

Automated Sentiment Analysis

Eugene Butler
University of Minnesota Morris
600 E 4th Street
Morris, Minnesota 56267
butle250@morris.umn.edu

ABSTRACT

Automated Sentiment Analysis refers to the computerized processing of text in order to determine the sentiments, the attitudes, thoughts, and judgments, of the author. Unfortunately, a number of pitfalls confound the accurate analysis of the sentiments that are conveyed by online statements. In many cases, a statement's phrasing is complex or ambiguous. The degree to which certain statements represent positive or negative sentiment is frequently informed by context that is not easily sensed by many common modes of sentiment analysis. For example, normally negative words may indicate positive sentiment when voiced about the quality or believability of a fictional villain. Also, positive statements about an actor or a particular tourist destination may be meant to contrast with an overall negative impression of the whole movie or vacation.

This paper gives a brief overview of the sentiment analysis field and introduces four modern means of sentiment analysis: Semantic Orientation from Pointwise Mutual Information measurement (SO-PMI), Semantic Orientation from Latent Semantic Analysis (SO-LSA), adjective conjunction measurement, and the Natural Language Processing Combined Method. It outlines the means by which each method determines the sentiment inherent in a given sample and the efficiency and effectiveness of each mode of sentiment analysis. It will then show that context-awareness, specifically the type offered in statement-oriented natural language processing, is essential to genuinely reliable sentiment analysis.

Categories and Subject Descriptors

I.27 [Natural Language Processing]: Text analysis; H.31 [Content Analysis and Indexing]: Linguistic processing

Keywords

sentiment analysis, semantic orientation, semantic association, text mining

This work is licensed under the Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/3.0/us/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.

1. INTRODUCTION

Traditional means of judging public opinion, such as customer satisfaction surveys, are expensive, cumbersome, and of limited effectiveness [4]. Among the host of problems that plague traditional customer surveys are the difficulty of reaching a sufficiently large survey population and developing effective survey questions. Reaching a sufficiently large survey population is time-consuming, expensive, and creates logistical difficulties. Developing questionnaires that reveal actionable information about the average customer's opinion, without biasing the respondent, overlooking important factors, or both, is also very difficult. Fortunately, the existence of the internet has presented a cheaper, more effective way of judging public opinion.

The advent of the internet has meant that individuals are increasingly able to post their opinions on public forums, resulting in a new wealth of free, yet largely uncatalogued, opinion information. Furthermore, these opinions, in the form of online customer feedback, influence the decisions of other customers [2], making the ability to easily know the general sentiment of this feedback extremely valuable. While manual attempts to augment traditional survey data with internet opinion are costly and time-intensive, the mining of internet opinion through automated sentiment analysis is genuinely effective [2]. By using a computer to determine the positive or negative sentiment inherent in each piece of feedback, one can rapidly develop an understanding of how reviewers, in general, feel about the product.

The purpose of this paper is to examine the efficacies of four modern means of sentiment analysis: Semantic Orientation from Pointwise Mutual Information measurement (SO-PMI), Semantic Orientation from Latent Semantic Analysis (SO-LSA), adjective conjunction measurement, and Natural Language Processing. The efficacy of each mode is defined by its ability to accurately analyze the text in order to interpret its "semantic orientation." Semantic orientation refers to the positive or negative connotation of the text and the degree to which these connotations are carried. 'A phrase has a positive semantic orientation when it has good associations (e.g. "subtle nuances") and a negative semantic orientation when it has bad associations (e.g., "very cavalier").' [6]. Thus, the semantic orientation decoded from a text is a direct reflection of the sentiments of the writer. Because of this, once the semantic orientation of a work is understood, one can subsequently explain how, and how strongly, a writer feels about the subject.

The history of sentiment analysis work is defined by a constant struggle against the complexity, ambiguity, and imprecision of human communication. From the simple enumeration of “good”- and “bad”-indicating words to the analysis of opinion-target pairs to attempts to analyze the language of entire sentences at a time, each sentiment analysis technique has attempted to more reliably control for these confounding factors. However, even efforts to make the simplest judgments, to determine a binary “like” or “dislike” of a particular subject, are heavily complicated by these problems.

