

Determining Term Subjectivity and Term Orientation for Opinion Mining

Andrea Esuli

Istituto di Scienza e Tecnologie dell'Informazione
Consiglio Nazionale delle Ricerche
Via G Moruzzi, 1 – 56124 Pisa, Italy
andrea.esuli@isti.cnr.it

Fabrizio Sebastiani

Dipartimento di Matematica Pura e Applicata
Università di Padova
Via GB Belzoni, 7 – 35131 Padova, Italy
fabrizio.sebastiani@unipd.it

ABSTRACT

Opinion mining is a recent subdiscipline of information retrieval which is concerned not with the topic a document is about, but with the opinion it expresses. To aid the extraction of opinions from text, recent work has tackled the issue of determining the *orientation* of “subjective” terms contained in text, i.e. deciding whether a term that carries opinionated content has a positive or a negative connotation; this is believed to be of key importance for identifying the orientation of documents, i.e. determining whether a document expresses a positive or negative opinion about its subject matter. We contend that the plain determination of the orientation of terms is not a realistic problem, since it starts from the non-realistic assumption that we already know whether a term is subjective or not; this would imply that a linguistic resource that marks terms as “subjective” or “objective” is available, which is usually not the case. In this paper we confront the task of deciding whether a given term has a positive connotation, or a negative connotation, or has no subjective connotation at all; this problem thus subsumes the problem of determining subjectivity/objectivity *and* the problem of determining orientation. We tackle this problem by testing three different variants of the semi-supervised method for orientation detection presented in [2]. Our results show that determining subjectivity *and* orientation is a much harder problem than determining orientation alone.

Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Information filtering; Search process*;
H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing—*Linguistic processing*; I.2.7 [Artificial Intelligence]: Natural Language Processing—*Text analysis*; I.5.2 [Pattern Recognition]: Design Methodology—*Classifier design and evaluation*

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
Copyright 200X ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

General Terms

Algorithms, Experimentation

Keywords

Text Classification, Opinion Mining, Sentiment Classification, Semantic Orientation, Polarity Detection, Subjectivity Detection

1. INTRODUCTION

Opinion mining is a recent subdiscipline of information retrieval which is concerned not with the topic a document is about, but with the opinion it expresses. Opinion-driven content management has several important applications, such as determining critics’ opinions about a given product by classifying online product reviews, or tracking the shifting attitudes of the general public towards a political candidate by mining online forums.

Within opinion mining, several subtasks can be identified, all of them having to do with tagging a given document according to expressed opinion:

1. *determining document subjectivity*, as in deciding whether a given text has a factual nature (i.e. describes a given situation or event, without expressing a positive or a negative opinion on it) or expresses an opinion on its subject matter. This amounts to performing binary text categorization under categories **Objective** and **Subjective** [10, 17];
2. *determining document orientation (or polarity)*, as in deciding if a given **Subjective** text expresses a **Positive** or a **Negative** opinion on its subject matter [10, 14];
3. *determining the strength of document orientation*, as in deciding e.g. whether the **Positive** opinion expressed by a text on its subject matter is **Weakly Positive**, **Mildly Positive**, or **Strongly Positive** [11, 16].

To aid these tasks, recent work [2, 3, 7, 8, 15] has tackled the issue of identifying the orientation of subjective *terms* contained in text, i.e. determining whether a term that carries opinionated content has a positive or a negative connotation (e.g. deciding that — using Turney and Littman’s [15] examples — **honest** and **intrepid** have a positive connotation while **disturbing** and **superfluous** have a negative connotation). This is believed to be of key importance for

identifying the orientation of documents, since it is by considering the combined contribution of these terms that one may hope to solve Tasks 1, 2 and 3 above. The conceptually simplest approach to this latter problem is probably Turney’s [14], who has obtained interesting results on Task 2 by considering the algebraic sum of the orientations of terms as representative of the orientation of the document they belong to; but more sophisticated approaches are also possible [4, 12, 16].

Implicit in most works dealing with term orientation is the assumption that, for many languages for which one would like to perform opinion mining, there is no available lexical resource where terms are tagged as having either a **Positive** or a **Negative** connotation, and that in the absence of such a resource the only available route is to generate such a resource automatically.

