents as a text classification problem. Each document of interest is represented as a "bag of words" feature vector $x \in \Re^{|V|}$, where the entries of x are the frequencies with which the words in the vocabulary set V appear in the document (perhaps normalized in some way [8]). We wish to learn a vector $c \in \Re^{|V|}$ such that the classifier orient = $sign(c^T x)$ accurately estimates the sentiment orientation of document x, returning +1 (−1) for documents expressing positive (negative) sentiment about the topic of interest.
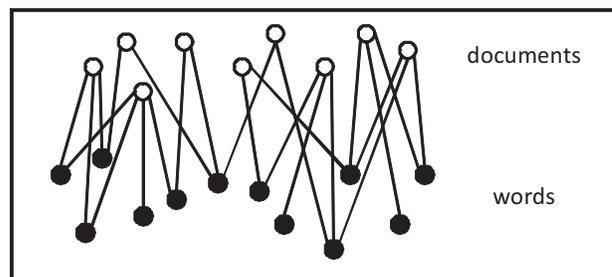
Knowledge-based classifiers leverage prior domain information to construct the vector c. One way to obtain such a classifier is to assemble lexicons of positive words $V^+ \subseteq V$ and negative words $V^- \subseteq V$, and then to set $c_i = +1$ if word $i \in V^+$, $c_i = -1$ if $i \in V^-$, and $c_i = 0$ if i is not in either lexicon; this classifier simply sums the positive and negative sentiment words in the document and assigns document orientation accordingly. While this scheme can provide acceptable performance in certain settings, it is unable to improve its performance or adapt to new domains, and it is usually labor-intensive to construct lexicons which are sufficiently complete to enable useful sentiment classification performance to be achieved.

Alternatively, learning-based methods attempt to generate the classifier vector c from examples of positive and negative sentiment. To obtain a learning-based classifier, one can begin by assembling a set of $n_l$ *labeled* documents $\{(x_i, d_i)\}$, where $d_i \in \{+1, -1\}$ is the sentiment label for document i. The vector c then can be learned through "training" with the set $\{(x_i, d_i)\}$, for instance by solving the following set of equations for c:

$$[X^T X + \gamma I_{|V|}] c = X^T d, \qquad (1)$$

where matrix $X \in \Re^{nl \times |V|}$ has document vectors for rows, $d \in \Re^{nl}$ is the vector of document labels, $I_{|V|}$ denotes the $|V| \times |V|$ identity matrix, and $\gamma \geq 0$ is a constant; this corresponds to regularized least squares (RLS) learning [16]. Many other strategies can be used to compute c, including Naïve Bayes (NB) statistical inference [8]. Learning-based classifiers have the potential to improve their performance and to adapt to new situations, but realizing these capabilities requires that fairly large training sets of labeled documents be obtained and this is usually an expensive undertaking.

Sentiment analysis of social media content for business applications is often characterized by the existence of only modest levels of prior knowledge regarding the domain of interest, reflected in the availability of a few labeled documents and small lexicon of sentiment-laden words, and by the need to rapidly learn and adapt to new domains. As a consequence, standard knowledge-based and learning-based sentiment analysis methods are typically ill-suited for business informatics. In order to address this challenge, the sentiment analysis methods developed in this paper enable limited labeled data to be combined with readily available "auxiliary" information to produce accurate sentiment estimates. More specifically, the first proposed method is a *semi-supervised* algorithm [e.g., 11,12] which leverages a source of supplementary data which is abundant online: unlabeled documents and words. Our second algorithm is a novel *transfer learning* method [e.g., 13] which permits the knowledge present in data that has been previously labeled in a related domain (say movie reviews) to be transferred to a new domain (electronics reviews).



**Figure 1.** Cartoon of bipartite graph model $G_b$, in which documents (white vertices) are connected to the words (black vertices) they contain, and the link weights (black edges) reflect word frequencies.

