

# On Learning Parsimonious Models for Extracting Consumer Opinions

Xue Bai and Rema Padman

The John Heinz III School of Public Policy and Management  
and Center for Automated Learning and Discovery  
Carnegie Mellon University  
{xbai, rpadman}@andrew.cmu.edu

Edoardo Airoldi

Data Privacy Laboratory  
School of Computer Science  
Carnegie Mellon University  
eairoldi@cs.cmu.edu

## Abstract

*Extracting sentiments from unstructured text has emerged as an important problem in many disciplines. An accurate method would enable us, for example, to mine on-line opinions from the Internet and learn customers' preferences for economic or marketing research, or for leveraging a strategic advantage. In this paper, we propose a two-stage Bayesian algorithm that is able to capture the dependencies among words, and, at the same time, finds a vocabulary that is efficient for the purpose of extracting sentiments. Experimental results on the Movie Reviews data set show that our algorithm is able to select a parsimonious feature set with substantially fewer predictor variables than in the full data set and leads to better predictions about sentiment orientations than several state-of-the-art machine learning methods. Our findings suggest that sentiments are captured by conditional dependence relations among words, rather than by keywords or high-frequency words.*

## 1. Introduction

Traditionally, researchers have used surveys to collect a limited amount of data in a structured form for their analyses. In recent years, the advent of the Internet, and the widespread use of advanced information technologies in general, have resulted in a surge of information that is freely available on-line in an *unstructured format*. For example, many discussion groups and review sites exist where people post their opinions about a product. The automatic understanding of *sentiments* expressed within the texts of such posts could lead to a number of new applications in the fields of marketing and information retrieval.

Researchers have been investigating the problem of automatic text categorization for the past two decades. Satisfactory solutions have been found for the cases of topic categorization and of authorship attribution; briefly, topics are captured by sets of keywords [21], whereas authors are identi-

fied by their choices about the use of non-contextual, high-frequency words [22, 23, 1]. Pang et al [26] showed that such solutions, or extensions of them, yield cross-validated accuracies and areas under the curve (AUC) in the low 80% when ported to sentiment extraction. We conjecture that one reason for the failure of such approaches maybe attributed to the fact that the features used in the classification (e.g. the words) are assumed to be pairwise independent. The goal of this paper is to present a machine learning technique for learning predominant sentiments of on-line texts, available in unstructured format, that:

- is able to capture dependencies among words, and
- is able to find a minimal vocabulary, sufficient for categorization purposes.

Our two-stage Markov Blanket Classifier (MBC) learns conditional dependencies among the words and encodes them into a *Markov Blanket Directed Acyclic Graph* (MB DAG) for the sentiment variable (first stage), and then uses a *Tabu Search* (TS) meta-heuristic strategy to fine tune the MB DAG (second stage) in order to yield a higher cross-validated accuracy. Learning dependencies allows us to capture semantic relations and dependency patterns among the words, thus approximating the meaning of sentences, with important applications for many real world situations. Furthermore, performing the classification task using a Markov Blanket (MB) for the sentiment variable (in a Bayesian network) has important properties: (a) it specifies a statistically efficient prediction of the probability distribution of the sentiment variable from the smallest subset of predictors, and, (b) it provides accuracy while avoiding over-fitting due to redundant predictors. We test our algorithm on the publicly available Movie Reviews data set [20] and achieve a cross-validated accuracy of 87.5% and a cross-validated AUC of 96.85% respectively, against best performances of competing state-of-the-art classifiers in the low 80%.

This paper is organized as follows: Section 2 surveys

related work. Section 3 provides some background about Bayesian networks, Markov Blankets, and Tabu Search. Section 4 contains details about our proposed methodology. Section 5 describes the data and presents the experimental results. Finally, Section 6 discusses our findings and presents our conclusions.

## 2. Related Work on Sentiments

The problem of sentiment extraction is also referred to as opinion extraction or semantic classification in the literature. A related problem is that of studying the semantic orientation, or polarity, of words as defined by Osgood et al. [25]. Hatzivassiloglou and McKeown [13] built a log-linear model to predict the semantic orientation of conjoined adjectives using the conjunctions between them. Huettner and Subasic [14] hand-crafted a cognitive linguistic model for *affection* sentiments based on fuzzy logic. Das and Chen [7] used domain knowledge to manually construct lexicon and grammar rules that aim to capture the “pulse” of financial markets as expressed by on-line news about traded stocks. They categorized news as *buy*, *sell* or *neutral* using five classifiers and various voting schemes to achieve an accuracy of 62% (random guesses would top 33%). Turney and Littman [34] proposed a compelling semi-supervised method to learn the polarity of adjectives starting from a small set of adjectives of known polarity, and Turney [33] used this method to predict the opinions of consumers about various objects (movies, cars, banks) and achieved accuracies between 66% and 84%. Pang et al. [26] used off-the-shelf classification methods on frequent, non-contextual words in combination with various heuristics and annotators, and achieved a maximum cross-validated accuracy of 82.9% on data from IMDb [16]. Dave et al. [8] categorized positive versus negative movie reviews using support vector machines on various types of semantic features based on substitutions and proximity, and achieved an accuracy of at most 88.9% on data from Amazon and Cnn.Net. Last, Liu et al. [19] proposed a framework to categorize emotions based on a large dictionary of common sense knowledge and on linguistic models.

