

# Unsupervised Learning for Lexicon-Based Classification

Jacob Eisenstein

Georgia Institute of Technology  
801 Atlantic Drive NW  
Atlanta, Georgia 30318

## Abstract

In lexicon-based classification, documents are assigned labels by comparing the number of words that appear from two opposed lexicons, such as positive and negative sentiment. Creating such words lists is often easier than labeling instances, and they can be debugged by non-experts if classification performance is unsatisfactory. However, there is little analysis or justification of this classification heuristic. This paper describes a set of assumptions that can be used to derive a probabilistic justification for lexicon-based classification, as well as an analysis of its expected accuracy. One key assumption behind lexicon-based classification is that all words in each lexicon are equally predictive. This is rarely true in practice, which is why lexicon-based approaches are usually outperformed by supervised classifiers that learn distinct weights on each word from labeled instances. This paper shows that it is possible to learn such weights without labeled data, by leveraging co-occurrence statistics across the lexicons. This offers the best of both worlds: light supervision in the form of lexicons, and data-driven classification with higher accuracy than traditional word-counting heuristics.

## Introduction

**Lexicon-based classification** refers to a classification rule in which documents are assigned labels based on the count of words from lexicons associated with each label (Taboada et al. 2011). For example, suppose that we have opposed labels  $Y \in \{0, 1\}$ , and we have associated lexicons  $\mathcal{W}_0$  and  $\mathcal{W}_1$ . Then for a document with a vector of word counts  $\mathbf{x}$ , the lexicon-based decision rule is,

$$\sum_{i \in \mathcal{W}_0} x_i \geq \sum_{j \in \mathcal{W}_1} x_j, \quad (1)$$

where the  $\geq$  operator indicates a decision rule. Put simply, the rule is to select the label whose lexicon matches the most word tokens.

Lexicon-based classification is widely used in industry and academia, with applications ranging from sentiment classification and opinion mining (Pang and Lee 2008; Liu 2015) to the psychological and ideological analysis of texts (Laver and Garry 2000; Tausczik and Pennebaker 2010). The popularity of this approach can be explained by

its relative simplicity and ease of use: for domain experts, creating lexicons is intuitive, and, in comparison with labeling instances, it may offer a faster path towards a reasonably accurate classifier (Settles 2011). Furthermore, classification errors can be iteratively debugged by refining the lexicons.

However, from a machine learning perspective, there are a number of drawbacks to lexicon-based classification. First, while intuitively reasonable, lexicon-based classification lacks theoretical justification: it is not clear what conditions are necessary for it to work. Second, the lexicons may be incomplete, even for designers with strong substantive intuitions. Third, lexicon-based classification assigns an equal weight to each word, but some words may be more strongly predictive than others.<sup>1</sup> Fourth, lexicon-based classification ignores multi-word phenomena, such as negation (e.g., *not so good*) and discourse (e.g., *the movie would be watchable if it had better acting*). Supervised classification systems, which are trained on labeled examples, tend to outperform lexicon-based classifiers, even without accounting for multi-word phenomena (Liu 2015; Pang and Lee 2008).

Several researchers have proposed methods for **lexicon expansion**, automatically growing lexicons from an initial seed set (Hatzivassiloglou and McKeown 1997; Qiu et al. 2011). There is also work on handling multi-word phenomena such as negation (Wilson, Wiebe, and Hoffmann 2005; Polanyi and Zaenen 2006), and discourse (Somasundaran, Wiebe, and Ruppenhofer 2008; Bhatia, Ji, and Eisenstein 2015). However, the theoretical foundations of lexicon-based classification remain poorly understood, and we lack principled means for automatically assigning weights to lexicon items without resorting to labeled instances.

This paper elaborates a set of assumptions under which lexicon-based classification is equivalent to Naïve Bayes classification. I then derive expected error rates under these assumptions. These expected error rates are not matched by observations on real data, suggesting that the underlying assumptions are invalid. Of key importance is the assumption that each lexicon item is equally predictive. To relax this as-

<sup>1</sup>Some lexicons attach coarse-grained predefined weights to each word. For example, the OpinionFinder Subjectivity lexicon labels words as “strongly” or “weakly” subjective (Wilson, Wiebe, and Hoffmann 2005). This poses an additional burden on the lexicon creator.

sumption, I derive a principled method for estimating word probabilities under each label, using a method-of-moments estimator on cross-lexical co-occurrence counts.

Overall, this paper makes the following contributions:

- justifying lexicon-based classification as a special case of multinomial Naïve Bayes;
- mathematically analyzing this model to compute the expected performance of lexicon-based classifiers;
- extending the model to justify a popular variant of lexicon-based classification, which incorporates word presence rather than raw counts;
- deriving a method-of-moments estimator for the parameters of this model, enabling lexicon-based classification with unique weights per word, without labeled data;
- empirically demonstrating that this classifier outperforms lexicon-based classification and alternative approaches.