Each of the algorithms proposed in this paper assumes the availability of a modest lexicon of sentiment-laden words. This lexicon is encoded as a vector $w \in \Re^{|V_l|}$, where $V_l = V^+ \cup V^-$ is the sentiment lexicon and the entries of w are set to +1 or −1 according to the polarity of the corresponding words. The development of the algorithms begins by mod-

eling the problem data as a bipartite graph $G_b$ of documents and words (see Figure 1). It is easy to see that the adjacency matrix A for graph $G_b$ is given by

$$A = \begin{bmatrix} 0 & X \\ X^T & 0 \end{bmatrix} \tag{2}$$

where the matrix $X \in \mathfrak{R}^{n \times |V|}$ is constructed by stacking the document vectors as rows, and each '0' is a matrix of zeros. In both the semi-supervised and transfer learning algorithms, integration of labeled and "auxiliary" data is accomplished by exploiting the relationships between documents and words encoded in the bipartite graph model. The basic idea is to assume that, in the bipartite graph $G_b$, positive/negative documents will tend to be connected to (contain) positive/negative words, and positive/negative words will tend to be connected to positive/negative documents.

## 3. Semi-supervised Sentiment Analysis

We now derive our first sentiment estimation algorithm for social media content. Consider the common situation in which only limited prior knowledge is available about the way sentiment is expressed in the domain of interest, in the form of small sets of documents and words for which sentiment labels are known, but where abundant unlabeled documents can be easily collected (e.g., via Web crawling). In this setting it is natural to adopt a semi-supervised approach, in which labeled and unlabeled data are combined and leveraged in the analysis process. In what follows we present a novel bipartite graph-based approach to semi-supervised sentiment analysis.

Assume the initial problem data consists of a corpus of n documents, of which $n_l \ll n$ are labeled, and a modest lexicon $V_l$ of sentiment-laden words, and suppose that this label information is encoded as vectors $d \in \mathfrak{R}^{nl}$ and $w \in \mathfrak{R}^{|V_l|}$, respectively. Let $d_{est} \in \mathfrak{R}^n$ be the vector of estimated sentiment orientations for the documents in the corpus, and define the "augmented" classifier $c_{aug} = [d_{est}{}^T \quad c^T]^T \in \mathfrak{R}^{n+|V|}$ which estimates the polarity of both documents and words. Note that the quantity $c_{aug}$ is introduced for notational convenience in the subsequent development and is not directly employed for classification. More specifically, in the proposed methodology we learn $c_{aug}$, and therefore c, by solving an optimization problem involving the labeled and unlabeled training data, and then use c to estimate the sentiment of any new document of interest with the simple linear classifier orient = $sign(c^T x)$. We refer to this classifier as *semi-supervised* because it is learned using both labeled and unlabeled data. Assume for ease of notation that the documents and words are indexed so the first $n_l$ elements of $d_{est}$ and $|V_l|$ elements of c correspond to the labeled data.

We wish to learn an augmented classifier $c_{aug}$ with the following three properties: 1.) if a document is labeled, then the corresponding entry of $d_{est}$ should be close to this $\pm 1$ label; 2.) if a word is in the sentiment lexicon, then the

corresponding entry of c should be close to this $\pm 1$ sentiment polarity; and 3.) if there is an edge $X_{ij}$ of $G_b$ that connects a document x and a word $v \in V$ and $X_{ij}$ possesses significant weight, then the estimated polarities of x and v should be similar. These objectives are encoded in the following minimization problem:

$$\min_{c_{aug}} \ c_{aug}^T L c_{aug} + \beta_1 \sum_{i=1}^{n_l} (d_{est,i} - d_i)^2 + \beta_2 \sum_{i=1}^{|V_l|} (c_i - w_i)^2 \tag{3}$$

where $L = D - A$ is the graph Laplacian matrix for $G_b$, with D the diagonal degree matrix for A (i.e., $D_{ii} = \Sigma_j A_{ij}$), and $\beta_1, \beta_2$ are nonnegative constants. Minimizing (3) enforces the three properties we seek for $c_{aug}$, with the second and third terms penalizing "errors" in the first two properties. To see that the first term enforces the third property, observe that this expression is a sum of components of the form $X_{ij}(d_{est,i} - c_j)^2$. The constants $\beta_1, \beta_2$ can be used to balance the relative importance of the three properties.