## 3. Theoretical Background

### 3.1. Bayesian Networks and Markov Blanket

A Bayesian network is a graphical representation of the joint probability distribution of a set of random variables as nodes in a graph, connected by directed edges. The orientations of the edges encapsulate the notion of parents, ancestors, children, and descendants of any node [27, 30].

More formally, a *Bayesian network* for a set of variables  $X = \{X_1, \dots, X_n\}$  consists of: (i) a network structure  $S$

that encodes a set of conditional independence assertions among variables in  $X$ ; and (ii) a set  $P = \{p_1, \dots, p_n\}$  of local conditional probability distributions associated with each node and its parents. Specifically,  $S$  is a directed acyclic graph (DAG) which, along with  $P$ , entails a joint probability distribution  $p$  over the nodes.

We say that  $P$  satisfies the *Markov condition* for  $S$  if every node  $X_i$  in  $S$  is independent of its non-descendants, conditional on its parents. The Markov Condition implies that the joint distribution  $p$  can be factorized as a product of conditional probabilities, by specifying the distribution of each node conditional on its parents. In particular, for given a structure  $S$ , the joint probability distribution for  $X$  can be written as

$$p(X) = \prod_{i=1}^n p_i(X_i | pa_i), \quad (1)$$

where  $pa_i$  denotes the set of parents of  $X_i$ .

Given the set of variables  $X$  and target variable  $Y$ , a *Markov Blanket* (MB) for  $Y$  is the smallest subset  $Q$  of variables in  $X$  such that  $Y$  is independent of  $X \setminus Q$ , conditional on the variables in  $Q$ . Intuitively, given a Bayesian network  $(S, P)$ , the Markov Blanket for  $Y$  consists of  $pa_Y$ , the set of parents of  $Y$ ;  $ch_Y$ , the set of children of  $Y$ ; and  $pa_{ch_Y}$ , the set of parents of children of  $Y$ .

For example, consider the two DAGs in Figure 1 and Figure 2, below. The factorization of  $p$  entailed by the Bayesian network  $(S, P)$  is

$$p(Y, X_1, \dots, X_6) = C \cdot p(Y|X_1) \cdot p(X_4|X_2, Y) \times \\ \times p(X_5|X_3, X_4, Y) p(X_2|X_1) \cdot p(X_3|X_1) \cdot p(X_6|X_4), \quad (2)$$

where  $C$  is a normalizing constant.

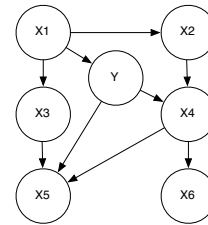
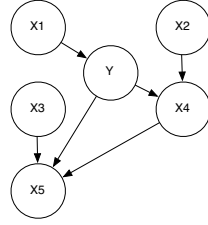


Figure 1. Bayesian network  $(S, P)$ .

The factorization of the conditional probability  $p(Y|X_1, \dots, X_6)$  entailed by the Markov blanket for  $Y$  corresponds to the product of those local factors in (2) which contain the term  $Y$ , that is

$$p(Y|X_1, \dots, X_6) = C' \cdot p(Y|X_1) \cdot p(X_4|X_2, Y) \times \\ \times p(X_5|X_3, X_4, Y), \quad (3)$$

where  $C'$  is a different normalizing constant.



**Figure 2. Markov Blanket for  $Y$  in  $(S, P)$ .**

Different MB DAGs that entail the same factorization for  $p(Y|X_1, \dots, X_6)$  belong to the same *Markov equivalence class*. Our algorithm searches the space of Markov equivalent classes, rather than that of DAGs, thus boosting its efficiency. Markov Blanket classifiers have been recently rediscovered and applied to several domains, but very few studies focus on how to learn the structure of the Markov Blanket from data. Further, the applications in the literature have been limited to data sets with few variables. Theoretically sound algorithms for finding DAGs are known (e.g. see [4]), but none have been tailored to the problem of finding MB DAGs.

