

# JOINT TOPIC-DOCUMENT MODELING VIA LOW-DIMENSIONAL SPARSE MODELS

Kerui Min<sup>\*</sup> Yi Ma<sup>\*†</sup>

<sup>\*</sup>Dept. of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign

<sup>†</sup>Visual Computing Group, Microsoft Research Asia

E-mail: kmin12@illinois.edu, mayi@microsoft.com

## ABSTRACT

*Topic modeling* is a well-known approach for document analysis. In this paper, we propose a new model, and corresponding optimization algorithm for topic modeling. Experimental results on polarity classification demonstrate that the new model provides a more accurate characterization for document corpus, and archived higher classification accuracy compared to *Latent Dirichlet Allocation* (LDA).

**Index Terms**— topic modeling, non-negative matrix factorization

## 1. INTRODUCTION

Topic modeling is a well-known computational tool for finding thematic information from a given corpus. In recent years, this approach and its variants have been proved to be very useful in a large number of applications, including text and image classification, clustering, retrieval etc., see, *e.g.* the discussion paper [1] and references therein.

Traditionally, topic modeling was based on the assumption that each document (or more generally, each observed vector) is a combination of a small number of topic vectors [2, 3]. Based on this assumption, quite a few algorithms have been developed, trying to solve the topic vectors given a corpus [4, 3, 5, 6, 7]. As we will see in Section 3, by verifying this assumption on real datasets, we found that the traditional model is highly inaccurate. For example, a document often contains certain keywords (person's name, brand name, etc.) that cannot be explained by common topics. We instead propose a *joint topic-document model* to better model the structure within a document corpus.

## 2. RELATED WORK

The idea of representing a document by a linear combination of few “topics” can be track back to [4]. The theoretical computer science community studied the behavior of the above algorithm [2]. In machine learning, many probabilistic topic models were developed in recent years, *e.g.* *Latent Dirichlet Allocation* (LDA) [3] and many of its variant like *Correlated Topic Model* (CTM) [5]. Accompanying

with the develop of topic models, varying algorithms have been proposed, including Variational Bayesian [3], Gibbs sampling [6] and stochastic gradient descent [8]. Despite its empirical success, it is also shown that the inference of LDA model is NP-Hard [9].

The inference of topic models has an intimate relationship with *non-negative matrix factorization* (NMF), as the non-negativity factors can be associated with a probabilistic interpretation. Very recently, [7] studied how to use NMF for topic modeling, and provided an algorithm with provably approximation under a separability assumption. Many other local search based algorithms also exist, *e.g.* [10, 11].

## 3. OUR APPROACH

### 3.1. Issues with current approach

Topic modeling is based on the assumption that each document in a corpus is a convex combination of relatively small number of topic vectors. Therefore, the corpus, written as a matrix by stacking the document vectors together, has a low-rank structure. Mathematically, it assumes that the corpus  $D$  has the following structure:

$$D \approx BS, \quad (1)$$

where  $B \in \mathbb{R}^{m \times k}$  consisting of  $k$  topic vectors, acting as the basis of documents, and  $S \in \mathbb{R}^{k \times n}$  is the coefficient matrix,  $k \ll n$ . To obtain a probabilistic interpretation,  $B$  and  $S$  should satisfying the following condition:

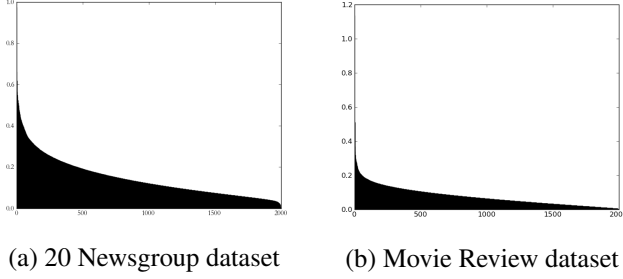
$$B \geq 0, \quad S \geq 0, \quad \|B_i\|_1 = 1 \text{ for } i \in [k], \quad (2)$$

where  $B_i$  is the  $i$ -th column of the basis matrix  $B$ . We note that for applications such as image analysis, the *sum-to-1* condition may not be necessary.