Common attempts at sentiment analysis involve the scanning of text for words thought to indicate “grammatical polarity,” the positive or negative attitude of the text with regard to its subject (also referred to as “semantic orientation”). For example, “My lunch was disgusting and unsatisfying.” would be weighted with a negative polarity because of the words “disgusting” and “unsatisfying.” [6] In [6], Turney found that this type of analysis was 84% accurate for automobile reviews, 80% for bank reviews, but only 65% accurate for movie reviews, a discrepancy that could not simply be explained through the use of less superlative or enthusiastic language. One reason for this discrepancy was the use of “negative” words in a positive context. For instance, a positive example of a movie villain is one with negative qualities capable of disgusting, horrifying, or discomforting the viewer. An example is the phrase “The slow, methodical way he spoke. I loved it! It made him seem more arrogant and even more evil.” While the passage contains “loved it!”, the use of the words “more evil” resulted in the review being judged to have a negative grammatical polarity. Unfortunately, the review was actually a five-star recommendation of the movie *The Matrix* that cited the believability of the villain Agent Smith [6].

Another problem, one that confounds multiple schools of sentiment analysis, is the existence of positive parts in negative wholes. A bad movie, for example, can have good actors. Or, a bad movie may be an unusual departure for a normally good actor, director, or writer. The review “Well as usual Keanu Reeves is nothing special, but surprisingly, the very talented Laurence Fishbourne is not so good either, I was surprised[sic]” was a two-star non-recommendation, but was scored as a recommendation. Turney speculates that this may also be to blame for the lower accuracy of some travel reviews: “good beaches do not necessarily add up to a good vacation. On the other hand, good automotive parts usually do add up to a good automobile and good banking services add up to a good bank.” [6] The inability to determine contextual polarity is a major weakness of the SO-PMI and SO-LSA methods. Fortunately, attempts have been made to mitigate this problem through the development of adjective-conjunction judging and through the augmentation of SO-PMI through the Natural Language Processing Combined Method.

This paper will review four means of sentiment analysis: Semantic Orientation from Pointwise Mutual Information measurement (SO-PMI), Semantic Orientation from Latent Semantic Analysis (SO-LSA), adjective conjunction measurement, and the Natural Language Processing Combined Method. It will outline the advantages, disadvantages, efficiency, and effectiveness of each method, and the ways in

which each method contributes to the field of sentiment analysis. It will then show, through example, that the type of context-awareness offered in statement-oriented natural language processing is essential to truly reliable sentiment analysis.

2. SEMANTIC ASSOCIATION

Semantic association is predicated on the idea that “a word is characterized by the company it keeps.” [8] In other words, a word’s semantic orientation tends to correspond with the semantic orientation of its neighbors [7]. In order to be able to find the semantic orientation of a new word, one must first have a list of “known-positive” or “known-negative” words against which it can be compared. Examples include the “positive” words “good, nice, excellent, positive, fortunate, correct, and superior,” and the “negative” words “bad, nasty, poor, negative, unfortunate, wrong, and inferior.” These words were specifically chosen by [7] for their lack of sensitivity to context, and the near-universality with which they were applied by human test subjects asked to label a number of test words. In confirmation of the results of [1], “the average agreement among subjects was 98% and the average agreement between the subjects and our benchmark labels was 94% (35 subjects, 28 words). This level of agreement compares favourably with validation studies in similar tasks, such as word sense disambiguation.” [7]

In order to compare these “known” words with neighbors, one must also provide a *corpus*, a large collection of text. An example is the 100 billion word AltaVista Advanced Search engine English language page index. The purpose of a corpus is to provide a sufficient number of instances of the target, known-positive, and known-negative words so that the frequency with which it is found near the words of known orientation can be found. The size of the corpus positively correlates with the accuracy of attempts to find semantic association, as it provides a larger sample from which to draw information about word nearness.