### 3.2. Tabu Search

Tabu Search (TS) is a powerful meta-heuristic strategy that helps local search heuristics explore the space of solutions by guiding them out of local optima [11]. It has been applied successfully to a wide variety of continuous and discrete combinatorial optimization problems, and has been shown to be capable of reducing the complexity of the search process and accelerating the rate of convergence [12].

The basic Tabu Search starts with a feasible solution and iteratively chooses the *best move*, according to a specified evaluation function, while assuring that solutions previously generated are not revisited in the short-term. This is accomplished by keeping a *tabu list* of restrictions on possible moves, updated at each step, which discourage the repetition of selected moves. Typically tabu restrictions are based on a short-term memory function, called the *tabu tenure*, to prevent loops in the search, but intermediate and long-term memory functions may also be adopted to intensify and diversify the search.

We were motivated to use Tabu Search because its adaptive memory capability - both short term and long term - appears particularly suited to the Bayesian Networks and Markov Blanket approaches. Our choice of TS was also motivated by the extensively documented situations where its adaptive memory capability has proved successful, both directly and embedded within other "hybrid" methods such as those involving genetic algorithms, evolutionary compu-

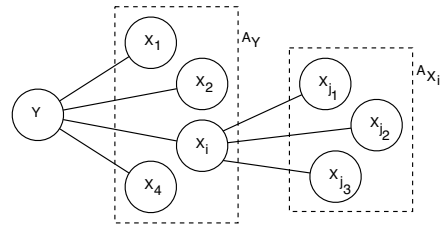
tation methods[11, 32] and scatter search [29]. Certainly it would be worthwhile to investigate additional metaheuristic procedures, and also to investigate more advanced forms of tabu search. Yet our results with the version of TS we propose have proved significantly better than those previously obtained by researchers in this area, and on this basis we conclude that our work provides a useful contribution even though we may subsequently find ways to improve upon it - as in fact we hope to do.

## 4. Methodology: Markov Blanket Classifier

### 4.1. 1<sup>st</sup> Stage: Learning Initial Dependencies

The first stage generates an initial MB DAG for  $Y$  from the data. This procedure involves the following: It begins by selecting those variables in  $\{X_1, \dots, X_N\}$  that are associated with  $Y$  within two hops in the graphical representation; that is, it finds potential parents and children ( $L_Y$ ) of  $Y$ , and potential parents and children ( $\cup_i L_{X_i}$ ) of nodes  $X_i \in L_Y$ , using conditional independence tests, representing adjacencies by undirected edges. At this point, the list  $Y \cup L_Y \cup \cup_i L_{X_i}$  is a skeleton (an undirected graph) which contains the MB for  $Y$  (See 3.1 for the precise definition of  $MB(Y)$  in terms of  $pa_Y$ ,  $ch_Y$ , and  $pa\ ch_Y$ .) The algorithm then **orients** the edges using six edge orientation rules described in Bai et al. [2]. Finally, it **prunes** the remaining undirected edges and bi-directed edges to avoid cycles, puts them in a list  $L$  for Tabu Search, and returns the MB DAG.

The core of the first stage lies in the search for the nodes ( $L_Y$ ) associated with  $Y$ , and for those ( $\cup_i L_{X_i}$ ) associated with the nodes in  $L_Y$ , based on causal discovery theory [27, 30]. This search is non trivial and is performed by two recursive calls to the function **findAdjacencies**( $Y$ ), as shown in Figure 3: independence tests between  $Y$  and each  $X_i$  are performed to identify a list ( $A_Y$ ) of variables associated with  $Y$ ; then, for  $X_i \in A_Y$  and for all distinct subsets  $S \subset \{A_Y \setminus X_i\}^d$ , where  $d$  controls the size of  $S$ , conditional independence tests between  $Y$  and  $X_i$  given  $S$



**Figure 3. Illustration of **findAdjacencies** ( $Y$ ).  $A_Y$  and  $A_{X_i}$  are shown.**

are performed to remove unfaithful associations; for more details about *unfaithful* associations and distributions see Spirtes et al. [30]. Then, for all pairs  $(X_i, X_j)_{i \neq j}$ , independence tests are performed to identify lists of variables  $(A_{X_i}, i=1, \dots, N)$  associated with each  $X_i$ ; last, for  $X_i \in A_Y$  and for all distinct subsets  $S \subset \{A_{X_i}\}^d$ , conditional independence tests between  $Y$  and each  $X_i$  given  $S$  are again performed to prune unfaithful associations.