Although the above model is widely accepted, and has been used explicitly or implicitly in [4, 3, 5], this generative model is far from being accurate, empirically. To see this, we analyze the spectral distribution of two datasets: the 20 Newsgroup dataset<sup>1</sup> and the Movie Review dataset<sup>2</sup>. The

<sup>1</sup>Available at <http://qwone.com/~jason/20Newsgroups/>

<sup>2</sup>Available at <http://www.cs.cornell.edu/people/pabo/movie-review-data/>



**Fig. 1.** Singular value distribution of two typical corpora. Both exhibit long tail distribution.

Variance	30%	50%	70%	90%
Rank	100	244	490	1001

**Fig. 2.** The minimum ranks required to capture 30%, 50%, 70%, and 90% of the variance of the TF-IDF matrix of the 20 Newsgroup dataset.

Movie Review dataset consists of 2000 reviews from the IMDB database, and the 20 Newsgroup is a collection of around 20,000 posts across 20 newsgroups. We randomly select 2000 posts from 20 Newsgroup dataset, and employ the standard pre-process to obtain the TF-IDF matrix  $D$  for each dataset. The singular value distributions of the two datasets are shown in Fig 1. The singular value distributions of both datasets qualitatively exhibit very similar heavy-tail distribution, which implies that the data matrices don't have an accurate low-rank approximation. As shown in Fig 2, to capture 90% of the variance of  $D$ , we should choose  $k$ , the number of topics, to be at least 1001. Notice that the rank function only provides an lower-bound of non-negative rank, as we further require  $B \geq 0$  and  $S \geq 0$ .

Empirical studies such as [3, 8] often choose  $k$  to be 50 to 200. As a consequence, the basic model  $D \approx BS$  is no longer appropriate as there remains large amount of variation in the data left unexplained by the model. To address this issue, [12] proposed the *joint topic-document model*,  $D = L + E$ , where  $L$  is the low-rank matrix capturing the background, or topic information, and  $E$  is a sparse matrix, representing the document-specific keywords or keyphrases that cannot be explained by the (low-rank) topic model.

The *joint topic-document model* provides a more accurate model for text corpus. Empirically, however, we found that when we penalize  $L$  to be a relatively low-rank matrix, the document-specific counterpart  $E$  is not so sparse. In addition, *joint topic-document model* only gives the span of topic distribution. Very often it is necessary to obtain the distribution for each topic and the corresponding mixture-of-topic coefficients for tasks like document classification and clustering.

### 3.2. Noisy joint topic-document model

Based on the above observations and discussions, we would like to come up with a model that can describe more accurately the corpus data. In particular we would like to model the topics from a corpus as well as those document-specific keywords, and remaining terms that cannot be represented by either the background topic or the document-specific keywords. More specifically, we would like to model our observed data matrix  $D \in \mathbb{R}^{m \times n}$  as:

$$D = BS + E + N, \quad (3)$$

where  $B$  is topic basis (or topic vector set),  $E$  corresponds to the document-specific keywords, and  $N$  are the remaining noises. We call this model the *noisy joint topic-document model*. By allowing  $N$  to be reasonably large in variance, the document-specific  $E$  can be made sufficiently sparse. Motivated by the observation that each document should only contain few topics, we can penalize the mixture-of-topic coefficients  $S$  by  $\ell_1$ -norm. We penalize the sparse keywords part  $E$  by  $\ell_1$ -norm for the same reason.

Hence one can learn the above model from the data via solving the following optimization problem

$$\begin{aligned} \min_{B, S, E, N} \quad & \sum_i \|S_i\|_1^2 + \lambda \|E\|_1 + \frac{\tau}{2} \|N\|_F^2 \\ \text{subject to} \quad & D = BS + E + N, \quad B \geq 0, \\ & S \geq 0, \quad \|B_j\|_1 = 1 \text{ for } j \in [k]. \end{aligned} \quad (4)$$

The reason to choose  $\|S\|_{1,2}^2 \triangleq \sum_i \|S_i\|_1^2$  as the regularization instead of  $\|S\|_1$  is twofold:

1. By choosing  $\sum_i \|S_i\|_1^2$ , we enforce *each* document in the given corpus to be a linear combination of a few topics. The norm  $\|S\|_1$  enforces the corpus *as a whole* can be sparsely represented by  $B$ , but does not ensure representations for the documents are evenly sparse.
2. Computationally, the regularizer  $\sum_i \|S_i\|_1^2$  allows us to have an efficient algorithm for the optimization problem. We will discuss this issue in detail in the next section.