The $c_{aug}$ which minimizes the objective function (3) can be obtained by solving the following set of linear equations:

$$\begin{bmatrix} L_{11} + \beta_1 I_{nl} & L_{12} & L_{13} & L_{14} \\ L_{21} & L_{22} & L_{23} & L_{24} \\ L_{31} & L_{32} & L_{33} + \beta_2 I_{|V_l|} & L_{34} \\ L_{41} & L_{42} & L_{43} & L_{44} \end{bmatrix} c_{aug} = \begin{bmatrix} \beta_1 d \\ 0 \\ \beta_2 w \\ 0 \end{bmatrix} \tag{4}$$

where the $L_{ij}$ are matrix blocks of L of appropriate dimension. The system (4) is sparse because the data matrix X is sparse, and therefore large-scale problems can be solved efficiently. Note that in situations where the set of available labeled documents and words is *very* limited, sentiment classifier performance can be improved by replacing L in (4) with the normalized Laplacian $L_n = D^{-1/2} L D^{-1/2}$, or with a power of this matrix $L_n{}^k$ (for k a positive integer).

We summarize this discussion by sketching an algorithm for learning the proposed semi-supervised classifier:
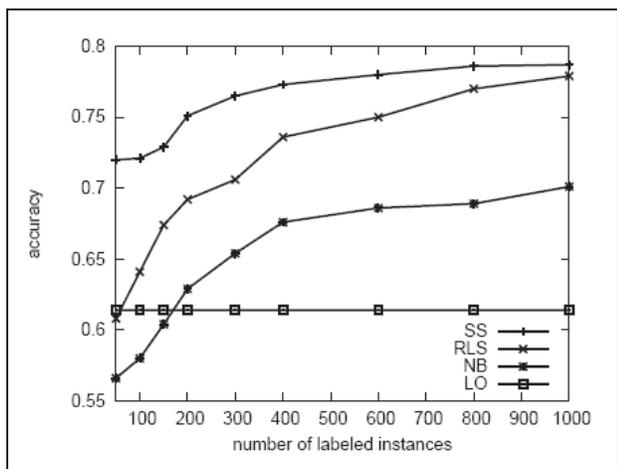
**Algorithm SS:**

1. Construct the set of equations (4), possibly by replacing the graph Laplacian L with $L_n{}^k$.

2. Solve equations (4) for $c_{aug} = [d_{est}{}^T \quad c^T]^T$ (for instance using the Conjugate Gradient method).

3. Estimate the sentiment orientation of any new document x of interest as: orient = $sign(c^T x)$.

The utility of Algorithm SS is now examined through a case study involving sentiment estimation for online movie reviews, a task which is representative of many business informatics applications.

## 4. Case Study One: Movie Reviews

This case study examines the performance of Algorithm SS for the problem of estimating sentiment of online movie reviews. The data used in this study is a publicly available set of 2000 movie reviews, 1000 positive and 1000 negative, collected from the Internet Movie Database and archived at the website [17]. The Lemur Toolkit [18] was employed to construct the data matrix X and vector of document labels d from these reviews. A lexicon of ~1400 domain-independent sentiment-laden words was obtained from [19] and employed to build the lexicon vector w.

This study compares the movie review orientation classification accuracy of Algorithm SS with that of three other schemes: 1.) lexicon-only, in which the lexicon vector w is used as the classifier as summarized in Section 2, 2.) a classical NB classifier obtained from [20], and 3.) a well-tuned version of the RLS classifier (1). Algorithm SS is implemented with the following parameter values: $\beta_1$ = 0.1, $\beta_2$ = 0.5, and k = 10. A focus of the investigation is evaluating the extent to which good sentiment estimation performance can be achieved even if only relatively few labeled documents are available for training; thus we examine training sets which incorporate a range of numbers of labeled documents: $n_l$ = 50, 100, 150, 200, 300, 400, 600, 800, 1000.



**Figure 2.** Results for the movie reviews case study. The plot shows how sentiment estimation accuracy (vertical axis) varies with number of available labeled movie reviews (horizontal axis) for four different classifiers: lexicon only (LO), NB, RLS, and Algorithm SS (SS).

Sample results from this study are depicted in Figure 2. Each data point in the plots represents the average of ten trials. In each trial, the movie reviews are randomly partitioned into 1000 training and 1000 test documents, and a randomly selected subset of training documents of size $n_l$ is "labeled" (i.e., the labels for these reviews are made available to the learning algorithm). As shown in Figure 2, Al-

gorithm SS outperforms the other three methods. Note that, in particular, the accuracy obtained with the proposed approach is significantly better than the other techniques when the number of labeled training documents is small. It is expected that this property of Algorithm SS will be of considerable value in business informatics applications that involve social media data.