#### 4.2. 2<sup>nd</sup> Stage: Tabu Search Optimization

Tabu Search (TS) is then applied to improve the initial MB DAG. Our algorithm searches for solutions in the space of logical Markov equivalence classes, instead of searching the space of MB DAGs; that is, moves that yield Markov Blankets within the same Markov equivalent class are not considered, and moves that result in cyclic graphs are not valid moves.

Briefly, four kinds of moves are allowed in the TS procedure: edge addition, edge deletion, edge reversal, and edge reversal with node pruning. At each stage, and for each allowed move, the corresponding MB DAG is computed, its conditional probability factored, its predictions scored, and the best move is then selected and applied. Best solution and best score at each step are tracked. The tabu list keeps a record of  $m$  previous moves, so that moves in the tabu list will not be repeated till their corresponding tabu tenure expires. Details can be found in [3].

#### 4.3. A Sketch of the Algorithm

We present a sketch of the algorithm below. The parameters are:  $D$ , a data set with  $N$  variables and  $K$  examples;  $Y$ , the class variable;  $d$ , the maximum number of nodes for the conditional independence tests;  $\alpha$ , the significance level for the  $G^2$  statistical independence tests (for a definition of  $G^2$  see [30]). The final output is the graphical Markov Blanket structure (MB) for  $Y$ .

**initialMBsearch** (Data  $D$ , Target  $Y$ , Depth  $d$ , Significance  $\alpha$ )

1.  $L_Y = \text{findAdjacencies}(Y, \{X_1, \dots, X_N\}, d, \alpha)$
2. **for**  $X_i \in L_Y$ 
  - 2.1.  $L_{X_i} = \text{findAdjacencies}(X_i, \{X_1, \dots, X_N\} \setminus X_i, d, \alpha)$
3.  $G = \text{orient}(Y \cup L_Y \cup L_{X_i})$
4.  $\{\text{MB DAG}, L\} = \text{prune}(G)$
5. **return**  $\{\text{MB DAG}, L\}$

**tabuSearch** (Data  $D$ , Target  $Y$ )

1. **init** ( $\text{bestSolution} = \text{currentSolution} = \text{MB DAG}$ ,  $\text{bestScore} = 0, \dots$ )
2. **repeat until** ( $\text{bestScore}$  does not improve for  $k$  consecutive iterations)

- 2.1. form  $\text{candidateMoves}$  for  $\text{currentSolution}$
- 2.2. **find**  $\text{bestMove}$  among  $\text{candidateMoves}$  according to function **score**
- 2.3. **if** ( $\text{bestScore} < \text{score}(\text{bestMove})$ )
  - 2.3.1. **update**  $\text{bestSolution}$  and  $\text{bestScore}$  by applying  $\text{bestMove}$
  - 2.3.2. **add**  $\text{bestMove}$  to  $\text{tabuList}$  // not re-considered in the next  $t$  iterations
- 2.4. **update**  $\text{currentSolution}$  by applying  $\text{bestMove}$
3. **return**  $\text{bestSolution}$  // an MB DAG

**findAdjacencies** (Node  $Y$ , Node List  $L$ , Depth  $d$ , Significance  $\alpha$ )

1.  $A_Y := \{X_i \in L: X_i \text{ is dependent of } Y \text{ at level } \alpha\}$
2. **for**  $X_i \in A_Y$  and **for** all distinct subsets  $S \subset \{A_Y \setminus X_i\}^d$ 
  - 2.1. **if**  $X_i$  is independent of  $Y$  given  $S$  at level  $\alpha$
  - 2.2. **then** remove  $X_i$  from  $A_Y$
3. **for**  $X_i \in A_Y$ 
  - 3.1.  $A_{X_i} := \{X_j \in L: X_j \text{ is dependent of } X_i \text{ at level } \alpha, j \neq i\}$
  - 3.2. **for** all distinct subsets  $S \subset \{A_{X_i}\}^d$ 
    - 3.2.1. **if**  $X_i$  is independent of  $Y$  given  $S$  at level  $\alpha$
    - 3.2.2. **then** remove  $X_i$  from  $A_Y$
4. **return**  $A_Y$

### 5. Experiments

#### 5.1. Movie Reviews Data

We tested our method on the data set used in Pang et al [26]. This data set contains approximately 29,000 posts to the rec.arts.movies.reviews newsgroup archived at the Internet Movie Database (IMDb). The original posts are available in the form of HTML pages. Some pre-processing was performed to produce the version of the data we used. Specifically, only reviews where authors' ratings were expressed explicitly (either by stars or by numerical values) were selected. Then explicit ratings were removed and converted into one of three categories: positive, negative, or neutral. Finally, 700 positive reviews and 700 negative reviews, which the authors of the corpus judged to be more extreme, were selected for our study. Various versions of the data are available on-line [20].