### Lexicon-Based Classification as Naïve Bayes

I begin by showing how the lexicon-based classification rule shown in (1) can be derived as a special case of Naïve Bayes classification. Suppose we have a prior probability  $P_Y$  for the label  $Y$ , and a likelihood function  $P_{X|Y}$ , where  $X$  is a random variable corresponding to a vector of word counts. The conditional label probability can be computed by Bayesian inversion,

$$P(y | \mathbf{x}) = \frac{P(\mathbf{x} | y)P(y)}{\sum_{y'} P(\mathbf{x} | y')P(y')}. \quad (2)$$

Assuming that the costs for each type of misclassification error are identical, then the minimum Bayes risk classification rule is,

$$\begin{aligned} \log \Pr(Y = 0) + \log P(\mathbf{x} | Y = 0) \\ \geq \log \Pr(Y = 1) + \log P(\mathbf{x} | Y = 1), \end{aligned} \quad (3)$$

moving to the log domain for simplicity of notation. I now show that lexicon-based classification can be justified under this decision rule, given a set of assumptions about the probability distributions.

Let us introduce some assumptions about the likelihood function,  $P_{X|Y}$ . The random variable  $X$  is defined over vectors of counts, so a natural choice for the form of this likelihood is the multinomial distribution, corresponding to a multinomial Naïve Bayes classifier. For a specific vector of counts  $X = \mathbf{x}$ , write  $P(\mathbf{x} | y) \triangleq P_{\text{multinomial}}(\mathbf{x}; \boldsymbol{\theta}_y, N)$ , where  $\boldsymbol{\theta}_y$  is a probability vector associated with label  $y$ , and  $N = \sum_{i=1}^V x_i$  is the total count of tokens in  $\mathbf{x}$ , and  $x_i$  is the count of word  $i \in \{1, 2, \dots, V\}$ . The multinomial likelihood is proportional to a product of likelihoods of categorical variables corresponding to individual words (tokens),

$$\Pr(W = i | Y = y; \boldsymbol{\theta}) = \theta_{y,i}, \quad (4)$$

where the random variable  $W$  corresponds to a single token, whose probability of being word type  $i$  is equal to  $\theta_{y,i}$  in a document with label  $y$ . The multinomial log-likelihood can be written as,

$$\begin{aligned} \log P(\mathbf{x} | y) &= \log P_{\text{multinomial}}(\mathbf{x}; \boldsymbol{\theta}_y, N) \\ &= K(\mathbf{x}) + \sum_{i=1}^V x_i \log \Pr(W = i | Y = y; \boldsymbol{\theta}) \\ &= K(\mathbf{x}) + \sum_{i=1}^V x_i \log \theta_{y,i}, \end{aligned} \quad (5)$$

where  $K(\mathbf{x})$  is a function of  $\mathbf{x}$  that is constant in  $y$ .

The first necessary assumption about the likelihood function is that the lexicons are **complete**: words that are in neither lexicon have identical probability under both labels. Formally, for any word  $i \notin \mathcal{W}_0 \cup \mathcal{W}_1$ , we assume,

$$\Pr(W = i | Y = 0) = \Pr(W = i | Y = 1), \quad (6)$$

which implies that these words are irrelevant to the classification boundary.

Next, we must assume that each in-lexicon word is **equally predictive**. Specifically, for words that are in lexicon  $y$ ,

$$\frac{\Pr(W = i | Y = y)}{\Pr(W = i | Y = \neg y)} = \frac{1 + \gamma}{1 - \gamma}, \quad (7)$$

where  $\neg y$  is the opposite label from  $y$ . The parameter  $\gamma$  controls the predictiveness of the lexicon: for example, if  $\gamma = 0.5$  in a sentiment classification problem, this would indicate that words in the positive sentiment lexicon are three times more likely to appear in documents with positive sentiment than in documents with negative sentiment, and vice versa. The word *atrocious* might be less likely overall than *good*, but still three times more likely in the negative class than in the positive class. In the limit,  $\gamma = 0$  implies that the lexicons do not distinguish the classes at all, and  $\gamma = 1$  implies that the lexicons distinguish the classes perfectly, so that the observation of a single in-lexicon word would completely determine the document label.

The conditions enumerated in (6) and (7) are ensured by the following definition,

$$\theta_{y,i} = \begin{cases} (1 + \gamma)\mu_i, & i \in \mathcal{W}_y \\ (1 - \gamma)\mu_i, & i \in \mathcal{W}_{\neg y} \\ \mu_i, & i \notin \mathcal{W}_y \cup \mathcal{W}_{\neg y}, \end{cases} \quad (8)$$

where  $\neg y$  is the opposite label from  $y$ , and  $\boldsymbol{\mu}$  is a vector of baseline probabilities, which are independent of the label.

Because the probability vectors  $\boldsymbol{\theta}_0$  and  $\boldsymbol{\theta}_1$  must each sum to one, we require an assumption of **equal coverage**,

$$\sum_{i \in \mathcal{W}_0} \mu_i = \sum_{j \in \mathcal{W}_1} \mu_j. \quad (9)$$

Finally, assume that the labels have **equal prior likelihood**,  $\Pr(Y = 0) = \Pr(Y = 1)$ . It is trivial to relax this assumption by adding a constant term to one side of the decision rule in (1).