## 5. Transfer Learning Sentiment Analysis

This section develops the second proposed sentiment estimation algorithm for social media content. Many business informatics applications are characterized by the presence of limited labeled data for the domain of interest but ample labeled information for a related domain. For instance, a firm may wish to ascertain the sentiment of online discussions about its new line of kitchen appliances, and may have in hand a large set of labeled examples of positive and negative reviews for its electronics products (e.g., from studies of previous product launches). In this setting it is natural to adopt a transfer learning approach, in which knowledge concerning the way sentiment is expressed in one domain, the so-called *source* domain, is transferred to permit sentiment estimation in a new *target* domain. In what follows we present a new bipartite graph-based approach to transfer learning-based sentiment analysis.

Assume that the initial problem data consists of a corpus of $n = n_T + n_S$ documents, where $n_T$ is the (small) number of labeled documents available for the target domain of interest and $n_S \gg n_T$ is the number of labeled documents from some related source domain; in addition, suppose that a modest lexicon $V_l$ of sentiment-laden words is known. Let this label data be encoded as vectors $d_T \in \mathfrak{R}^{nT}$, $d_S \in \mathfrak{R}^{nS}$, and $w \in \mathfrak{R}^{|Vl|}$, respectively. Denote by $d_{T,est} \in \mathfrak{R}^{nT}$, $d_{S,est} \in \mathfrak{R}^{nS}$, and $c \in \mathfrak{R}^{|Vl|}$ the vectors of estimated sentiment orientations for the target and source documents and the words, and define the augmented classifier as $c_{aug} = [d_{S,est}^T \quad d_{T,est}^T \quad c^T]^T \in \mathfrak{R}^{n+|V|}$. Note that the quantity $c_{aug}$ is introduced for notational convenience in the subsequent development and is not directly employed for classification.

In what follows we derive an algorithm for learning $c_{aug}$, and therefore c, by solving an optimization problem involving the labeled source and target training data, and then use c to estimate the sentiment of any new document of interest via the simple linear classifier orient = $sign(c^T x)$. This classifier is referred to as *transfer learning-based* because c is learned, in part, by transferring knowledge about the way sentiment is expressed from a domain which is related to (but need not be identical to) the domain of interest.

We wish to learn an augmented classifier $c_{aug}$ with the following four properties: 1.) if a source document is labeled, then the corresponding entry of $d_{S,est}$ should be close to this ±1 label; 2.) if a target document is labeled, then the corresponding entry of $d_{T,est}$ should be close to this ±1 label, and the information encoded in $d_T$ should be emphasized

relative to that in the source labels $d_S$; 3.) if a word is in the sentiment lexicon, then the corresponding entry of c should be close to this $\pm 1$ sentiment polarity; and 4.) if there is an edge $X_{ij}$ of $G_b$ that connects a document x and a word $v \in V$ and $X_{ij}$ possesses significant weight, then the estimated polarities of x and v should be similar.

The four objectives listed above may be realized by solving the following minimization problem:

$$\min_{c_{aug}} \quad c_{aug}^T L c_{aug} + \beta_1 \left\| d_{S,est} - k_S d_S \right\|^2 + \beta_2 \left\| d_{T,est} - k_T d_T \right\|^2$$

$$+ \beta_3 \left\| c - w \right\|^2 \qquad (5)$$

where $L = D - A$ is the graph Laplacian matrix for $G_b$, as before, and $\beta_1$, $\beta_2$, $\beta_3$, $k_S$, and $k_T$ are nonnegative constants. Minimizing (5) enforces the four properties we seek for $c_{aug}$. More specifically, the second, third, and fourth terms penalize "errors" in the first three properties, and choosing $\beta_2 > \beta_1$ and $k_T > k_S$ favors target label data over source labels. To see that the first term enforces the fourth property, note that this expression is a sum of components of the form $X_{ij} (d_{T,est,i} - c_j)^2$ and $X_{ij} (d_{S,est,i} - c_j)^2$. The constants $\beta_1$, $\beta_2$, $\beta_3$ can be used to balance the relative importance of the four properties.