However, we think this approach lacks realism since it is also true that, for the very same languages, there is no available lexical resource where terms are tagged as having either a **Subjective** or an **Objective** connotation. Thus, the availability of an algorithm that tags **Subjective** terms as being either **Positive** or **Negative** is of little help, since determining if a term is **Subjective** is itself non-trivial.

In this paper we confront the task of determining whether a given term has a **Positive** connotation (e.g. **honest, intrepid**), or a **Negative** connotation (e.g. **disturbing, superfluous**), or has instead no **Subjective** connotation at all (e.g. **white, triangular**); this problem thus subsumes the problem of deciding between **Subjective** and **Objective** *and* the problem of deciding between **Positive** and **Negative**. We tackle this problem by testing three different variants of the semi-supervised method for orientation detection proposed in [2]. Our results show that determining subjectivity *and* orientation is a much harder problem than determining orientation alone.

1.1 Outline of the paper

The rest of the paper is structured as follows. Section 2 reviews related work on opinion mining, with special emphasis on works dealing with term orientation and/or subjectivity detection. Section 3 briefly reviews the semi-supervised method for orientation detection presented in [2]. Section 4 describes in detail three different variants of it we propose for determining, at the same time, subjectivity *and* orientation, and describes the general setup of our experiments. In Section 5 we discuss the results we have obtained. Section 6 concludes.

2. RELATED WORK

2.1 Determining term orientation

Previously appeared papers dealing with the properties of terms within an opinion mining perspective have concentrated, as mentioned above, on determining term orientation.

Hatzivassiloglou and McKeown [3] attempt to predict the orientation of subjective *adjectives* by analysing pairs of adjectives (conjoined by **and**, **or**, **but**, **either-or**, or **neither-nor**) extracted from a large unlabelled document set. The underlying intuition is that the act of conjoining adjectives is subject to linguistic constraints on the orientation of the adjectives involved; e.g. **and** usually conjoins adjectives of equal orientation, while **but** conjoins adjectives of opposite orientation. The authors generate a graph where terms are nodes

connected by “equal-orientation” or “opposite-orientation” edges, depending on the conjunctions extracted from the document set. A clustering algorithm then partitions the graph into a **Positive** cluster and a **Negative** cluster, based on a relation of similarity induced by the edges.

Turney and Littman [15] determine term orientation by bootstrapping from two small sets of subjective “seed” terms (with the seed set for **Positive** containing terms such as **good** and **nice**, and the seed set for **Negative** containing terms such as **bad** and **nasty**). Their method is based on computing the *pointwise mutual information* (PMI) of the target term t with each seed term t_i as a measure of their semantic association. Given a target term t , its orientation value $O(t)$ (where positive value means positive orientation, and higher absolute value means stronger orientation) is given by the sum of the weights of its semantic association with the seed positive terms minus the sum of the weights of its semantic association with the seed negative terms. For computing PMI, term frequencies and co-occurrence frequencies are measured by querying a document set by means of the AltaVista search engine¹ with a “ t ” query, a “ t_i ” query, and a “ t NEAR t_i ” query, and using the number of matching documents returned by the search engine as estimates of the probabilities needed for the computation of PMI.

Kamps et al. [7] consider instead the subgraph defined on adjectives by the WordNet² synonymy relation, and determine the orientation of a target adjective t contained in the subgraph by comparing the lengths of (i) the shortest path between t and the seed term **good**, and (ii) the shortest path between t and the seed term **bad**: if the former is shorter than the latter, than t is deemed to be **Positive**, otherwise it is deemed to be **Negative**.

The system of Kim and Hovy [8] tackles orientation detection by attributing, to each term, a positivity score *and* a negativity score; interestingly, terms may thus be deemed to have both a positive and a negative correlation, maybe with different degrees, and some terms may be deemed to carry a stronger positive (or negative) orientation than others. Their system starts from a set of positive and negative seed terms, and expands the positive (resp. negative) seed set by adding to it the synonyms of positive (resp. negative) seed terms and the antonyms of negative (resp. positive) seed terms. The system classifies then a target term t into either **Positive** or **Negative** by means of two alternative learning-free methods based on the probabilities that synonyms of t also appear in the respective expanded seed sets. A problem with this method is that it can classify only terms that share some synonyms with the expanded seed sets.