#### 5.2. Feature Definition

In our study, we used words as features, where *words* are strings of letters enclosed by non-letters to the left and to the right. Note that our definition excludes punctuation sign

even though exclamation signs and question marks may be helpful for our task. Intuitively the task of sentiment extraction is a hybrid task between authorship attribution and topic categorization; we look for frequent words, possibly not related to the context, that help express lexical patterns, as well as low frequency words which may be specific to a few review styles, but very indicative of an opinion. We considered all the words that appeared in more than 8 documents as our input features, whereas words with lower counts were discarded since they appear too rarely to be helpful in the classification of many reviews. We were left with a total number of 7,716 words, as input features. In our experiments, we represented each document as a vector,  $X := [X_1, \dots, X_{7716}]$ , of the size of the initial vocabulary, where each  $X_i$  is a binary random variable that takes the value of 1 if the  $i^{th}$  word in the vocabulary is present in the document and the value of 0 otherwise.

### 5.3. Experimental Set-Up

In order to compute unbiased estimates for AUC and accuracy we used a nested, stratified, five-fold cross-validation scheme. The parameters in our experiments were the scoring criteria, the maximum size of the condition set to consider for conditional independence tests when learning the MB DAG (i.e. the depth  $d$ ), and the  $\alpha$  level to decide whether to accept or reject each of these tests. We explored 24 configurations of parameter combinations, shown in Table 1. We found the *dominant configuration* of the

**Table 1. Experimental Configurations.**

Scoring Criteria	Depth of Search	Alpha	C.V. Folds
AUC	1, 2, 3	0.001, 0.005,	5-fold
Accuracy		0.01, 0.05	

parameters on the training data and estimated the performance on the testing data, according to the (outer) five-fold cross-validation scheme. In order to find this configuration, within each fold  $i$ , we further split the training data in two ( $TR_{i1}$  and  $TR_{i2}$ ), trained the MB classifier on  $TR_{i1}$  for each parameter configuration, and tested the performance on  $TR_{i2}$ . The configuration that led to the best MB, in terms of accuracy on  $TR_{i2}$  across all five folds  $i = 1, \dots, 5$ , was chosen as the best configuration.

### 5.4. Results and Analysis

We compare the performances of our two-stage MB classifier first with single-stage MB classifier, the Markov Blanket classification without TS enhancing procedure, and then

with those of four widely used classifiers: a naïve Bayes classifier based on the multivariate Bernoulli distribution with Laplace prior for unseen words, discussed in Nigam et al. [24], a support vector machine (SVM) classifier along with a TF-IDF re-weighting of the vectors of word counts, discussed by Joachims [17], an implementation of the voted Perceptron, discussed in Freund and Schapire [10], and a maximum entropy conditional random field learner, introduced by Lafferty et al. [18].

Table 2 compares the two-stage MBC with the performances of the other classifiers using the *whole feature set* as input. As shown in the first two rows of Table 2, although MB procedure itself can identify a discriminating subset of predictors, MB procedure coupled with TS improves both AUC and accuracy, and pushes the solution towards the optimality. In this data set, it happens accidentally that the number of selected features by MB procedure is the same as those by two-stage MB Classifier, whereas the results obtained from other data sets on sentiment classification (financial news group, mergers and acquisitions news group, and mixed topic news group [9]) show that two-stage MBC is able to find an even smaller set of features with better independence structures, and produce higher accuracies, than single-stage MBC. The comparative results of two-stage MBC against the other four methods are as expected: more features did not necessarily lead to better results, as the classifiers were not able to distinguish discriminating words from noise. In such a situation we also expected the SVM with TFIDF re-weighting and the voted perceptron to perform better than the other classifiers. As shown in table 2, the two-stage MB classifier selects 22 relevant words out of 7,716 words in the vocabulary. The feature reduction ratio is 99.71%; the cross-validated AUC based on the 22 words and their dependencies is 96.85%, which is 14.3% higher than the best of the other four methods; the corresponding cross-validated accuracy is 87.5%, which is 3.5% higher than the best of the other four methods.

**Table 2. Average performances on the whole feature set.**

Method	AUC (%)	Accuracy (%)	# Selected Features	Size Reduction
Single-stage MB	71.24	65.00	22	99.71%
Two-stage MB	96.85	87.52	22	99.71%
Naïve Bayes	82.61	66.22	7,716	0%
SVM + TFIDF	81.32	84.07	7,716	0%
Voted perceptron	77.09	70.00	7,716	0%
Max. entropy	75.79	79.43	7,716	0%

We notice that the two-stage MB classifier is able to au-

tomatically identify a very discriminating subset of features (or words) that are relevant to the target variable ( $Y$ , the label of the review). Specifically, the selected features are those that form the Markov Blanket for  $Y$ . Further, the two-stage MB classifier yields the best results in terms of both cross-validated AUC and accuracy. Other methods perform worse on the whole feature set and need to be paired with a variable selection strategy.