The above optimization problem also extends the *stable principal component pursuit* [13] by providing a specific generative model for the low-rank topic component.

## 4. ALGORITHM

In this section, we study how to solve problem (6) efficiently. It is easily seen that the optimization problem is non-trivial due to the non-linear constraint  $D = BS + E + N$ , where both  $B$  and  $S$  are unknown. Notice that the bilinear term also occurs in *non-negative matrix factorization*, *dictionary learning*, and *blind source separation*. Our problem formulation

(6) is closely related to dictionary learning, and the solution is inspired by non-negative matrix factorization.

First, we can take care of the non-linear constraint under the *Augmented Lagrangian Method* [14] framework to solve the problem iteratively. Specifically, we solve the following problem

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{S}, \mathbf{E}, \mathbf{N}} \quad & \|\mathbf{S}\|_{1,2}^2 + \lambda \|\mathbf{E}\|_1 + \frac{\tau}{2} \|\mathbf{N}\|_F^2 \\ & + \frac{\mu}{2} \|\mathbf{D} - \mathbf{BS} - \mathbf{E} - \mathbf{N}\|_F^2 \\ & + \langle \mathbf{Y}, \mathbf{D} - \mathbf{BS} - \mathbf{E} - \mathbf{N} \rangle \end{aligned} \quad (5)$$

$$\text{subject to } \mathbf{B} \geq \mathbf{0}, \mathbf{S} \geq \mathbf{0}, \|\mathbf{B}_j\|_1 = 1 \text{ for } j \in [k], \quad (6)$$

where  $\mathbf{Y}$  is the Lagrange multiplier. The remaining constraints are on  $(\mathbf{B}, \mathbf{S})$ . Once  $(\mathbf{B}, \mathbf{S})$  is fixed, the problem can be reduce to a typical  $\ell_1$  minimization problem using a fixed-point algorithm [15]. Alternatively, we can further fix  $\mathbf{N}$  and optimize  $\mathbf{E}$ , and then fix  $\mathbf{E}$  and optimize  $\mathbf{N}$ . Therefore, we solve the following three sub-problems iteratively.

**Sub-problem 1** Given  $(\mathbf{B}, \mathbf{S}, \mathbf{N}, \mathbf{Y}, \mu)$ , optimize  $\mathbf{E}$ . This is equivalent to solve the following unconstrained optimization program

$$\arg \min_{\mathbf{E}} \frac{\lambda}{\mu} \|\mathbf{E}\|_1 + \frac{1}{2} \left\| \mathbf{X}^{(1)} - \mathbf{E} \right\|_F^2, \quad (7)$$

where  $\mathbf{X}^{(1)}$  is given by  $\mathbf{X}^{(1)} \triangleq \mathbf{D} - \mathbf{BS} - \mathbf{N} + \frac{1}{\mu} \mathbf{Y}$ . It is well-known (see, e.g. [15]) that the unique optimal solution is given by the *shrinkage operator*  $\mathcal{S}_{\lambda\mu^{-1}}(\mathbf{X}^{(1)})$ , where

$$\mathcal{S}_\nu(\cdot) \triangleq \text{sgn}(\cdot) \max\{|\cdot| - \nu, 0\}, \quad (8)$$

and  $\text{sgn}(\cdot)$  is the sign function. We extend it to matrix domain by applying it element-wise to all entries.

**Sub-problem 2** Given  $(\mathbf{B}, \mathbf{S}, \mathbf{E}, \mathbf{Y}, \mu)$ , optimize  $\mathbf{N}$ . Similarly, this is equivalent to solve

$$\arg \min_{\mathbf{N}} \frac{\tau}{2} \|\mathbf{N}\|_F^2 + \frac{\mu}{2} \left\| \mathbf{X}^{(2)} - \mathbf{N} \right\|_F^2, \quad (9)$$

where  $\mathbf{X}^{(2)} \triangleq \mathbf{D} - \mathbf{BS} - \mathbf{E} + \frac{1}{\mu} \mathbf{Y}$ .