The approach we have proposed for determining term orientation [2] is described in more detail in Section 3, since it will be extensively used in this paper.

All these works evaluate the performance of the proposed algorithms by checking them against precompiled sets of **Positive** and **Negative** terms, i.e. checking how good the algorithms are at classifying a term known to be subjective into either **Positive** or **Negative**. When tested on the same benchmarks, the methods of [2, 15] have performed with comparable accuracies (however, the method of [2] is much more efficient than the one of [15]), and have outperformed

¹<http://www.altavista.com/>

²<http://wordnet.princeton.edu/>

the method of [3] by a wide margin and the one by [7] by a very wide margin. The methods described in [3, 7] are also limited by the fact that they can only decide the orientation of *adjectives*; the method of [7] is further limited in that it can only work on adjectives that are present in WordNet. The method of [8] is instead difficult to compare with the other ones mentioned so far since it was not evaluated on publicly available datasets.

2.2 Determining term subjectivity

To the best of our knowledge, the work of Riloff et al. [12] is the only one so far that has attempted to determine whether a term has a Subjective or Objective connotation. Their method identifies patterns for the extraction of subjective *nouns* from text, bootstrapping from a seed set of 20 terms that the authors judge to be strongly subjective and have found to have high frequency in the text collection from which the subjective nouns must be extracted.

The results of this method are not easy to compare with the ones we present in this paper because of the different evaluation methodologies. While we adopt the evaluation methodology used in all of the papers reviewed so far (i.e. checking how good our system is at replicating an existing, independently motivated lexical resource), the authors of [12] test their method on the set of terms that the algorithm itself extracts. This evaluation methodology only allows to test precision, and not accuracy, since no quantification can be made of false negatives (i.e. the subjective terms that the algorithm should have spotted but has not spotted). In Section 5 this will prevent us from drawing comparisons between this method and our own.

3. DETERMINING TERM SUBJECTIVITY AND TERM ORIENTATION BY SEMI-SUPERVISED LEARNING

The method we use in this paper for determining term subjectivity and term orientation is a variant of the method proposed in [2] for determining term orientation alone.

This latter method relies on training, in a semi-supervised way, a binary classifier that labels terms as being either Positive or Negative. A *semi-supervised* method is a learning process whereby only a small subset $L \subset Tr$ of the training data Tr are human-labelled. In origin the training data in $U = Tr - L$ are instead unlabelled; it is the process itself that labels them, automatically, by using L (with the possible addition of other publicly available resources) as input. The method of [2] starts from two small seed (i.e. training) sets L_p and L_n of known Positive and Negative terms, respectively, and expands them into the two final training sets $Tr_p \supset L_p$ and $Tr_n \supset L_n$ by adding them new sets of terms U_p and U_n found by navigating the WordNet graph along the synonymy and antonymy relations³. This process is based on the hypothesis that synonymy and antonymy, in addition to defining a relation of meaning, also define a relation of orientation, i.e. that

³Several other WordNet lexical relations, and several combinations of them, are tested in [2]. In the present paper we only use synonymy (e.g. *use / utilize*) and (direct) antonymy (e.g. *light / dark*), with the restriction that the two related terms must have the same POS, since this is the combination that has been shown to perform best in [2]. The version of WordNet used here and in [2] is 2.0.

two synonyms typically have the same orientation and two antonyms typically have opposite orientation⁴. The method is iterative, generating two sets Tr_p^k and Tr_n^k at each iteration k , where $Tr_p^k \supset Tr_p^{k-1} \supset \dots \supset Tr_p^1 = L_p$ and $Tr_n^k \supset Tr_n^{k-1} \supset \dots \supset Tr_n^1 = L_n$. At iteration k , Tr_p^k is obtained by adding to Tr_p^{k-1} all synonyms of terms in Tr_p^{k-1} and all antonyms of terms in Tr_n^{k-1} ; similarly, Tr_n^k is obtained by adding to Tr_n^{k-1} all synonyms of terms in Tr_n^{k-1} and all antonyms of terms in Tr_p^{k-1} . If a total of K iterations are performed, then $Tr = Tr_p^K \cup Tr_n^K$.