Table 3 compares the performance of the two-stage MBC with others classifiers using the *same number of features* selected using information gain criterion. We notice that feature selection using information gain criterion does not tell us how many features have to be selected, but rather allows us to rank the features from most to least discriminating instead. Again, the two-stage MB classifier dominates the other methods both in terms of AUC and accuracy, though it is not clear whether the extra performance comes from the different feature selection strategies, or from the dependencies encoded by the MB.

**Table 3. Average performances on the same number of features.**

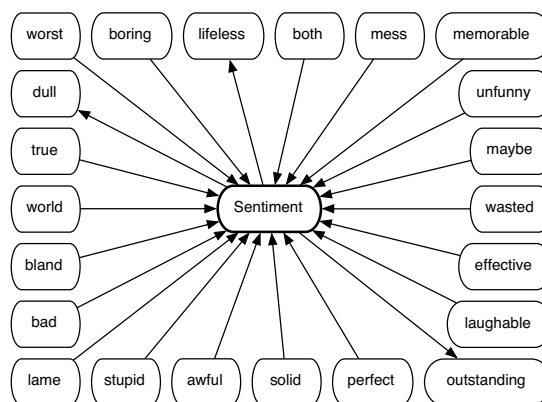
Method	AUC (%)	Accuracy (%)	# Selected Features	Size Reduction
Single-stage MB	71.24	65.00	22	99.71%
Two-stage MB	96.85	87.52	22	99.71%
Naïve Bayes	78.85	72.07	22	99.71%
SVM + TFIDF	67.30	70.43	22	99.71%
Voted perceptron	78.68	71.71	22	99.71%
Max. entropy	68.42	71.93	22	99.71%

To investigate this point, in Table 4 we compare the performance of the two-stage MBC with others classifiers using the *same exact features*. We find that a small part of the difference between the accuracy of the MBC and that of other classifiers in Table 3 arises from the fact that we selected features using information gain; in fact all the four competing classifiers performed better on the set of features in the Markov blanket. We also find that the major portion of such differences is due to the MB classification method itself. We attribute the jump in the accuracy and AUC to the fact that the MB classifier encodes and takes advantage of conditional dependencies among words, which all other methods fail to capture.

Finally, in Figure 4 below we show the best MB DAG learned by the two-stage MB classifier. All the directed edges are robust over at least 4 out of five cross validation runs; the variation is very small. The structure of the final MB DAG does not indicate independence of the words conditional on the sentiment variable, which is the strong assumption underlying all the competing classifiers.

**Table 4. Average performances on the same exact features.**

Method	AUC (%)	Accuracy (%)	# Selected Features	Size Reduction
Single-stage MB	71.24	65.00	22	99.71%
Two-stage MB	96.85	87.52	22	99.71%
Naïve Bayes	81.81	73.36	22	99.71%
SVM + TFIDF	69.47	72.00	22	99.71%
Voted perceptron	80.61	73.93	22	99.71%
Max. entropy	69.81	73.44	22	99.71%



**Figure 4. Best Fitting MB DAG for the Movie Dataset.**

These experiments, as well as more results we have obtained on other news and medical data sets (see [2] and next section), suggest that for problems where the independence assumption is not appropriate, the two-stage MB classifier is a better choice and leads to more robust predictions by: (i) selecting statistically discriminating features for the class variable, and, (ii) learning a more realistic model that allows for dependencies among the predictors. Further, according to the empirical findings in Pang et al [26], the baseline accuracy for human-selected vocabularies can be set at about 70%. Comparing the human intuition to our fully automated machine learning technique (two-stage MBC), we observe a non-negligible improvement.

## 5.5. Extensions: News Corpora

Here we present the results of more experiments that explore the accuracy of the Markov-Blanket classifier on different news topics and on problems where we try to capture opinions with more than two categories.

Recently, we tested our algorithm on more data sets, pro-

vided by courtesy of Infonic [15]. The new collections consist of five sets of 1000 news articles each, which originally appeared on news web sites, on the following topics:

- Mergers and acquisitions (M&A, 1000 documents)
- Finance (1000 documents)
- Mixed news (3 topics  $\times$  1000 documents).