With these assumptions in hand, it is now possible to simplify the decision rule in (3). Thanks to the assumption of

equal prior probability, we can drop the priors  $P(Y)$ , so that the decision rule is a comparison of the likelihoods,

$$\log P(\mathbf{x} | Y = 0) \geq \log P(\mathbf{x} | Y = 1) \quad (10)$$

$$K(\mathbf{x}) + \sum_i x_i \log \theta_{0,i} \geq K(\mathbf{x}) + \sum_i x_i \log \theta_{1,i}. \quad (11)$$

Canceling  $K(\mathbf{x})$  and applying the definition from (8),

$$\begin{aligned} & \sum_{i \in \mathcal{W}_0} x_i \log((1 + \gamma)\mu_i) + \sum_{i \in \mathcal{W}_1} x_i \log((1 - \gamma)\mu_i) \\ & \geq \sum_{i \in \mathcal{W}_0} x_i \log((1 - \gamma)\mu_i) + \sum_{i \in \mathcal{W}_1} x_i \log((1 + \gamma)\mu_i). \end{aligned} \quad (12)$$

The  $\mu_i$  terms cancel after distributing the log, leaving,

$$\sum_{i \in \mathcal{W}_0} x_i \log \frac{1 + \gamma}{1 - \gamma} \geq \sum_{i \in \mathcal{W}_1} x_i \log \frac{1 + \gamma}{1 - \gamma}. \quad (13)$$

For any  $\gamma \in (0, 1)$ , the term  $\log \frac{1 + \gamma}{1 - \gamma}$  is a finite and positive constant. Therefore, (13) is identical to the counting-based classification rule in (1). In other words, lexicon-based classification is minimum Bayes risk classification in a multinomial probability model, under the assumptions of equal prior likelihood, lexicon completeness, equal predictiveness of words, and equal coverage.

### Analysis of Lexicon-Based Classification

One advantage of deriving a formal foundation for lexicon-based classification is that it is possible to analyze its expected performance. For a label  $y$ , let us write the count of in-lexicon words as  $m_y = \sum_{i \in \mathcal{W}_y} x_i$ , and the count of opposite-lexicon words as  $m_{-y} = \sum_{i \in \mathcal{W}_{-y}} x_i$ . Lexicon-based classification makes a correct prediction whenever  $m_y > m_{-y}$  for the correct label  $y$ . To assess the likelihood that  $m_y > m_{-y}$ , it is sufficient to compute the expectation and variance of the difference  $m_y - m_{-y}$ ; under the central limit theorem, we can treat this difference as approximately normally distributed, and compute the probability that the difference is positive using the Gaussian cumulative distribution function (CDF).

Let us use the convenience notation  $s_\mu$ ,

$$s_\mu \triangleq \sum_{i \in \mathcal{W}_0} \mu_i = \sum_{i \in \mathcal{W}_1} \mu_i. \quad (14)$$

Recall that we have already taken the assumption that the sums of baseline word probabilities for the two lexicons are equal. Under the multinomial probability model, given a document with  $N$  tokens, the expected counts are,

$$E[m_y] = N \sum_{i \in \mathcal{W}_y} \theta_{y,i} = N(1 + \gamma)s_\mu \quad (15)$$

$$E[m_{-y}] = N \sum_{i \in \mathcal{W}_{-y}} \theta_{-y,i} = N(1 - \gamma)s_\mu \quad (16)$$

$$E[m_y - m_{-y}] = 2N\gamma s_\mu. \quad (17)$$

Next we compute the variance of this margin,

$$V[m_y - m_{-y}] = V[m_y] + V[m_{-y}] + Cov(m_y, m_{-y}). \quad (18)$$

Each of these terms is the variance of a sum of counts. Under the multinomial distribution, the variance of a single count is  $V[x_i] = N\theta_i(1 - \theta_i)$ . The variance of the sum  $m_y$  is,

$$\begin{aligned} V[m_y] &= \sum_{i \in \mathcal{W}_y} N\theta_i(1 - \theta_i) - \sum_{j \in \mathcal{W}_y, j \neq i} N\theta_i\theta_j \\ &= \sum_{i \in \mathcal{W}_y} N\theta_i - N\theta_i^2 - \sum_{j \in \mathcal{W}_y, j \neq i} N\theta_i\theta_j \\ &\leq N \sum_{i \in \mathcal{W}_y} \theta_i = N \sum_{i \in \mathcal{W}_y} (1 + \gamma)\mu_i \\ &= N(1 + \gamma)s_\mu. \end{aligned} \quad (19)$$

An equivalent upper bound can be computed for the variance of the count of opposite lexicon words,

$$V[m_{-y}] \leq N(1 - \gamma)s_\mu. \quad (21)$$

These bounds are fairly tight because the products of probabilities  $\theta_i^2$  and  $\theta_i\theta_j$  are nearly always small, due to the fact that most words are rare. Because the covariance  $Cov(m_y, m_{-y})$  is negative (and also involves a product of word probabilities), we can further bound the variance of the margin, obtaining the upper bound,

$$V[m_y - m_{-y}] \leq N(1 + \gamma)s_\mu + N(1 - \gamma)s_\mu = 2Ns_\mu. \quad (22)$$

By the central limit theorem, the margin  $m_y - m_{-y}$  is approximately normally distributed, with mean  $2N\gamma s_\mu$  and variance upper-bounded by  $2Ns_\mu$ . The probability of making a correct prediction (which occurs when  $m_y > m_{-y}$ ) is then equal to the cumulative density of a standard normal distribution  $\Phi(z)$ , where the  $z$ -score is equal to the ratio of the expectation and the standard deviation,

$$z = \frac{E[m_y - m_{-y}]}{\sqrt{V[m_y - m_{-y}]}} \geq \frac{2N\gamma s_\mu}{\sqrt{2Ns_\mu}} = \gamma\sqrt{2Ns_\mu}. \quad (23)$$

Note that by upper-bounding the variance, we obtain a lower bound on the  $z$ -score, and thus a lower bound on the expected accuracy.