The second main feature of the method presented in [2] is that terms are given vectorial representations based on their WordNet *glosses* (i.e. textual definitions). For each term t_i in $Tr \cup Te$ (Te being the test set, i.e. the set of terms to be classified), a textual representation of t_i is generated by collating all the glosses of t_i as found in WordNet⁵. Each such representation is converted into vectorial form by standard text indexing techniques (in [2] and in the present work, stop words are removed and the resulting features are weighted by cosine-normalized *tfidf*; no stemming is performed)⁶. This representation method is based on the assumption that terms with a similar orientation tend to have “similar” glosses: for instance, that the glosses of **honest** and **intrepid** will both contain appreciative expressions, while the glosses of **disturbing** and **superfluous** will both contain derogative expressions. Note that this method allows *any* term to be classified, provided there is a gloss for it in the lexical resource.

4. EXPERIMENTS

In this paper we extend the method of [2] to the determination of term subjectivity *and* term orientation.

4.1 Test sets

The benchmark (i.e. test set) we use for our experiments is the General Inquirer (GI) lexicon [13] This is a lexicon of terms labelled according to a large set of categories⁷, each one denoting the presence of a specific trait in the term. The two main categories, and the ones we will be concerned with, are Positive/Negative, which contain 1,915/2,291 terms having a positive/negative orientation (in what follows we will also refer to the category Subjective, defined as the union of the two categories Positive and Negative). Examples of Positive terms are *advantage*, *fidelity* and *worthy*, while examples of negative terms are *badly*, *cancer*, *stagnant*. In opinion mining research the GI was first used by Turney and Littman [15], who reduced the list of terms to 1,614/1,982

⁴This intuition is basically the same as that of Kim and Hovy [8], whose work we did not know at the time of writing [2].

⁵In general a term t_i may have more than one gloss, since it may have more than one sense; dictionaries normally associate one gloss to each sense.

⁶Several combinations of subparts of a WordNet gloss are tested as term representations in [2]. Of all those combinations, in the present paper we always use the DGS \neg combination, since this is the one that has been shown to perform best in [2]. DGS \neg corresponds to using the entire gloss and performing *negation propagation* on its text, i.e. replacing all the terms that occur after a negation in a sentence with negated versions of the term (see [2] for details).

⁷The definitions of all such categories are available at <http://www.webuse.umd.edu:9090/>

entries after removing 17 terms appearing in both categories (e.g. *deal*) and reducing all the multiple entries of the same term in a category, caused by multiple senses, to a single entry. Likewise, we take all the 7,582 General Inquirer terms that are not labelled as either **Positive** or **Negative**, as being (implicitly) labelled as **Objective**, and reduce them to 5,009 terms after combining multiple entries of the same term, caused by multiple senses, to a single entry.

The effectiveness of our classifiers will thus be evaluated in terms of their ability to assign the total 8,605 GI terms to the correct category among **Positive**, **Negative**, and **Objective**⁸.

4.2 Seed sets and training sets

Similarly to [2], our training set is obtained by expanding initial seed sets by means of WordNet lexical relations, with the difference that, unlike in [2], our training set is now the union of *three* sets of training terms $Tr = Tr_p^K \cup Tr_n^K \cup Tr_o^K$ obtained by expanding, through K iterations, three seed sets Tr_p^1, Tr_n^1, Tr_o^1 , one for each of the categories **Positive**, **Negative**, and **Objective**, respectively.

Concerning categories **Positive** and **Negative**, we have used the seed sets, expansion policy, and number of iterations, that have performed best in the experiments of [2], i.e. the seed sets $Tr_p^1 = \{\text{good}\}$ and $Tr_n^1 = \{\text{bad}\}$ expanded by using the union of synonymy and direct antonymy, restricting the relations only to terms with the same POS of the original terms (i.e. adjectives), for a total of $K = 4$ iterations. The final expanded sets contain 6,053 **Positive** terms and 6,874 **Negative** terms.