To obtain a closed-form solution, we take the derivative of the objective function with respect to  $\mathbf{N}$  and set it to zero.

$$\tau \mathbf{N} + \mu (\mathbf{N} - \mathbf{X}^{(2)}) = \mathbf{0}. \quad (10)$$

Hence, we get  $\mathbf{N} = \mu(\tau + \mu)^{-1} \mathbf{X}^{(2)}$ .

**Sub-problem 3** Given  $(\mathbf{E}, \mathbf{N}, \mathbf{Y}, \mu)$ , optimize  $(\mathbf{B}, \mathbf{S})$ . The original objective function is reduced to the following form

$$\begin{aligned} \arg \min_{\mathbf{B}, \mathbf{S}} \quad & \|\mathbf{S}\|_{1,2}^2 + \frac{\mu}{2} \left\| \mathbf{X}^{(3)} - \mathbf{BS} \right\|_F^2 \\ \text{subject to } & \mathbf{B} \geq \mathbf{0}, \mathbf{S} \geq \mathbf{0}, \|\mathbf{B}_j\|_1 = 1 \text{ for } j \in [k], \end{aligned}$$

where  $\mathbf{X}^{(3)} \triangleq \mathbf{D} - \mathbf{E} - \mathbf{N} + \frac{1}{\mu} \mathbf{Y}$ .

We note that unlike the previous two sub-problems where a closed-form optimal solution can be found, this problem is more difficult as the objective function contains the bilinear term  $\mathbf{BS}$ . This is a dictionary learning problem with non-negative constraints, or alternatively, a non-negative matrix factorization problem with sparse regularizer. To make the problem slightly simpler, we can solve the following problem instead.

$$\begin{aligned} \arg \min_{\mathbf{B}, \mathbf{S}} \quad & \|\mathbf{S}\|_{1,2}^2 + \gamma \|\mathbf{B}\|_F^2 + \frac{\mu}{2} \left\| \mathbf{X}^{(3)} - \mathbf{BS} \right\|_F^2 \\ \text{subject to } & \mathbf{B} \geq \mathbf{0}, \mathbf{S} \geq \mathbf{0}, \end{aligned} \quad (11)$$

where the weight  $\gamma$  is used to strike a balance between the norms of  $\mathbf{B}$  and  $\mathbf{S}$ . Let  $(\mathbf{B}', \mathbf{S}')$  be a solution of the above optimizations, and define  $\mathbf{T} \in \mathbb{R}^{k \times k}$  to be

$$\mathbf{T} = \text{diag} \left( \|\mathbf{B}_1\|_1^{-1}, \dots, \|\mathbf{B}_k\|_1^{-1} \right). \quad (12)$$

Then, the solution  $(\mathbf{B}'\mathbf{T}, \mathbf{T}^{-1}\mathbf{S}')$  satisfies the sum-to-1 constraints. Moreover,  $\mathbf{B}'\mathbf{T}$  has the same  $\ell_0$ -norm as  $\mathbf{B}'$ , so the modified solution is sparse as long as the original  $\mathbf{B}'$  is sparse.

The problem (11) was formulated and studied in [11], where an efficient algorithm, called SNMF, was proposed based on non-negative least squares. Empirically, we find that SNMF provides sufficiently good solutions for sub-problem 3, and the algorithm is insensitive to the choice of  $\gamma$ . Therefore, we fix  $\gamma$  to be a constant throughout our experiments.

**Update Lagrange Multiplier** The final step in each iteration is to update the Lagrange Multiplier  $\mathbf{Y}$ . Let  $\mathbf{Y} \leftarrow \mathbf{Y} + \mu(\mathbf{D} - \mathbf{BS} - \mathbf{E} - \mathbf{N})$ , and update  $\mu$  to be  $\rho\mu$  for a constant  $\rho > 1$ .