According to this approximation, accuracy is expected to increase with the predictiveness  $\gamma$ , the document length  $N$ , and the lexicon coverage  $s_\mu$ . This helps to explain a dilemma in lexicon design: as more words are added, the coverage increases, but the average predictiveness of each word decreases (assuming the most predictive words are added first). Thus, increasing the size of a lexicon by adding marginal words may not improve performance.

The analysis also predicts that longer documents should be easier to classify. This is because the expected size of the gap  $m_y - m_{-y}$  grows with  $N$ , while its standard deviation grows only with  $\sqrt{N}$ . This prediction can be tested empirically, and on all four datasets considered in this paper, it is false: longer documents are harder to classify accurately. This is a clue that the underlying assumptions are not valid. The decreased accuracy for especially long reviews may be due to these reviews being more complex, perhaps requiring modeling of the discourse structure (Somasundaran, Wiebe, and Ruppenhofer 2008).

## Justifying the Word-Appearance Heuristic

An alternative heuristic to lexicon-based classification is to consider only the **presence** of each word type, and not its count. This corresponds to the decision rule,

$$\sum_{i \in \mathcal{W}_0} \delta(x_i > 0) \geq \sum_{j \in \mathcal{W}_1} \delta(x_j > 0), \quad (24)$$

where  $\delta(\cdot)$  is a delta function that returns one if the Boolean condition is true, and zero otherwise. In the context of supervised classification, Pang, Lee, and Vaithyanathan (2002) find that word presence is a more predictive feature than word frequency. By ignoring repeated mentions of the same word, heuristic (24) emphasizes the diversity of ways in which a document covers a lexicon, and is more robust to document-specific idiosyncrasies — such as a review of *The Joy Luck Club*, which might include the positive words *joy* and *luck* many times even if the review is negative.

The word-appearance heuristic can also be explained in the framework defined above. The multinomial likelihood  $P_{X|Y}$  can be replaced by a **Dirichlet-compound multinomial** (DCM) distribution, also known as a multivariate Polya distribution (Madsen, Kauchak, and Elkan 2005). This distribution is written  $P_{\text{dcm}}(\mathbf{x}; \boldsymbol{\alpha}_y)$ , where  $\boldsymbol{\alpha}_y$  is a vector of parameters associated with label  $y$ , with  $\alpha_{y,i} > 0$  for all  $i \in \{1, 2, \dots, V\}$ . The DCM is a “compound” distribution because it treats the parameter of the multinomial as a latent variable to be marginalized out,

$$P_{\text{dcm}}(\mathbf{x}; \boldsymbol{\alpha}_y) = \int_{\boldsymbol{\nu}} P_{\text{multinomial}}(\mathbf{x} | \boldsymbol{\nu}) P_{\text{Dirichlet}}(\boldsymbol{\nu} | \boldsymbol{\alpha}_y) d\boldsymbol{\nu}. \quad (25)$$

Intuitively, one can think of the DCM distribution as encoding a model in which each document has its own multinomial distribution over words; this document-specific distribution is itself drawn from a prior that depends on the class label  $y$ .

Suppose we set the DCM parameter  $\boldsymbol{\alpha} = \tau \boldsymbol{\theta}$ , with  $\boldsymbol{\theta}$  as defined in (8). The constant  $\tau > 0$  is then the **concentration** of the distribution: as  $\tau$  grows, the probability distribution over  $\boldsymbol{\alpha}$  is more closely concentrated around the prior expectation  $\boldsymbol{\theta}$ . Because  $\sum_i \theta_i = 1$ , the likelihood function under this model is,

$$P_{\text{dcm}}(\mathbf{x} | y) = \frac{\Gamma(\tau)}{\Gamma(N + \tau)} \prod_i \frac{\Gamma(x_i + \tau \theta_{y,i})}{\Gamma(\tau \theta_{y,i})}, \quad (26)$$

where  $\Gamma(\cdot)$  is the gamma function. Minimum Bayes risk classification in this model implies the decision rule:

$$\sum_{i \in \mathcal{W}_0} \log \frac{r_{\text{in}}(x_i)}{r_{\text{out}}(x_i)} \geq \sum_{i \in \mathcal{W}_1} \log \frac{r_{\text{in}}(x_{t,i})}{r_{\text{out}}(x_i)} \quad (27)$$

where,

$$r_{\text{in}}(x_i) \triangleq \frac{\Gamma(x_i + \tau(1 + \gamma)\mu_i)}{\Gamma(\tau(1 + \gamma)\mu_i)} \quad (28)$$

$$r_{\text{out}}(x_i) \triangleq \frac{\Gamma(x_i + \tau(1 - \gamma)\mu_i)}{\Gamma(\tau(1 - \gamma)\mu_i)}. \quad (29)$$

As  $\tau \rightarrow \infty$ , the prior on  $\boldsymbol{\nu}$  is tightly linked to  $\boldsymbol{\theta}$ , so that the model reduces to the multinomial defined above. Another