Concerning the category **Objective**, the process we have followed is similar, with a few differences. These are motivated by the fact that **Objective** terms are more varied and diverse in meaning than the terms in the other categories. To obtain a representative expanded set Tr_o^K , we have chosen the seed set $Tr_o^1 = \{\text{entity}\}$ and we have expanded it by using, along with synonymy and antonymy, the WordNet relation of hyponymy (e.g. *vehicle* / *car*), and without imposing the restriction that the two related terms must have the same POS. These choices are strictly related to each other: the term **entity** is the root term of the largest generalisation hierarchy in WordNet (with more than 40,000 terms) [1], thus allowing to reach a very large number of terms by using the hyponymy relation⁹. Moreover, it seems reasonable to assume that terms that refer to *entities* are likely to have an “objective” nature, and that hyponyms (and also synonyms and antonyms) of an objective term are also objective. Note that, at each iteration k , before adding a given term t to Tr_o^k we check if it already belongs to either Tr_p^k or Tr_n^k ; if it does, the term is not added to Tr_o^k and is discarded from consideration (i.e. is expanded no further). We experiment with two different choices for the Tr_o set, corresponding to the sets generated in $K = 3$ and $K = 4$ iterations, respectively; this yields sets Tr_o^3 and Tr_o^4 consisting of 8,353 and 33,870 training terms, respectively.

4.3 Learning approaches and evaluation measures

We experiment with three “philosophically” different learning approaches to the problem of distinguishing between

⁸This labelled term set can be downloaded from <http://patty.isti.cnr.it/~esuli/software/SentimentGI.tgz>

⁹The largest connected component for the synonymy relation consists instead of only 10,922 names [7].

Positive, **Negative**, and **Objective** terms.

Approach I is a two-stage method which consists in learning two binary classifiers: the first classifier places terms into either **Subjective** or **Objective**, while the second classifier places terms that have been classified as **Subjective** into either **Positive** or **Negative**. In the training phase, the terms in $Tr_p^K \cup Tr_n^K$ are used as training examples of category **Subjective**.

Approach II is again based on learning two binary classifiers. One of them must discriminate between terms that belong to the **Positive** category and ones that belong to its complement (**not Positive**), while the other must discriminate between terms that belong to the **Negative** category and ones that belong to its complement (**not Negative**). Terms that have been classified *both* into **Positive** by the former classifier and into (**not Negative**) by the latter are deemed to be positive, and terms that have been classified *both* into (**not Positive**) by the former classifier and into **Negative** by the latter are deemed to be negative. The terms that have been classified (i) into both (**not Positive**) and (**not Negative**), or (ii) into both **Positive** and **Negative**, are taken to be **Objective**. In the training phase of Approach II, the terms in $Tr_n^K \cup Tr_o^K$ are used as training examples of category (**not Positive**), and the terms in $Tr_p^K \cup Tr_o^K$ are used as training examples of category (**not Negative**).

Approach III consists in viewing **Positive**, **Negative**, and **Objective** as three categories with equal status, and in learning a ternary classifier that classifies each term into exactly one among the three categories.

There are several differences among these three approaches. A first difference, of a conceptual nature, is that only Approaches I and III view **Objective** as a category, or concept, in its own right, while Approach II views objectivity as a nonexistent entity, i.e. as the “absence of subjectivity” (in fact, in Approach II the training examples of **Objective** are only used as training examples of the *complements* of **Positive** and **Negative**). A second difference is that Approaches I and II are based on standard binary classification technology, while Approach III requires “multiclass” (i.e. 1-of- m) classification. As a consequence, while for the former we use well-known learners for binary classification (the naive Bayesian learner using the multinomial model (NB) [9], support vector machines using linear kernels [6], the Rocchio learner, and its PrTFIDF probabilistic version [5]), for Approach III we use their multiclass versions¹⁰.

Before running our learners we run a pass of feature selection, with the intent of retaining only those features that are good at discriminating our categories, while discarding those which are not. Feature selection is implemented by scoring each feature f_k (i.e. each term that occurs in the glosses of at least one training term) by means of the *mutual information* (MI) function, defined as

$$MI(f_k) = \sum_{c \in \{c_1, \dots, c_m\}} \sum_{f \in \{f_k, \bar{f}_k\}} \Pr(f, c) \cdot \log \frac{\Pr(f, c)}{\Pr(f) \Pr(c)}$$

and discarding the $x\%$ features f_k that minimize it. We will

¹⁰The naive Bayesian, Rocchio, and PrTFIDF learners we have used are from McCallum’s *Bow* package (<http://www-2.cs.cmu.edu/~mccallum/bow/>), while the SVMs learner we have used is version 6.01 of Joachims’ *SVMlight* (<http://svmlight.joachims.org/>). Both packages allow the respective learners to be run in “multiclass” fashion.

call $x\%$ the *reduction factor*. Note that the set $\{c_1, \dots, c_m\}$ from Equation 1 is interpreted differently in Approaches I to III, and always consistently with which the categories at stake are.