The $c_{aug}$ which minimizes the objective function (5) can be obtained by solving the following set of linear equations:

$$\begin{bmatrix} L_{11} + \beta_1 I_{nS} & L_{12} & L_{13} \\ L_{21} & L_{22} + \beta_2 I_{nT} & L_{23} \\ L_{31} & L_{32} & L_{33} + \beta_3 I_{|V_1|} \end{bmatrix} c_{aug} = \begin{bmatrix} \beta_1 k_S d_S \\ \beta_2 k_T d_T \\ \beta_3 w \end{bmatrix}$$

$$(6)$$

where the $L_{ij}$ are matrix blocks of L of appropriate dimension. The system (6) is sparse because the data matrix X is sparse, and therefore large-scale problems can be solved efficiently. In situations where the set of available labeled documents and words is very limited, sentiment classifier performance can be improved by replacing L in (6) with the normalized Laplacian $L_n$ or with a power of this matrix $L_n^k$.

We summarize the above discussion by sketching an algorithm for learning the proposed transfer learning classifier:
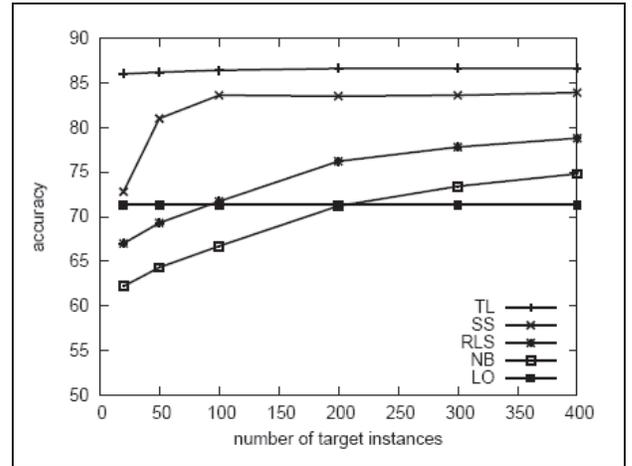
**Algorithm TL:**

1. Construct the set of equations (6), possibly by replacing the graph Laplacian L with $L_n^k$.
2. Solve equations (6) for $c_{aug} = [d_{S,est}^T \quad d_{T,est}^T \quad c^T]^T$.
3. Estimate the sentiment orientation of any new document x of interest as: orient = $\text{sign}(c^T x)$.

The utility of Algorithm TL is now examined through a case study involving sentiment estimation for online product reviews.

## 6. Case Study Two: Product Reviews

This case study examines the performance of Algorithm TL for the problem of estimating sentiment of online product reviews. The data used in this study is a publicly available set of 1000 reviews of electronics products, 500 positive and 500 negative, and 1000 reviews of kitchen appliances, 500 positive and 500 negative, collected from Amazon and archived at the website [21]. The Lemur Toolkit [18] was employed to construct the data matrix X and vectors of document labels $d_S$ and $d_T$ from these reviews. A lexicon of 150 domain-independent sentiment-laden words was constructed manually and employed to form the lexicon vector w.

This study compares the product review sentiment classification accuracy of Algorithm TL with that of four other strategies: 1.) lexicon-only, in which the lexicon vector w is used as the classifier as summarized in Section 2, 2.) a classical NB classifier obtained from [20], 3.) a well-tuned version of the RLS classifier (1), and 4.) Algorithm SS. Algorithm SS is implemented with the following parameter values: $\beta_1 = 1.0$, $\beta_2 = 3.0$, $\beta_3 = 5.0$, $k_S = 0.5$, $k_T = 1.0$, and $k = 5$. A focus of the investigation is evaluating the extent to which the knowledge present in labeled reviews from a related domain, here kitchen appliances, can be transferred to a new domain for which only limited labeled data is available, in this case electronics. Thus we assume that all 1000 labeled kitchen reviews are available to Algorithm TL (the only algorithm which is designed to exploit this information), and examine training sets which incorporate a range of numbers of labeled documents from the electronics domain: $n_T = 20, 50, 100, 200, 300, 400$.



**Figure 3.** Results for the consumer product reviews case study. The plot shows how sentiment estimation accuracy (vertical axis) varies with number of available labeled electronics reviews (horizontal axis) for five different classifiers: lexicon only (LO), NB, RLS, Algorithm SS (SS), and Algorithm TL (TL).