## 5. EXPERIMENTS

### 5.1. Experimental settings

In this section, we evaluate the proposed topic-document modeling algorithm for polarity classification. We use the Movie Review dataset, which contains 1000 positive reviews and 1000 negative reviews [16]. Therefore, it is a binary classification task. To make a fair comparison, only unigram feature (terms) is used, as it is difficult for traditional topic models like *Latent Dirichlet Allocation* to incorporate other features such as *part-of-speech* tags, or bigram features. Since terms with low document frequency (the number of documents it occurs within a corpus) have little contribution to the topic words, we only keep the top  $m$  terms with the highest document frequencies, and construct the TF-IDF matrices accordingly.

*Support vector machine* (SVM) is used as the classifier for all the subsequent tasks, since it has achieved the highest performance in [16]. For simplicity, we use the linear SVM. The performance accuracy will be evaluated using 5-fold cross-validation using the 2000 samples.

Our	LDA	Our	LDA	Our	LDA
elizabeth (0.12723)	elizabeth (0.00510)	wedding (0.05833)	sandler (0.00791)	batman (0.06913)	effects (0.00797)
queen (0.04298)	foster (0.00449)	singer (0.03481)	wedding (0.00721)	robin (0.02903)	star (0.00713)
england (0.01782)	giles (0.00277)	julia (0.03362)	dvd (0.00645)	mr (0.01424)	special (0.00637)
mary (0.01290)	game (0.00264)	adam (0.01425)	series (0.00410)	joel (0.01008)	batman (0.00401)
rush (0.01159)	anna (0.00263)	romantic (0.01236)	snake (0.00374)	comic (0.00615)	wars (0.00387)
love (0.01112)	fight (0.00256)	comedy (0.00833)	singer (0.00366)	wayne (0.00591)	series (0.00326)
court (0.00971)	video (0.00253)	billy (0.00788)	disc (0.00347)	city (0.00583)	planet (0.00265)
power (0.00900)	eddie (0.00246)	funny (0.00554)	fans (0.00333)	series (0.00557)	trek (0.00245)
country (0.00864)	king (0.00240)	happy (0.00548)	anaconda (0.00323)	villains (0.00529)	earth (0.00236)
political (0.00839)	mortal (0.00237)	steve (0.00546)	x-files (0.00304)	tim (0.00528)	alien (0.00233)

**Fig. 3.** Comparison of topics learned using our method with LDA for similar topics ( $k = 30$ ). The weight of each term is given in the bracket.

# of topics ( $k$ )	LDA	Our method
30	57.2%	<b>68.2%</b>
50	61.3%	<b>70.0%</b>

**Fig. 4.** The classification accuracy on the Movie Review dataset, using the distribution of topics information.

## 5.2. Verify the learned topic

In this experiment, we compare the polarity classification performance of our topic modeling algorithm with LDA using the distribution of topics information, *i.e.* we use  $\mathbf{S}_i \in \mathbb{R}^k$  to represent the  $i$ -th document. For LDA, we use the online LDA algorithm from [8] to learn the topics with default parameters. To ensure the convergence of online LDA, we pass the corpus 200 times. The result is given in Fig 4. Geometrically, the new features are given by projecting the original high-dimensional features to a low-dimensional simplex. We can see that the  $k$ -dimensional vector learned by our method significantly outperformed the LDA counterpart, implying that the learned information is semantically more meaningful. In Fig 3, we provide examples of the learned topics compared to LDA.

## 5.3. Classification using the structure information

The last experiments shows that the decomposed triple  $(\mathbf{BS}, \mathbf{E}, \mathbf{N})$  is effective for classification. Recall that our algorithm decompose the corpus  $\mathbf{D}$  into  $\mathbf{BS} + \mathbf{E} + \mathbf{N}$ , where  $\mathbf{BS}$  contains the topic information,  $\mathbf{E}$  represents those document-specific keywords, and  $\mathbf{N}$  is the remaining residues. Our intuition is that the above three terms all contain useful information and should be weighed differently to obtain better performance. To validate this hypothesis, we simply concatenate the decomposed parts together, *i.e.* we use the vector  $[(\mathbf{BS}_i)^T, \mathbf{E}_i^T, \mathbf{N}_i^T] \in \mathbb{R}^{3m}$  as the feature of the  $i$ -th document, instead of  $\mathbf{D}_i \in \mathbb{R}^m$ . If the structure of the model does not help extract any additional information, we should expect a performance drop, since the new dimension-