Since the task we aim to solve is manifold, we will evaluate our classifiers according to two evaluation measures:

- *SO-accuracy*, i.e. the accuracy of a classifier in separating **Subjective** from **Objective**, i.e. in deciding term subjectivity alone;
- *PNO-accuracy*, the accuracy of a classifier in discriminating among **Positive**, **Negative**, and **Objective**, i.e. in deciding both term orientation and subjectivity.

5. RESULTS

We present results obtained from running every combination of (i) the three approaches to classification described in Section 4.3, (ii) the four learners mentioned in the same section, (iii) five different reduction factors for feature selection (0%, 50%, 90%, 95%, 99%), and (iv) the two different training sets (Tr_o^3 and Tr_o^4) for **Objective** mentioned in Section 4.2. We discuss each of these four dimensions of the problem individually, for each one reporting results averaged across all the experiments we have run (see Table 1).

The first and most important observation is that, with respect to a pure term orientation task, accuracy drops significantly. In fact, the best *SO-accuracy* and the best *PNO-accuracy* results obtained across the 120 different experiments are .676 and .660, respectively (these were obtained by using Approach II with the PrTFIDF learner, and no feature selection (with $Tr_o = Tr_o^3$) for the .676 *SO-accuracy* result, with $Tr_o = Tr_o^4$ for the .660 *PNO-accuracy* result); this contrasts sharply with the accuracy obtained in [2] on discriminating **Positive** from **Negative** (where the best run obtained .830 accuracy), *with the same benchmarks and pretty much the same algorithms*. This suggests that good performance at orientation detection (as e.g. in [2, 3, 15]) may not be a guarantee of good performance at subjectivity detection, quite evidently a harder (and, as we have suggested, more realistic) task.

The second important observation is that there is very little variance in the results: across all 120 experiments, average *SO-accuracy* and *PNO-accuracy* results were .635 (with standard deviation $\sigma = .030$) and .603 ($\sigma = .036$), a mere 6.06% and 8.64% deterioration from the best results reported above. This seems to indicate that the levels of performance obtained may be hard to improve upon, especially if working in a similar framework.

Let’s analyse the individual dimensions of the problem. Concerning the three approaches to classification described in Section 4.3, Approach II outperforms the other two, but by an extremely narrow margin. As for the choice of learners, on average the best performer is NB, but again by a very small margin wrt the others. On average, the best reduction factor for feature selection turns out to be 50%, but the performance drop we witness in approaching 99% (a dramatic reduction factor) is extremely graceful. As for the choice of Tr_o^K , we note that Tr_o^3 and Tr_o^4 elicit comparable levels of performance, with the former performing best at *SO-accuracy* and the latter performing best at *PNO-accuracy*.

An interesting observation on the learners we have used is that NB, PrTFIDF and SVMs, unlike Rocchio, generate

classifiers that depend on $P(c_i)$, the prior probabilities of the classes, which are normally estimated as the proportion of training documents that belong to c_i . In many classification applications this is reasonable, as we may assume that the training data are sampled from the same distribution from which the test data are sampled, and that these proportions are thus indicative of the proportions that we are going to encounter in the test data. However, in our application this is not the case, since we do not have a “natural” sample of training terms. What we have is one human-labelled training term for each category in $\{\text{Positive, Negative, Objective}\}$, and as many machine-labelled terms as we deem reasonable to include, in possibly different numbers for the different categories; and we have no indication whatsoever as to what the “natural” proportions among the three might be. This means that the proportions of **Positive**, **Negative**, and **Objective** terms we decide to include in the training set will strongly bias the classification results if the learner is one of NB, PrTFIDF and SVMs. We may notice this by looking at Table 2, which shows the average proportion of test terms classified as **Objective** by each learner, depending on whether we have chosen Tr_o to coincide with Tr_o^3 or Tr_o^4 ; note that the former (resp. latter) choice means having roughly as many (resp. roughly five times as many) **Objective** training terms as there are **Positive** and **Negative** ones. Table 2 shows that, the more **Objective** training terms there are, the more test terms NB, PrTFIDF and (in particular) SVMs will classify as **Objective**; this is not true for Rocchio, which is basically unaffected by the variation in size of Tr_o .