Once sets of positive and negative words are established, the orientation of a word can be calculated as follows:

$$SO-A(word) = \sum_{pword \in Pwords} A(word, pword) - \sum_{nword \in Nwords} A(word, nword)$$

where A is a measure of distance between the words, $pword$ and $nword$ are known positive and negative words, and $Pwords$ and $Nwords$ are the sets of known positive and negative words.

The semantic orientation is defined by the degree to which the word is associated with positive words minus the degree to which it is associated with negative words. For example, the phrase “wonderful service” could have a semantic orientation of 2.35. This would be composed of the sum of its nearness to each of the set of positive words, .83 with “good”, .35 with “nice”, etc, minus its nearness to each of the negative words, .01 to “bad”, etc. The result being positive indicates a positive association and, therefore, positive orientation. When the number is negative, it has a negative orientation. The absolute value of SO-A represents the strength of the orientation.

2.1 Semantic Orientation through Pointwise Mutual Information Measurement and Information Retrieval

In statistics, Pointwise Mutual Information is a way to measure the association of two outcomes by considering their coincidence (or lack thereof).

$$PMI(x, y) = \log \frac{p(x, y)}{p(x)p(y)}$$

where p is the probability of a given occurrence or concurrence.

The probability that the outcomes will coincide is defined by $p(x, y)$. If it is assumed that the outcomes are statistically independent, the probability that they'll coincide is $p(x) * p(y)$. The log of the ratio of these two values, $PMI(x, y)$ is positive when x and y correlate, or are likely to coincide, is zero if x and y are independent, and negative if they negatively correlate, if the occurrence of one makes the occurrence of the other less likely.

This technique can be extended to questions of semantic orientation through the use of the variables x and y to represent word occurrence.

$$PMI(word_1, word_2) = \log \frac{p(word_1, word_2)}{p(word_1)p(word_2)}$$

where p is the probability of a given occurrence or concurrence and where $word_1$ and $word_2$ represent words drawn from the corpus or lexicon.

The probabilities can be calculated through the use of PMI-IR, or pointwise mutual information through information retrieval. "PMI-IR estimates PMI by issuing queries to a search engine (hence the IR in PMI-IR) and noting the number of hits (matching documents)." [7] Turney's experiment used the AltaVista search engine because of its support of the NEAR operator, which constrains search results to those documents in which the words supplied are within ten words of each other, on either side, limiting the search results to those in which the elements are close together, and can therefore be "characterized by the company they keep."

The notation $p(word_1, word_2)$ describes the probability that $word_1$ and $word_2$ correlate. It is one over the total number of pages relevant to the sample times the number of times $word_1$ is near $word_2$. Put simply, it is the number of search results, "hits" for $word_1$ within 10 words proximity of $word_2$ divided by the total number of pages on which each word appears. The denominator describes the probability that $word_1$ and $word_2$ would correlate if assumed to be independent.

$$PMI(word_1, word_2) = \log \frac{\frac{1}{N} hits(word_1 \text{ NEAR } word_2)}{\frac{1}{N} hits(word_1) \frac{1}{N} hits(word_2)}$$

where N is the total number of pages in which the words appear, $word_1$ and $word_2$ represent words drawn from the corpus or lexicon, and $hits$ is the number of hits produced by a search for $word_1$ NEAR (within 10 words of) $word_2$.

Finally:

$$SO-PMI(word) = \sum_{pword \in Pwords} PMI(word, pword) - \sum_{nword \in Nwords} PMI(word, nword)$$

where PMI is as previously defined, $pword$ and $nword$ are known positive and negative words within the sets $Pwords$ and $Nwords$, the sets of known positive and negative words.

SO-PMI, being the sum of PMI comparison of known positive and negative words across the corpus, represents the semantic orientation of the text as a whole. When SO-PMI is positive, the semantic orientation of text is positive. When SO-PMI is negative, the semantic orientation of the text is negative. The degree to which SO-PMI is positive or negative, its absolute value, is the strength which the orientation of the text is conveyed. Higher absolute values represent stronger semantic orientation and lower absolute values represent weaker ones.