Sample results from this study are depicted in Figure 3. Each data point in the plots represents the average of ten trials. In each trial, the electronics reviews are randomly partitioned into 500 training and 500 test documents, and a randomly selected subset of reviews of size $n_T$ is extracted from the 500 labeled training instances and made available to the learning algorithms. As shown in Figure 3, Algorithm TL outperforms the other four methods. Note that, in particular, the accuracy obtained with the transfer learning approach is significantly better than the other techniques when the number of labeled training documents in the target domain is small. It is expected that the ability of Algorithm TL to exploit knowledge from a related domain to quickly learn an effective sentiment classifier for a new domain will be of considerable value in business informatics applications involving social media data.
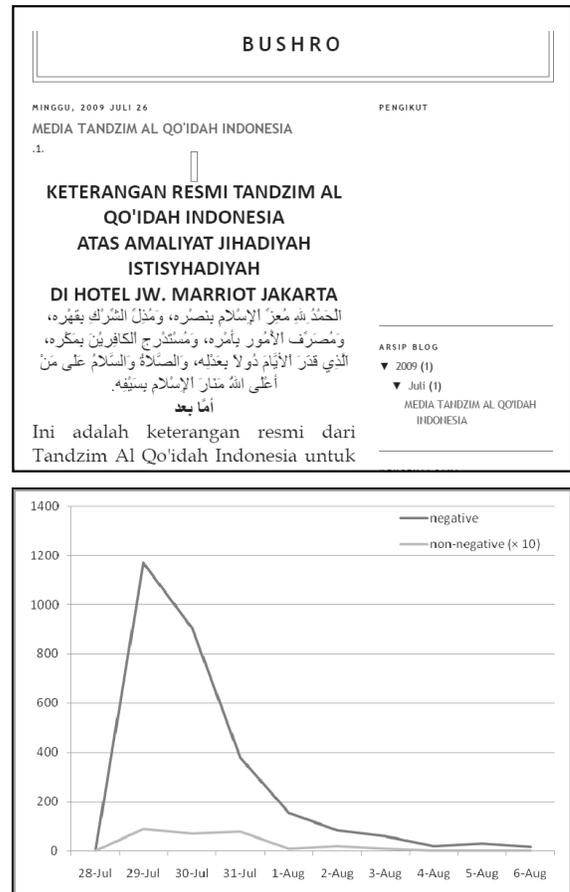
## 7. Case Study Three: Public Opinion

On 17 July 2009 the JW Marriott and Ritz-Carlton Hotels in Jakarta, Indonesia were hit by suicide bombing attacks within five minutes of each other. A little over a week later, on 26 July 2009, a document claiming responsibility for the attacks and allegedly written by N.M. Top was posted on the blog [22]; see Figure 4 for a screenshot of a portion of this blog post. In subsequent discussions we will refer to this post as the "Top post" for convenience, with the understanding that the authorship of the post is uncertain. At the time, businesses in the travel and lodging sectors expressed interest in understanding sentiment in the region regarding the bombings and the alleged claim of responsibility by a well-known extremist.

To enable a preliminary assessment along these lines, we collected two sets of social media data related to the Top post: 1.) ~3000 comments made to the post during the two week period immediately following its publication, and 2.) several hundred posts made to other Indonesian language blogs in which the Top post was discussed. We manually labeled the sentiment of a small subset of these documents, and translated into Indonesian the generic sentiment lexicon used in Case Study One for implementation in this study. Observe that this approach to constructing a sentiment lexicon is far from perfect. However, because our proposed algorithms employ several sources of information to estimate the sentiment of content, it is expected that they will exhibit robustness to imperfections in any single data source. This study therefore offers the opportunity to explore the utility of a very simple approach to multilingual sentiment analysis: translate a small lexicon of sentiment-laden words into the language of interest and then apply Algorithm SS or Algorithm TL directly within that language (treating words as tokens). The capability to perform automated, multilingual content analysis is of substantial interest in modern business operations.