Finally, we have measured the time required by each experiments, consisting in performing feature selection on the term glosses, training the classifier, and applying it to the test set. These timings do not include the time required to compute the expanded Tr_x^K sets (2 minutes for $K = 3$, 6 minutes for $K = 4$), and to extract and process WordNet glosses of training and test terms (32 minutes for the training set that uses Tr_o^3 , 72 minutes for the one that uses Tr_o^4 , and 14 minutes for the test set). The average times for the various learners and approaches are shown in Table 3. The first observation is that SVMs are more computationally demanding than the other learners, which require no more than two minutes for a complete run. The multiclass version of SVMs, used in Approach III, is the least efficient learner, ten times less efficient (on average) than the single-class version used in Approach I. The other three learners have a similar behaviour, that is mostly driven by the amount of data to be processed during the learning phase. In Approach III the training data are processed only once, which means this is the fastest approach. In Approach I the training data are fully processed once, to learn the **Subjective** vs. **Objective** classifier, and then the **Positive** and **Negative** training examples are re-processed to learn the **Positive** vs. **Negative** classifier. In Approach II the training data are processed twice, in order to build the two binary classifiers for the **Positive** and **Negative** categories, respectively, which means this is the slowest approach.

6. CONCLUSIONS

We have presented a method for determining *both* term subjectivity *and* term orientation for opinion mining applications. This is a valuable advance with respect to the state of the art, since past work in this area had mostly confined to determining term orientation alone, a task that (as we have

Table 1: Average and best accuracy values over the four dimensions analysed in the experiments.

Dimension	SO-accuracy		PNO-accuracy	
	Avg (σ)	Best	Avg (σ)	Best
<i>Approach</i>				
I	.635 (.020)	.668	.595 (.029)	.635
II	.636 (.033)	.676	.614 (.037)	.660
III	.635 (.036)	.674	.600 (.039)	.648
<i>Learner</i>				
NB	.653 (.014)	.674	.619 (.022)	.647
SVMs	.627 (.033)	.671	.601 (.037)	.658
Rocchio	.624 (.030)	.654	.585 (.033)	.616
PrTFIDF	.637 (.031)	.676	.606 (.042)	.660
<i>TSR</i>				
0%	.649 (.025)	.676	.619 (.027)	.660
50%	.650 (.022)	.670	.622 (.022)	.657
80%	.646 (.023)	.674	.621 (.021)	.647
90%	.642 (.024)	.667	.616 (.024)	.651
95%	.635 (.027)	.671	.606 (.031)	.658
99%	.612 (.036)	.661	.570 (.049)	.647
<i>Tr_o^K set</i>				
Tr _o ³	.645 (.006)	.676	.608 (.007)	.658
Tr _o ⁴	.633 (.013)	.674	.610 (.018)	.660

Table 2: Average proportion of test terms classified as Objective, for each of the four learners and for each choice of the Tr_o^K set.

Learner	Tr _o ^K set		Variation
	Tr _o ³	Tr _o ⁴	
NB	.564 ($\sigma = .069$)	.693 (.069)	+23.0%
SVMs	.601 (.108)	.814 (.083)	+35.4%
Rocchio	.572 (.043)	.544 (.061)	-4.8%
PrTFIDF	.636 (.059)	.763 (.085)	+20.0%

argued) has limited practical significance in itself, given the generalized absence of lexical resources that tag terms as being either Subjective or Objective. Our algorithms have tagged by orientation *and* subjectivity the entire General Inquirer lexicon, a complete general-purpose lexicon that is the *de facto* standard benchmark for researchers in this field. Our results thus constitute, for this task, the first baseline for other researchers to improve upon.

Unfortunately, our results have shown that an algorithm that had shown excellent, state-of-the-art performance in deciding term orientation [2], once modified for the purposes of deciding term subjectivity, performs more poorly. This has been shown by testing several variants of the basic algorithm, some of them involving radically different supervised learning policies. The results suggest that deciding term subjectivity is a substantially harder task than deciding term orientation alone.