2.2 SO-PMI Results

The results of SO-PMI tests are encouraging in some ways and discouraging in others. Given a large corpus, such as the 350 million page set of all English-language pages indexed by AltaVista (conservatively estimated at at least 100 billion words)[5], SO-PMI's accuracy vs manual tagging is 87.13% over 100% of the result set. When applied to only the most confident judgments, as informed by a high absolute value for SO-PMI, SO-PMI's accuracy increases to 98.20% over the most confident quartile of its results. When a smaller corpus is used, such as the Touchstone Applied Science Associates (TASA) 10 million word [3] set of short English documents, the accuracy plummets to 61.26% over 100% of the result set, 47.33% over the most confident 50%, and 69.74% over the most confident quartile. Not only is the accuracy lower, but the stability of SO-PMI, the degree to which its accuracy correlates with confidence, drops as well[7]. The accuracy of SO-PMI over the TASA corpus varies without meaningful correlation with confidence.

When test data is grouped by topic rather than corpus size, SO-PMI shows 83.78% accuracy over car reviews, but only 65.83% accuracy over movie reviews [6]. This discrepancy has been attributed to the increased role of context in movie reviews and SO-PMI's lack of contextual awareness. Movie recommendations often contain descriptions of unpleasant scenes or characters while negative reviews may mention pleasant scenes. Unfortunately, SO-PMI is unable to distinguish between the "goodness" or "badness" of elements and wholes and simply sums the PMI of every descriptor to find the whole.

Finally, this method, for maximum accuracy and stability, relies on an extremely large corpus. The disk space and processing time costs associated with running SO-PMI across all English-language pages indexed by AltaVista is substantial, and accuracy and stability quickly decrease with smaller corpus size. That said, it is more accurate than more efficient methods when given large corpora [7].

2.3 Semantic Orientation through Latent Semantic Analysis

Latent Semantic Analysis (LSA) is another means of finding the semantic association between a pair of words. LSA uses Singular Value Decomposition (SVD) to analyze the statistical relationships between words in a corpus [7]. The value produced by the LSA of a pair of words can be passed to SO-A to find the total semantic orientation of the text.

The first step in LSA is to construct a matrix X such that the row vectors represent words and the column vectors represent “chunks of text.” The chunks can be sentences, paragraphs, or documents. Each cell in the matrix is made to represent the “weight” of each word in the corresponding text. This weight is usually the pair’s tf-idf score. The tf-idf represents the “term frequency” times the “inverse document frequency”.

$$tf\text{-}idf_{i,j} = tf_{i,j} * idf_i$$

where tf is the term frequency of a term within a particular document, i and j are word indices, and idf is the inverse document frequency, or the importance of the term i defined by its rarity.

$$idf_{i,j} = \frac{n_{i,j}}{\sum_k n_{i,j}}$$

where the numerator, $n_{i,j}$, is the number of times the term (t_i) appears in the document (d_j) and the denominator, $\sum_k n_{i,j}$, is the total size of the document, the sum of the number of occurrences of every item in the document.

$$idf_i = \log \frac{|D|}{|\{d : t_i \in d\}|}$$

where $|D|$ is the number of documents in the corpus and $|\{d : t_i \in d\}|$ represents the number of documents in which $n_{i,j}$ is nonzero. This represents the number of documents in which the word i appears.

Finding the tf-idf score rather than basic word count is necessary for two main reasons. First, using word frequency prevents a bias towards large documents, which would likely have higher counts of relevant words regardless of their actual relevance. Second, some words appear much more frequently throughout the corpus than others. Therefore, the importance of each word is determined by its Inverse Document Frequency, the degree to which it does not appear in every document. Each word’s frequency times its importance determines its tf-idf score, its “weight” in the matrix X .