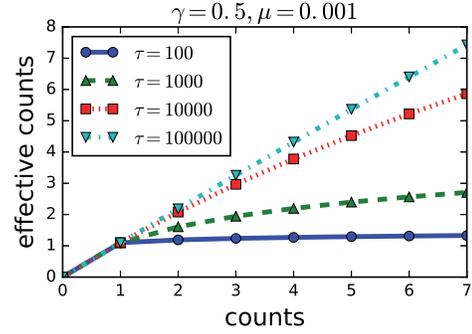


Figure 1: Effective counts for varying values of  $\tau$ . For the datasets considered in this paper,  $\tau$  usually falls in the range between 500 and 1000.

way to see this is to apply the equality  $\Gamma(x + 1) = x\Gamma(x)$  to (28) and (29) when  $\tau\mu_i \gg x_i$ . As  $\tau \rightarrow 0$ , the prior on  $\boldsymbol{\nu}$  becomes increasingly diffuse. Repeated counts of any word are better explained by document-specific variation from the prior, than by properties of the label. This situation is shown in Figure 1, which plots the “effective counts” implied by the classification rule (27) for a range of values of the concentration parameter  $\tau$ , holding the other parameters constant ( $\mu = 10^{-3}, \gamma = 0.5$ ). For high values of  $\tau$ , the effective counts track the observed counts linearly, as in the multinomial model; for low values of  $\tau$ , the effective counts barely increase beyond 1.

Minka (2012) presents a number of estimators for the concentration parameter  $\tau$  from a corpus of text. When the label  $y$  is unknown, we cannot apply these estimators directly. However, as described above, out-of-lexicon words are assumed to have identical probability under both labels. This assumption can be exploited to estimate  $\tau$  exclusively from the first and second moments of these out-of-lexicon words. Analysis of the expected accuracy of this model is left to future work.

## Estimating Word Predictiveness

A crucial simplification made by lexicon-based classification is that all words in each lexicon are equally predictive. In reality, words may be more or less predictive of class labels, for reasons such as sense ambiguity (e.g., *well*) and degree (e.g., *good* vs *flawless*). By introducing a per-word predictiveness factor  $\gamma_i$  into (8), we arrive at a model that is a restricted form of Naïve Bayes. (The restriction is that the probabilities of non-lexicon words are constrained to be identical across classes.) If labeled data were available, this model could be estimated by maximum likelihood. This section shows how to estimate the model without labeled data, using the method of moments.

First, note that the baseline probabilities  $\mu_i$  can be estimated directly from counts on an unlabeled corpus; the challenge is to estimate the parameters  $\gamma_i$  for all words in the two lexicons. The key intuition that makes this possible is that highly predictive words should rarely appear with words in the opposite lexicon. This idea can be formalized

in terms of **cross-label counts**: the cross-label count  $c_i$  is the co-occurrence count of word  $i$  with all words in the opposite lexicon,

$$c_i = \sum_{t=1}^T \sum_{j \in \mathcal{W}_{-y}} x_i^{(t)} x_j^{(t)}, \quad (30)$$

where  $\mathbf{x}^{(t)}$  is the vector of word counts for document  $t$ , with  $t \in \{1 \dots T\}$ . Under the multinomial model defined above, for a single document with  $N$  tokens, the expected product of counts for a word pair  $(i, j)$  is equal to,

$$\begin{aligned} E[x_i x_j] &= E[x_i] E[x_j] + \text{Cov}(x_i, x_j) \\ &= N \theta_i N \theta_j - N \theta_i \theta_j \\ &= N(N-1) \theta_i \theta_j. \end{aligned} \quad (31)$$

Let us focus on the expected products of counts for cross-lexicon word pairs  $(i \in \mathcal{W}_0, j \in \mathcal{W}_1)$ . The parameter  $\theta$  depends on the document label  $y$ , as defined in (8). As a result, we have the following expectations,

$$\begin{aligned} E[x_i x_j | Y = 0] &= N(N-1) \mu_i (1 + \gamma_i) \mu_j (1 - \gamma_j) \\ &= N(N-1) \mu_i \mu_j (1 + \gamma_i - \gamma_j - \gamma_i \gamma_j) \end{aligned} \quad (32)$$

$$\begin{aligned} E[x_i x_j | Y = 1] &= N(N-1) \mu_i (1 - \gamma_i) \mu_j (1 + \gamma_j) \\ &= N(N-1) \mu_i \mu_j (1 - \gamma_i + \gamma_j - \gamma_i \gamma_j) \end{aligned} \quad (33)$$

$$\begin{aligned} E[x_i x_j] &= P(Y = 0) E[x_i x_j | Y = 0] \\ &\quad + P(Y = 1) E[x_i x_j | Y = 1] \\ &= N(N-1) \mu_i \mu_j (1 - \gamma_i \gamma_j). \end{aligned} \quad (34)$$

Summing over all words  $j \in \mathcal{W}_1$  and all documents  $t$ ,

$$\begin{aligned} E[c_i] &= \sum_{t=1}^T \sum_{j \in \mathcal{W}_1} E[x_i^{(t)} x_j^{(t)}] \\ &= \sum_{t=1}^T N_t (N_t - 1) \mu_i \sum_{j \in \mathcal{W}_1} \mu_j (1 - \gamma_i \gamma_j) \end{aligned} \quad (35)$$