7. REFERENCES

- [1] A. Devitt and C. Vogel. The topology of WordNet: Some metrics. In *Proceedings of GWC-04, 2nd Global WordNet Conference*, pages 106–111, Brno, CZ, 2004.
- [2] A. Esuli and F. Sebastiani. Determining the semantic orientation of terms through gloss analysis. In *Proceedings*

Table 3: Average time (in seconds) required by an experiment (i.e. learning plus classification) for each of the four learners and the three Approaches used.

Learner	Approach		
	I	II	III
NB	91 ($\sigma = 31.8$)	119 (53.0)	75 (34.1)
SVMs	364 (188.0)	594.2 (301.2)	3749 (625.0)
Rocchio	89 (31.0)	111 (51.9)	69 (12.3)
PrTFIDF	86 (28.4)	114 (49.3)	74 (22.7)

of *CIKM-05, the ACM SIGIR Conference on Information and Knowledge Management*, Bremen, DE, 2005.

- [3] V. Hatzivassiloglou and K. R. McKeown. Predicting the semantic orientation of adjectives. In *Proceedings of ACL-97, 35th Annual Meeting of the Association for Computational Linguistics*, pages 174–181, Madrid, ES, 1997.
- [4] V. Hatzivassiloglou and J. M. Wiebe. Effects of adjective orientation and gradability on sentence subjectivity. In *Proceedings of COLING-00, 18th International Conference on Computational Linguistics*, pages 174–181, 2000.
- [5] T. Joachims. A probabilistic analysis of the Rocchio algorithm with TFIDF for text categorization. In *Proceedings of ICML-97, 14th International Conference on Machine Learning*, pages 143–151, Nashville, US, 1997.
- [6] T. Joachims. Text categorization with support vector machines: learning with many relevant features. In *Proceedings of ECML-98, 10th European Conference on Machine Learning*, pages 137–142, Chemnitz, DE, 1998.
- [7] J. Kamps, M. Marx, R. J. Mokken, and M. De Rijke. Using WordNet to measure semantic orientation of adjectives. In *Proceedings of LREC-04, 4th International Conference on Language Resources and Evaluation*, volume IV, pages 1115–1118, Lisbon, PT, 2004.
- [8] S.-M. Kim and E. Hovy. Determining the sentiment of opinions. In *Proceedings of COLING-04, 20th International Conference on Computational Linguistics*, pages 1367–1373, Geneva, CH, 2004.
- [9] A. K. McCallum and K. Nigam. A comparison of event models for naive Bayes text classification. In *Proceedings of the AAAI Workshop on Learning for Text Categorization*, pages 41–48, Madison, US, 1998.
- [10] B. Pang and L. Lee. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of ACL-04, 42nd Meeting of the Association for Computational Linguistics*, pages 271–278, Barcelona, ES, 2004.
- [11] B. Pang and L. Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of ACL-05, 43rd Meeting of the Association for Computational Linguistics*, Ann Arbor, US, 2005.
- [12] E. Riloff, J. Wiebe, and T. Wilson. Learning subjective nouns using extraction pattern bootstrapping. In *Proceedings of CONLL-03, 7th Conference on Natural Language Learning*, pages 25–32, Edmonton, CA, 2003.
- [13] P. J. Stone, D. C. Dunphy, M. S. Smith, and D. M. Ogilvie. *The General Inquirer: A Computer Approach to Content Analysis*. MIT Press, Cambridge, US, 1966.
- [14] P. Turney. Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews. In *Proceedings of ACL-02, 40th Annual Meeting of the Association for Computational Linguistics*, pages 417–424, 2002.
- [15] P. D. Turney and M. L. Littman. Measuring praise and criticism: Inference of semantic orientation from association. *ACM Transactions on Information Systems*,

21(4):315–346, 2003.

- [16] T. Wilson, J. Wiebe, and R. Hwa. Just how mad are you? Finding strong and weak opinion clauses. In *Proceedings of AAAI-04, 21st Conference of the American Association for Artificial Intelligence*, San Jose, US, 2004.
- [17] H. Yu and V. Hatzivassiloglou. Towards answering opinion questions: Separating facts from opinions and identifying the polarity of opinion sentences. In *Proceedings of EMNLP-03, 8th Conference on Empirical Methods in Natural Language Processing*, pages 129–136, 2003.