The second step involves decomposing $m \times n$ matrix X , composed of the tf-idf scores, through Singular Value Decomposition (SVD). SVD describes the factorization of the matrix X into:

$$X = U \Sigma V^T$$

where $m \times m$ matrix U and $n \times n$ matrix V (V^T being the transpose of V) are orthonormal matrices and Σ is a nonnegative diagonal matrix of singular values.

LSA uses the matrix:

$$\hat{X} = U_k \Sigma_k V_k^T$$

where, when \hat{X} is of rank r , ($k < r$), a “smoothed” or “compressed” version meant to reduce approximation errors [7]. It measures the similarity of two words, $word_1$ and $word_2$ by taking the cosine of the angles between the corresponding row vectors in \hat{X} . As with SO-PMI,

$$SO\text{-}LSA(word) = \sum_{pword \in Pwords} LSA(word, pword) - \sum_{nword \in Nwords} LSA(word, nword)$$

2.4 SO-LSA Results

Unfortunately, SO-LSA does not yet function on larger corpora, such as the AV-ENG (all English-language web pages indexed by AltaVista) or AV-CA (all English-language web pages in the Canadian TLD indexed by Alta-Vista), so comparisons can only be made when using SO-PMI’s weakest corpora. Over small corpora such as the TASA set, the results of SO-LSA are more reliable than those of SO-PMI. Over the most confident quartile of the test set, SO-LSA achieved 81.98% accuracy, while SO-PMI achieved 68.74%. Over the whole set, SO-LSA achieved 65.72% accuracy while SO-PMI achieved 61.26%. In addition, SO-LSA was much more stable over the reduced corpus. Accuracy steadily rose with confidence, while the same was not true of SO-PMI.

3. CONJUNCTION

The third attempt to increase the reliability of sentiment analysis comes in the form of the analysis of adjective conjunctions. It has been observed that conjunctions imply important information about the orientation of their arguments [1]. For example, in these three sentences:

1. “The tax proposal was simple and well-received by the public.”
2. “The tax proposal was simplistic but well-received by the public.”
3. “The tax proposal was simplistic and well-received by the public.”

The semantic orientation of the adjectives is informed by the use of the conjunctions “and” and “but.” The use of “and” implies that the two adjectives are of similar, complimentary, orientation. The use of “but” implies that the two adjectives are of opposite, dissonant orientations. In sentence number one, the use of “simple” and “well-received” is clear and intuitive. The tax proposal is both *simple* and *well-received* and both “simple” and “well-received” are “good” words. In sentence two, “but” is used to separate words of different orientation. “Well-received” is good, but “simplistic” is bad. The third example illustrates how a misused adjective conjunction is easily apparent. Connecting adjectives of alternate orientation with “and” is contradictory and grammatically incorrect. This is similar to the conjunction of words of consonant orientation with “but.” The statement “The subject is good but good.” has unnecessary redundancy. Because of these factors, uses of “and” with adjectives of dissonant orientation and “but” with consonant adjectives are extremely rare.

The conjunction-judging system in [1] has four main stages. In the first, it extracts conjunctions and adjectives from the corpus. Hatzivassiloglou used the 1987 Wall Street Journal corpus, a collection of 21 million words automatically annotated with part-of-speech tags, then selected all adjectives that appeared at least 20 times and removed all labeled adjectives that had no orientation, such as “domestic” and “medical.” [1] Then, all adjectives are given a binary, positive or negative orientation label. Adjectives that can ascribe multiple qualities based on context, such as “cheap” were discarded. The conjunctions were then extracted from the corpus using a finite-state grammar. 13,426 conjunctions were expanded to a total of 15,431 conjoined pairs, reduced to 15,048 tokens of 9,296 distinct adjective pairs after the morphological transformation of words such as “taller” and “tallest” to “tall”.

In the second phase, the system uses a log-linear regression model to combine information from different conjunctions to determine if the adjectives are of similar or different orientation. “77.84% of all links from conjunctions indicate same orientation,” [1] so the algorithm can achieve fair performance simply by always guessing that the adjectives within the conjunction have equivalent orientation. The performance is improved by the fact that “but” is primarily exhibited in conjunctions of opposite orientation.