Let us write  $\boldsymbol{\gamma}^{(1)}$  to indicate the vector of  $\gamma_j$  parameters for all  $j \in \mathcal{W}_1$ , and  $\boldsymbol{\gamma}^{(0)}$  for all  $i \in \mathcal{W}_0$ . The expectation in (35) is a linear function of  $\gamma_i$ , and a linear function of the vector  $\boldsymbol{\gamma}^{(1)}$ . Analogously, for all  $j \in \mathcal{W}_1$ ,  $E[c_j]$  is a linear function of  $\gamma_j$  and  $\boldsymbol{\gamma}^{(1)}$ . Our goal is to choose  $\boldsymbol{\gamma}$  so that the expectations  $E[c_i]$  closely match the observed counts  $c_i$ . This can be viewed as form of **method of moments** estimation, with the following objective,

$$J = \frac{1}{2} \sum_{i \in \mathcal{W}_0} (c_i - E[c_i])^2 + \frac{1}{2} \sum_{j \in \mathcal{W}_1} (c_j - E[c_j])^2, \quad (36)$$

which can be minimized in terms of  $\boldsymbol{\gamma}^{(0)}$  and  $\boldsymbol{\gamma}^{(1)}$ . However, there is an additional constraint: the probability distributions  $\boldsymbol{\theta}_0$  and  $\boldsymbol{\theta}_1$  must still sum to one. We can express this as a linear constraint on  $\boldsymbol{\gamma}^{(0)}$  and  $\boldsymbol{\gamma}^{(1)}$ ,

$$\boldsymbol{\mu}^{(0)} \cdot \boldsymbol{\gamma}^{(0)} - \boldsymbol{\mu}^{(1)} \cdot \boldsymbol{\gamma}^{(1)} = 0, \quad (37)$$

where  $\boldsymbol{\mu}^{(y)}$  is the vector of baseline probabilities for words  $i \in \mathcal{W}_y$ , and  $\boldsymbol{\mu}^{(0)} \cdot \boldsymbol{\gamma}^{(0)}$  indicates a dot product.

We therefore formulate the following constrained optimization problem,

$$\begin{aligned} \min_{\boldsymbol{\gamma}^{(0)}, \boldsymbol{\gamma}^{(1)}} & \frac{1}{2} \sum_{i \in \mathcal{W}_0} (c_i - E[c_i])^2 + \frac{1}{2} \sum_{j \in \mathcal{W}_1} (c_j - E[c_j])^2 \\ \text{s.t.} & \boldsymbol{\mu}^{(0)} \cdot \boldsymbol{\gamma}^{(0)} - \boldsymbol{\mu}^{(1)} \cdot \boldsymbol{\gamma}^{(1)} = 0 \\ & \forall i \in (\mathcal{W}_0 \cup \mathcal{W}_1), \gamma_i \in [0, 1]. \end{aligned} \quad (38)$$

This problem can be solved by **alternating direction method of multipliers** (Boyd et al. 2011). The equality constraint can be incorporated into an augmented Lagrangian,

$$\begin{aligned} L_\rho(\boldsymbol{\gamma}^{(0)}, \boldsymbol{\gamma}^{(1)}) &= \frac{1}{2} \sum_{i \in \mathcal{W}_0} (c_i - E[c_i])^2 + \frac{1}{2} \sum_{j \in \mathcal{W}_1} (c_j - E[c_j])^2 \\ &\quad + \frac{\rho}{2} (\boldsymbol{\mu}^{(0)} \cdot \boldsymbol{\gamma}^{(0)} - \boldsymbol{\mu}^{(1)} \cdot \boldsymbol{\gamma}^{(1)})^2, \end{aligned} \quad (39)$$

where  $\rho > 0$  is the penalty parameter. The augmented Lagrangian is biconvex in  $\boldsymbol{\gamma}^{(0)}$  and  $\boldsymbol{\gamma}^{(1)}$ , which suggests an iterative solution (Boyd et al. 2011, page 76). Specifically, we hold  $\boldsymbol{\gamma}^{(1)}$  fixed and solve for  $\boldsymbol{\gamma}^{(0)}$ , subject to  $\gamma_i \in [0, 1]$  for all  $i \in \mathcal{W}_0$ . We then solve for  $\boldsymbol{\gamma}^{(1)}$  under the same conditions. Finally, we update a dual variable  $u$ , representing the extent to which the equality constraint is violated. These updates are iterated until convergence. The unconstrained local updates to  $\boldsymbol{\gamma}^{(0)}$  and  $\boldsymbol{\gamma}^{(1)}$  can be computed by solving a system of linear equations, and the result can be projected back onto the feasible region. The penalty parameter  $\rho$  is initialized at 1, and then dynamically updated based on the primal and dual residuals (Boyd et al. 2011, pages 20-21). More details are available in the online supplement and source code.<sup>2</sup>

## Evaluation

An empirical evaluation is performed on four datasets in two languages. All datasets involve binary classification problems, and performance is quantified by the **area-under-the-curve (AUC)**, a measure of classification performance that is robust to unbalanced class distributions. A perfect classifier achieves  $\text{AUC} = 1$ ; in expectation, a random decision rule gives  $\text{AUC} = 0.5$ .