These three corpora have been designed to exhibit increasing levels of specificity; M&A is the most specific corpus, mixed news is the least specific one, and the news in the Financial corpus fall somewhere in between. Furthermore, the sentiments we consider are three: *positive*, *neutral* and *negative*. Each news has been manually labeled with a document-level sentiment by three independent trained annotators, and all the documents have at least a two-way consensus for their sentiment rating. Table 5 shows the cross-validated accuracy obtained with the MBC against that of the best known competing method, as in [9]. In all cases the cross-validation is three-fold.

**Table 5. Average performances on the news corpora, with three possible opinions**

Data	Method	Accuracy (%)	# Selected Features	Size Reduction
M&A	Two-stage MB	89.95	14	99.80%
	Best competitor (SVM)	68.50	7,166	0%
Finance	Two-stage MB	96.31	16	99.78%
	Best competitor (SVM)	69.00	7,166	0%
Mixed	Two-stage MB	91.97	12	99.83%
	Best competitor (SVM)	70.20	7,166	0%

These results further establish the robustness and dominance of our algorithm, and lead us to believe that a plausible solution of the problem of extracting opinions lies in learning a local dependency structure for the sentiment variable using causal reasoning and Bayesian inference.

## 6. Discussion and Conclusions

The two-stage Markov Blanket classifier that we have proposed in this paper

- is able to capture dependencies among words, and
- is a fully automated system able to select a parsimonious vocabulary, customized for the classification task in terms of size and relevant features.

Overall, the two-stage MB classifier significantly outperforms the four baseline methods and is able to extract the most discriminating features for classification purposes. The main drawbacks of the competing methods are that they cannot automatically select relevant features, and they cannot encode the dependencies among them. While the first issue is easily overcome by combining the classifiers with off-the-shelf feature selection methods as illustrated by results in Table 4, the second issue cannot be addressed. In fact, it is a direct consequence of the assumption of pairwise independence of features underlying all the competing methods. Further, many techniques have been tried in order to automatically capture the way people express their opinions, including models for the contextual effects of negations, the use of feature frequency counts instead of their presence or absence, the use of different probability distributions for different positions of the words in the text, the use of sequences of words or  $N$ -grams, the combination of words and part of speech tags, noun-phrase chunks, and so on. However, the empirical results in terms of prediction accuracy and AUC always remain in the same ballpark.

We performed three sets of experiments to compare the methods along various dimensions, in Tables 2, 3, 4. In particular, Table 4 shows that given the *same exact features*, which were identified by the MBC as belonging to the Markov blanket, the MBC leads to significantly higher AUC and accuracy, thus suggesting that taking into account dependencies among words is crucial to perform sentiment extraction. The comparison of results of Table 3 and Table 4 suggests that information gain is not the best criterion to select discriminating variables, but the statistical tests that measure association among features and causal reasoning are better tools to perform the selection. Similar results on three data sets from different news domains which contained more complex opinions (in Table 5). Although we acknowledge that these are experimental results, the high cross-validated accuracy achieved in eight hard classification problems provides strong evidence to support the superiority of our MBC classifier.

In conclusion, we believe that in order to capture sentiments we have to go beyond the search for richer feature sets and the independence assumption. Rather we need to capture those elements of the text that help identify context and meaning. We believe that a robust model is obtained by encoding dependencies among words, and by actively searching for a better dependency structure using heuristic and optimal strategies. The MBC achieves these goals by using causal reasoning and Bayesian inference, producing excellent results for extracting consumer opinions from unstructured text data.

## Acknowledgments

The authors thank Professor Fred Glover of the University of Colorado at Boulder, Professors Clark Glymour, Peter Spirtes, and Joseph Ramsey of Carnegie Mellon University for many helpful discussions and valuable suggestions and comments, and Professor William Cohen of Carnegie Mellon University for introducing us to the problem and for valuable discussion. The authors also thank Professor Roy Lipski for providing us with the annotated corpora.