In the third phase, the system places the adjectives into groups. A clustering algorithm separates the adjectives using an “associated dissimilarity value.” Dissimilar adjectives are given values that approach one and similar adjectives are given values that approach zero. The adjectives connected by “and” have low dissimilarities and the ones connected by “but” have high dissimilarities. The system uses the dissimilarity information to split the adjectives into two subsets of differing orientation and places as many adjectives as it can into the same subset.

Finally, in the fourth phase, the system assigns polarity to the two sets of unmarked adjectives. We know that “in oppositions of gradable adjectives where one member is semantically unmarked, the unmarked member is the most frequent one about 81% of the time.” [1] Because of this, we can compare the average frequencies of each semantically unmarked group of words and safely assume that the one with the highest frequency is the set of positive words. Hatzivassiloglou argues that this practice increases labeling precision even when some of the words are incorrectly assigned [1].

Results

Using conjunctions proves effective, “achieving 82% accuracy in this task when each conjunction is considered independently.” Through combining the conjunction-imposed constraints over multiple pairs of adjectives, the accuracy was further increased. When the average number of links for each adjective was 10.49, the accuracy of the system was 92.37% [1]. Hatzivassiloglou notes that graph connectivity, the average number of conjunctive links, is largely a function of corpus size, and can be increased by using larger corpora. Accuracy was as low as 78% on the “sparsest” set but was the aforementioned 92% accuracy on the densest set.

4. NATURAL LANGUAGE PROCESSING COMBINED METHOD

The natural language processing approach is predicated on the idea that the analysis of local statements is more reliable than attempts to judge overall opinion. An example is the sentence “Product A is good but expensive.” The sentence contains two different statements, “Product A is good,” and “Product A is expensive,” one of which constitutes a favorable statement and the other constitutes an unfavorable statement [4]. This is extremely reminiscent of Turney’s call for future work able to differentiate between “parts” and “wholes,” a problem that lead to decreased SO-PMI and SO-LSA accuracy.

Natural language processing attempts to identify and judge each individual statement rather than attempting to sum the whole. It does so by annotating each word or phrase in the text with its part of speech. Parts of speech whose meanings apply to each other are separated into individual statements. The polarity of these statements is then, independently, determined. Nasukawa outlines a sentiment notation for this purpose: On issues of polarity, positive polarity is denoted by the letter “g,” negative by the letter “b,” and neutral by the letter “n.” Verbs that transfer sentiment are denoted by “t.” Parts of speech are signified by letter pairs, with adjectives denoted by (JJ), adverbs by (RB), nouns by (NN), and verbs by (VB). The notation also includes the sentiment term in canonical form and any subjects (sub) and objects (obj) that receive sentiment from other arguments. This allows statements to be easily and understandably encoded.

Nasukawa provides examples of how this encoding is useful. The notation “gVB admire obj” indicates that “admire” indicates favorability towards the noun phrase in its object when the object contains a subject term. This is important because the statement that “XXXX admires YYYY” indicates favorability to YYYY rather than XXXX. Conversely, “bVB fail sub” shows a verb “fail” that indicates unfavorability towards the noun phrase in its subject that contains a subject term. The fact that XXXX fails to do YYYY doesn’t reflect poorly on YYYY. Finally, “tVB provide obj sub” shows that the verb “provide” passes the favorability or unfavorability of its object onto its target, assuming its object noun phrase contains favorability information and that the target term is in its subject [4]. The relationship implied in “XXX provides a good working environment” passes the favorable association of “good working environment” to the subject XXX.

The sentiment analysis algorithm searches for a subject term and considers the sentence surrounding the term as well as the subsequent parts of the paragraph in which the term takes place. This window includes a minimum of five words before and five after and a maximum of 50 words before and 50 after. A Markov-model-based tagger is used for part-of-speech tagging [4]. Once the relevant statements in the text are fully tagged, sentiment polarity is attributed according to a premade sentiment dictionary. When negative expressions such as “not” and “never” are found, the opposite sentiment polarity is attributed [4].