**Datasets** The proposed method relies on co-occurrence counts, and therefore is best suited to documents containing at least a few sentences each. With this in mind, the following datasets are used in the evaluation:

**Amazon** English-language product reviews across four domains; of these reviews, 8000 are labeled and another 19677 are unlabeled (Blitzer, Dredze, and Pereira 2007).

**Cornell** 2000 English-language film reviews (version 2.0), labeled as positive or negative (Pang and Lee 2004).

**CorpusCine** 3800 Spanish-language movie reviews, rated on a scale of one to five (Vilares, Alonson, and Gómez-Rodríguez 2015). Ratings of four or five are considered as positive; ratings of one and two are considered as negative. Reviews with a rating of three are excluded.

<sup>2</sup>Online supplement: <http://link.to/appendix>. Source code: <https://github.com/jacobeisenstein/probabilistic-lexicon-classification>

**IMDB** 50,000 English-language film reviews (Maas et al. 2011). This evaluation includes only the test set of 25,000 reviews, of which half are positive and half are negative.

**Lexicons** Preliminary evaluation compared several English-language sentiment lexicons. The Liu lexicon (Liu 2015) consistently obtained the best performance on all three English-language datasets, so it was made the focus of all subsequent experiments. Ribeiro et al. (2016) also found that the Liu lexicon is one of the strongest lexicons for review analysis. For the Spanish data, the ISOL lexicon was used (Molina-González et al. 2013). It is a modified translation of the Liu lexicon.

**Classifiers** The evaluation compares the following unsupervised classification strategies:

**LEXICON** basic word counting, as in decision rule (1);

**LEX-PRESENCE** counting word presence rather than frequency, as in decision rule (24);

**PROBLEX-MULT** probabilistic lexicon-based classification, as proposed in this paper, using the multinomial likelihood model;

**PROBLEX-DCM** probabilistic lexicon-based classification, using the Dirichlet Compound Multinomial likelihood to reduce effective counts for repeated words;

**PMI** An alternative approach, discussed in the related work, is to impute document labels from a seed set of words, and then compute “sentiment scores” for individual words from pointwise mutual information between the words and imputed labels (Turney 2002). The implementation of this method is based on the description from Kiritchenko, Zhu, and Mohammad (2014), using the lexicons as the seed word sets.

As an upper bound, a supervised logistic regression classifier is also considered. This classifier is trained using five-fold cross validation. It is the only classifier with access to training data. For the PROBLEX-MULT and PROBLEX-DCM methods, lexicon words which co-occur with the opposite lexicon at greater than chance frequency are eliminated from the lexicon in a preprocessing step.

**Results** Results are shown in Table 1. The superior performance of the logistic regression classifier confirms the principle that supervised classification is far more accurate than lexicon-based classification. Therefore, supervised classification should be preferred when labeled data is available. Nonetheless, the probabilistic lexicon-based classifiers developed in this paper (PROBLEX-MULT and PROBLEX-DCM) go a considerable way towards closing the gap, with improvements in AUC ranging from less than 1% on the CorpusCine data to nearly 8% on the IMDB data. The PMI approach performs poorly, improving over the simpler lexicon-based classifiers on only one of the four datasets. The word presence heuristic offers no consistent improvements, and the Bayesian adjustment to the classification rule (PROBLEX-DCM) offers only modest improvements on two of the four datasets.

	Amazon	Cornell	Cine	IMDB
LEXICON	.820	.765	.636	.807
LEX-PRESENCE	.820	.770	.638	.805
PMI	.793	.761	.638	.868
PROBLEX-MULT	.832	.810	.644	<b>.884</b>
PROBLEX-DCM	<b>.836</b>	<b>.824</b>	<b>.645</b>	<b>.884</b>
LOGREG	.897	.914	.889	.955

Table 1: Area-under-the-curve (AUC) for all classifiers. The best unsupervised result is shown in bold for each dataset.

## Related work

Turney (2002) uses pointwise mutual information to estimate the “semantic orientation” of all vocabulary words from co-occurrence with a small seed set. This approach has later been extended to the social media domain by using emoticons as the seed set (Kiritchenko, Zhu, and Mohammad 2014). Like the approach proposed here, the basic intuition is to leverage co-occurrence statistics to learn weights for individual words; however, PMI is a heuristic score that is not justified by a probabilistic model of the text classification problem. PMI-based classification underperforms PROBLEX-MULT and PROBLEX-DCM on all four datasets in our evaluation.

The method-of-moments has become an increasingly popular estimator in unsupervised machine learning, with applications in topic models (Anandkumar et al. 2014), sequence models (Hsu, Kakade, and Zhang 2012), and more elaborate linguistic structures (Cohen et al. 2014). Of particular relevance are “anchor word” techniques for learning latent topic models (Arora, Ge, and Moitra 2012). In these methods, each topic is defined first by a few keywords, which are assumed to be generated only from a single topic. From these anchor words and co-occurrence statistics, the topic-word probabilities can be recovered. A key difference is that the strong anchor word assumption is not required in this work: none of the words are assumed to be perfectly predictive of either label. We require only the much weaker assumption that words in a lexicon tend to co-occur less frequently with words in the opposite lexicon.

## Conclusion

Lexicon-based classification is a popular heuristic that has not previously been analyzed from a machine learning perspective. This analysis yields two techniques for improving unsupervised binary classification: a method-of-moments estimator for word predictiveness, and a Bayesian adjustment for repeated counts of the same word. The method-of-moments estimator yields substantially better performance than conventional lexicon-based classification, without requiring any additional annotation effort. Future work will consider the generalization to multi-class classification, and more ambitiously, the extension to multiword units.