## References

- [1] E. Airoldi, A. Anderson, S. Fienberg, and K. Skinner. Who wrote Ronald Reagan radio addresses?, 2004. Manuscript.
- [2] X. Bai, C. Glymour, R. Padman, P. Spirtis, and J. Ramsey. Mb fan search classifier for large data sets with few cases. Technical Report CMU-CALD-04-102, School of Computer Science, Carnegie Mellon University, 2004.
- [3] X. Bai and R. Padman. Mb fan search classifier for large data sets with few cases. Technical Report CMU-HEINZ-2004-5, The John Heinz III School of Public Policy and Management, Carnegie Mellon University, 2004.
- [4] D. Chickering. Learning equivalence classes of bayesian-network structures. *Journal of Machine Learning Research*, 3:507–554, 2002.
- [5] D. Chickering. Optimal structure identification with greedy search. *Journal of Machine Learning Research*, 3:507–554, 2002.
- [6] W. Cohen. Minor-third: Methods for identifying names and ontological relations in text using heuristics for inducing regularities from data. <http://minorthird.sourceforge.net>, 2004.
- [7] S. Das and M. Chen. Yahoo! for amazon: Sentiment parsing from small talk on the web. In *Proceedings of the Eighth Asia Pacific Finance Association Annual Conference*. APFA, 2001.
- [8] K. Dave, S. Lawrence, and D. Pennock. Mining the peanut gallery: Opinion extraction and semantic classification of product reviews. In *Proceedings of the Twelfth International Conference on World Wide Web*, pages 519–528, 2003.
- [9] C. Engström. Topic dependence in sentiment classification. Technical Report 07-22-2004, St Edmunds College, University of Cambridge, 2004.
- [10] Y. Freund and R. Schapire. Large margin classification using the perceptron algorithm. *Machine Learning*, 37:277–296, 1999.
- [11] F. Glover. *Tabu Search*. Kluwer Academic Publishers, 1997.
- [12] F. Glover and S. Hanafi. Tabu search and finite convergence. *Discrete Applied Mathematics*, 119:3–36, 2002.
- [13] V. Hatzivassiloglou and K. McKeown. Predicting the semantic orientation of adjectives. In *Proceedings of the Eighth Conference on European Chapter of the Association for Computational Linguistics*, pages 174–181. ACL, 1997.
- [14] A. Huettner and P. Subasic. Fuzzy typing for document management. In *Association for Computational Linguistics 2000 Companion Volume: Tutorial Abstracts and Demonstration Notes*, pages 26–27, 2000.
- [15] Infonic data. <http://www.infonic.com/>.
- [16] The internet movie database. <http://www.imdb.com/>.
- [17] T. Joachims. A statistical learning model of text classification with support vector machines. In *Proceedings of the Conference on Research and Development in Information Retrieval*, pages 128–136. ACM, 2001.
- [18] J. Lafferty, A. McCallum, and F. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning*, pages 282–289. In: Proceedings of the Eighteenth International Conference on Machine Learning., 2001.
- [19] H. Liu, H. Lieberman, and T. Selker. A model of textual affect sensing using real-world knowledge. In *Proceedings of the Eighth International Conference on Intelligent User Interfaces*, pages 125–132, 2003.
- [20] Movie review data. <http://www.cs.cornell.edu/people/pabo/movie-review-data/>.
- [21] T. Mitchell. *Machine Learning*. McGraw-Hill, 1997.
- [22] F. Mosteller and D. Wallace. *Inference and Disputed Authorship: The Federalist*. Addison-Wesley, 1964.
- [23] F. Mosteller and D. Wallace. *Applied Bayesian and Classical Inference: The Case of The Federalist Papers*. Springer-Verlag, 1984.
- [24] K. Nigam, A. McCallum, S. Thrun, and T. Mitchell. Text classification from labeled and unlabeled documents using em. *Machine Learning*, 39:103–134, 2000.
- [25] C. Osgood, G. Suci, and P. Tannenbaum. *The Measurement of Meaning*. University of Illinois Press, Chicago, Illinois, 1957.
- [26] B. Pang, L. Lee, and S. Vaithyanathan. Thumbs up? sentiment classification using machine learning techniques. In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing*, pages 79–86, 2002.
- [27] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2000.
- [28] F. Provost, T. Fawcett, and R. Kohavi. The case against accuracy estimation for comparing induction algorithms. In *Proceedings of the Fifteenth International Conference on Machine Learning*, pages 445–453, 1998.
- [29] C. Rego and B. Alidaee. *Metaheuristic Optimization via Memory and Evolution: Tabu Search and Scatter Search*. Kluwer Academic Publishers, 2004.
- [30] P. Spirtes, C. Glymour, and R. Scheines. *Causation, Prediction, and Search*. MIT Press, 2000.
- [31] P. Spirtes and C. Meek. Learning bayesian networks with discrete variables from data. In *Proceedings of the First International Conference on Knowledge Discovery and Data Mining*, pages 294–299. AAAI Press, 1995.
- [32] Tabu search web site. <http://www.tabusearch.net/>.
- [33] P. Turney. Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews. In *Proceedings Fortieth Annual Meeting of the Association for Computational Linguistics*, pages 417–424, 2002.
- [34] P. Turney and M. Littman. Unsupervised learning of semantic orientation from a hundred-billion-word corpus. Technical Report EGB-1094, National Research Council, Canada, 2002.