Results

The results are impressive but mixed. Without any modification of the dictionary, the Natural Language Processing Combined Method system achieved 94.5% (=241/255) accuracy with about 24% (=241/1,000) recall. Recall represents the Natural Language Processing system's ability to recognize statements relevant to the query. When, out of 1,000 total relevant statements, 241 are recognized and returned, the system has a recall of 24%. Accuracy drops to about 75% in some cases involving "well-written texts" such as organizational web pages and news articles. It's theorized that the accuracy hit is due to the long and complex sentences in those kinds of documents. Other failures are caused by larger context. For example "It's difficult to take a bad picture with this camera." is scored "-1 bad-picture (a bad picture)." While it's true that bad referred to picture, the statement referred to the camera positively.

5. CONCLUSIONS

In conclusion, more is better. The more a method "understands" the text, the more effective it is, and the ability to understand context is of vital importance to sentiment analysis efforts. Of the four systems, the Natural Language Processing Combined Method was the most effective, due to its ability to render sentences into their component parts of speech and sensibly apply sentiment over subject-object relationships. This statement-level awareness provides the Natural Language Processing Combined Method with an accuracy advantage over SO-PMI alone.

SO-PMI is the most accurate of the other methods, but only over large corpora, such as the 100 billion word AltaVista English corpus. Over smaller corpora, its accuracy and stability quickly drop. While SO-LSA more efficiently uses smaller corpora, it cannot yet utilize large corpora and shares SO-PMI's general problems. While using a different form of semantic association analysis, it too sums the associations of every weighted word throughout the passage to find the total weight. This way offers less insight than part of speech-aware natural language processing, and is less useful as a result.

Finally, adjective conjunction combination, despite making use of some contextual information, is the least effective of all of the above. While its ability to increase accuracy by weighting unweighted adjectives based on frequency is helpful, its inability to interpret the semantic orientations of verbs and nouns, such as SO-PMI and SO-LSA, or draw on the larger relationships of parts of speech within the statement, such as the NLP combined method, greatly hampers its effectiveness.

6. REFERENCES

- [1] V. Hatzivassiloglou and K. R. McKeown. Predicting the semantic orientation of adjectives. In *EACL '97: Proceedings of the eighth conference on European chapter of the Association for Computational Linguistics*, pages 174–181, Morristown, NJ, USA, 1997. Association for Computational Linguistics.
- [2] D. Lee, O.-R. Jeong, and S.-g. Lee. Opinion mining of customer feedback data on the web. In *ICUIMC '08: Proceedings of the 2nd international conference on Ubiquitous information management and*

communication, pages 230–235, New York, NY, USA, 2008. ACM.

- [3] M. Louwerse, X. Hu, Z. Cai, M. Ventura, and P. Jeuniaux. The embodiment of amodal symbolic knowledge representations. *Proceedings of the 18th International Florida Artificial Intelligence Research Society*, pages 542–547, 2005.
- [4] T. Nasukawa and J. Yi. Sentiment analysis: capturing favorability using natural language processing. In *K-CAP '03: Proceedings of the 2nd international conference on Knowledge capture*, pages 70–77, New York, NY, USA, 2003. ACM.
- [5] P. Turney and M. Littman. Unsupervised learning of semantic orientation from a hundred-billion-word corpus. National research council technical report erb-1094, Institute for Information Technology, National Research Council Canada, 2002.
- [6] P. D. Turney. Thumbs up or thumbs down?: semantic orientation applied to unsupervised classification of reviews. In *ACL '02: Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, pages 417–424, Morristown, NJ, USA, 2002. Association for Computational Linguistics.
- [7] P. D. Turney and M. L. Littman. Measuring praise and criticism: Inference of semantic orientation from association. *ACM Trans. Inf. Syst.*, 21(4):315–346, 2003.
- [8] H. Widdowson. J.R. Firth, 1957, papers in linguistics 1934-51. *International Journal of Applied Linguistics*, 17(3):402–413, 2007.