**Acknowledgment** This research was supported by the National Institutes of Health under award number R01GM112697-01, and by the Air Force Office of Scientific Research.

## References

- Anandkumar, A.; Ge, R.; Hsu, D.; Kakade, S. M.; and Telgarsky, M. 2014. Tensor decompositions for learning latent variable models. *The Journal of Machine Learning Research* 15(1):2773–2832.
- Arora, S.; Ge, R.; and Moitra, A. 2012. Learning topic models - going beyond SVD. In *FOCS*, 1–10.
- Bhatia, P.; Ji, Y.; and Eisenstein, J. 2015. Better document-level sentiment analysis from rst discourse parsing. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Blitzer, J.; Dredze, M.; and Pereira, F. 2007. Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification. In *Proceedings of the Association for Computational Linguistics (ACL)*, 440–447.
- Boyd, S.; Parikh, N.; Chu, E.; Peleato, B.; and Eckstein, J. 2011. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning* 3(1):1–122.
- Cohen, S. B.; Stratos, K.; Collins, M.; Foster, D. P.; and Ungar, L. 2014. Spectral learning of latent-variable PCFGs: Algorithms and sample complexity. *Journal of Machine Learning Research* 15:2399–2449.
- Hatzivassiloglou, V., and McKeown, K. R. 1997. Predicting the semantic orientation of adjectives. In *Proceedings of the Association for Computational Linguistics (ACL)*, 174–181.
- Hsu, D.; Kakade, S. M.; and Zhang, T. 2012. A spectral algorithm for learning hidden markov models. *Journal of Computer and System Sciences* 78(5):1460–1480.
- Kiritchenko, S.; Zhu, X.; and Mohammad, S. M. 2014. Sentiment analysis of short informal texts. *Journal of Artificial Intelligence Research* 50:723–762.
- Laver, M., and Garry, J. 2000. Estimating policy positions from political texts. *American Journal of Political Science* 619–634.
- Liu, B. 2015. *Sentiment Analysis: Mining Opinions, Sentiments, and Emotions*. Cambridge University Press.
- Maas, A. L.; Daly, R. E.; Pham, P. T.; Huang, D.; Ng, A. Y.; and Potts, C. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Madsen, R. E.; Kauchak, D.; and Elkan, C. 2005. Modeling word burstiness using the dirichlet distribution. In *Proceedings of the 22nd international conference on Machine learning*, 545–552. ACM.
- Minka, T. 2012. Estimating a dirichlet distribution. <http://research.microsoft.com/en-us/um/people/minka/papers/dirichlet/minka-dirichlet.pdf>.
- Molina-González, M. D.; Martínez-Cámara, E.; Martín-Valdivia, M.-T.; and Perea-Ortega, J. M. 2013. Semantic orientation for polarity classification in spanish reviews. *Expert Systems with Applications* 40(18):7250–7257.
- Pang, B., and Lee, L. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the Association for Computational Linguistics (ACL)*, 271–278.
- Pang, B., and Lee, L. 2008. Opinion mining and sentiment analysis. *Foundations and trends in information retrieval* 2(1-2):1–135.
- Pang, B.; Lee, L.; and Vaithyanathan, S. 2002. Thumbs up?: sentiment classification using machine learning techniques. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, 79–86.
- Polanyi, L., and Zaenen, A. 2006. Contextual valence shifters. In *Computing attitude and affect in text: Theory and applications*. Springer.
- Qiu, G.; Liu, B.; Bu, J.; and Chen, C. 2011. Opinion word expansion and target extraction through double propagation. *Computational linguistics* 37(1):9–27.
- Ribeiro, F. N.; Araújo, M.; Gonçalves, P.; Gonçalves, M. A.; and Benevenuto, F. 2016. Sentibench-a benchmark comparison of state-of-the-practice sentiment analysis methods. *EPJ Data Science* 5(1):1–29.
- Settles, B. 2011. Closing the loop: Fast, interactive semi-supervised annotation with queries on features and instances. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 1467–1478. Association for Computational Linguistics.
- Somasundaran, S.; Wiebe, J.; and Ruppenhofer, J. 2008. Discourse level opinion interpretation. In *Proceedings of the 22nd International Conference on Computational Linguistics-Volume 1*, 801–808. Association for Computational Linguistics.
- Taboada, M.; Brooke, J.; Tofiloski, M.; Voll, K.; and Stede, M. 2011. Lexicon-based methods for sentiment analysis. *Computational linguistics* 37(2):267–307.
- Tausczik, Y. R., and Pennebaker, J. W. 2010. The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of Language and Social Psychology* 29(1):24–54.
- Turney, P. 2002. Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews. In *Proceedings of the Association for Computational Linguistics (ACL)*, 417–424.
- Vilares, D.; Alonson, M. A.; and Gómez-Rodríguez, C. 2015. A syntactic approach for opinion mining on spanish reviews. *Natural Language Engineering* 21:139–163.
- Wilson, T.; Wiebe, J.; and Hoffmann, P. 2005. Recognizing contextual polarity in phrase-level sentiment analysis. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, 347–354.