We used Algorithm SS to estimate the sentiment expressed in the corpus of comments made to the Top post [22] and in the set of related discussions posted at other blogs. This analysis revealed that the comments made directly to the Top post are almost universally negative, condemning both the bombings and the justification for the bombings given in the Top post (see Figure 4). Manual examination of a subset of the comments confirms the results provided by Algorithm SS. Analysis of relevant posts made to other blogs also indicated that sentiment about the Top post is largely negative in these forums, although they are not as uniformly negative as the comments made on the blog site [22].



**Figure 4.** Results for public opinion case study. Image at top is a screenshot of the blog post, allegedly by N.M. Top, which claims responsibility for the July 2009 bombings of two hotels in Jakarta, Indonesia. Plot at bottom shows the estimated sentiment of comments made directly to the blog [22] in response to the Top post (dark grey is negative, light grey is "nonnegative", including positive posts, with light grey multiplied by a factor of 10 to be visible on the plot).

## Acknowledgements

# References

1. Glance, N., M. Hurst, K. Nigam, M. Siegler, R. Stockton, and T. Tomokiyo, "Deriving marketing intelligence from online discussion", *Proc. 11th ACM International Conference on Knowledge Discovery and Data Mining*, Chicago, August 2005.
2. Ziegler, C. and M. Skubacz, "Towards automated reputation and brand monitoring on the Web", *Proc. IEEE/ ACM International Conference on Web Intelligence*, Hong Kong, December 2006.
3. Melville, P., V. Sindhwani, and R. Lawrence, "Social media analytics: Channeling the power of the blogosphere for marketing insight", *Proc. Workshop on Information in Networks*, New York, September 2009.
4. Kaplan, A. and M. Haenlein, "Users of the world, unite! The challenges and opportunities of social media", *Business Horizons*, Vol. 53, pp. 59-68, 2010.
5. van der Lans, R., G. van Bruggen, J. Eliashberg, and B. Wierenga, "A viral branching model for predicting the spread of electronic word of mouth", *Marketing Science*, Vol. 29, pp. 348-365, 2010.
6. http://www.universalmccann.com, accessed July 2010.
7. http://www.nielson.com, accessed July 2010.
8. Pang, B. and L. Lee, "Opinion mining and sentiment analysis", *Foundations and Trends in Information Retrieval*, Vol. 2 , pp. 1-135, 2008.
9. Dhillon, I., "Co-clustering documents and words using bipartite spectral graph partitioning", *Proc. ACM International Conference on Knowledge Discovery and Data Mining*, San Francisco, August 2001.
10. Kim, S. and E. Hovy, "Determining the sentiment of opinions", *Proc. International Conference on Computational Linguistics*, 2004.
11. Sindhwani, V. and P. Melville, "Document-word co-regularization for semi-supervised sentiment analysis", *Proc. 2008 IEEE International Conference on Data Mining*, Pisa, Italy, December 2008.
12. Colbaugh, R. and K. Glass, "Estimating sentiment orientation in social media for intelligence monitoring and analysis", *Proc. 2010 IEEE International Conference on Intelligence and Security Informatics*, Vancouver, BC Canada, May 2010.
13. Pan, S. and Q. Yang, "A survey on transfer learning", *IEEE Trans. Knowledge and Data Engineering*, Vol. 22, pp. 1345-1359, 2010.
14. Blitzer, J., M. Dredze, and F. Perieia, "Biographies, bollywood, boom-boxes, and blenders: Domain adaptation for sentiment classification", *Proc. 45th Annual Meeting of the ACL*, Prague, June 2007.
15. He, J., Y. Liu, and R. Lawrence, "Graph-based transfer learning", *Proc. 18th ACM Conference on Information and Knowledge Management*, Hong Kong, November 2009.
16. Hastie, T., R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, Second Edition, Springer, New York, 2009.
17. http://www.cs.cornell.edu/People/pabo/movie-review-data/, accessed Dec. 2009.
18. http://www.lemurproject.org/, accessed Dec. 2009.
19. Ramakrishnan, G., A. Jadhav, A. Joshi, S. Chakrabarti, and P. Bhattacharyya, "Question answering via Bayesian inference on lexical relations", *Proc. 41st Annual Meeting of the ACL*, 2003.
20. http://www.borgelt.net/bayes.html, accessed Dec. 2009.
21. http://www.cs.jhu.edu/~mdredze/, accessed Dec. 2010.
22. www.mediaislam-bushro.blogspot.com, accessed Dec. 2009.