# of unigram feature $m$	$\mathbf{D}_i$	$[(\mathbf{BS}_i)^T, \mathbf{E}_i^T, \mathbf{N}_i^T]$
$m = 2000$	85.2%	<b>87.0%</b>
$m = 3000$	85.5%	<b>87.2%</b>

**Fig. 5.** Classification accuracy on the Moview Review dataset using the original and new feature space.

ality of the new feature space is three times large. Or more specifically, based on the classical statistical learning theory (see, *e.g.* [17]), we know that

$$err_{\text{true}}(h) \leq err_{\text{train}}(h) + \sqrt{\frac{VC(H)(\ln \frac{2n}{VC(H)} + 1) + \ln \frac{4}{\delta}}{n}}, \quad (13)$$

whree  $VC(H)$  indicates the VC-dimension of the hypothesis space  $H$ . For linear classifier, the original feature space has  $VC(H) = m + 1$ , and the new features space has  $VC(H') = 3m + 1$ . For sufficiently large  $n$ , the second term of the right-hand side of Eq (13) is about 1.73 times larger. The actual performance of the two approaches is given in Fig 5. Clearly, despite the fact that the new feature space is more prone to over-fitting, the structure information does help to improve the accuracy of classification.

Interestingly, although we have not done any special engineering towards this task, we found that the performance is on par with the best performance reported in [18], where a much more sophisticated algorithm based on sentence-level analysis and a subjectivity dataset was designed specifically for the sentimental analysis.

## 6. CONCLUSION

In this paper, to address issues of previous topic models, we proposed the *noisy joint topic-document model* and provide an efficient iterative algorithm for topic modeling based on it. The model is generic, we believe it can be used for document classification, clustering, and beyond.

## 7. REFERENCES

- [1] D. Blei, “Probabilistic topic models,” in *Communications of the ACM*, 2012, vol. 55, pp. 77–84.
- [2] C. H. Papadimitriou, H. Tamaki, P. Raghavan, and S. Vempala, “Latent semantic indexing: a probabilistic analysis,” in *Proceedings of the seventeenth ACM symposium on Principles of database systems*, 1998, pp. 159–168.
- [3] D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent dirichlet allocation,” in *J. Mach. Learn. Res.*, 2003, number 3, pp. 993–1022.
- [4] S. T. Dumais, G. W. Furnas, T. K. Landauer, S. Deerwester, and R. Harshman, “Using latent semantic analysis to improve access to textual information,” in *Proceedings of CHI*, 1988, pp. 281–286.
- [5] D. M. Blei and J. D. Lafferty, “Correlated topic models,” in *NIPS*, 2005.
- [6] I. Porteous, D. Newman, A. Ihler, A. Asuncion, P. Smyth, and M. Welling, “Fast collapsed gibbs sampling for latent dirichlet allocation,” in *Proc. of the 14th ACM SIGKDD*, 2009, pp. 569–577.
- [7] S. Arora, R. Ge, and A. Moitra, “Learning topic models - going beyond svd,” in *CoRR abs/1204.1956*, 2012.
- [8] M. D. Hoffman, D. M. Blei, and F. R. Bach, “Online learning for latent dirichlet allocation,” in *NIPS*, 2010, pp. 856–864.
- [9] D. Sontag and D. Roy, “Complexity of inference in latent dirichlet allocation,” in *NIPS*, 2011.
- [10] C.-J. Lin, “Projected gradient methods for non-negative matrix factorization,” in *Neural Computation*, 2007, vol. 19, pp. 2756–2779.
- [11] H. Kim and H. Park, “Nonnegative matrix factorization based on alternating nonnegativity constrained least squares and active set method,” in *SIAM J. Matrix Analysis Applications*, 2008, vol. 30, pp. 713–730.
- [12] K. Min, Z. Zhang, J. Wright, and Y. Ma, “Decomposing background topics from keywords by principal component pursuit,” in *CIKM*, 2010, pp. 269–278.
- [13] Z. Zhou, X. Li, J. Wright, E. J. Candes, and Y. Ma, “Stable principal component pursuit,” in *ISIT*, 2010, pp. 1518–1522.
- [14] D. Bertsekas, “Constrained optimization and lagrange multiplier method,” in *Academic Press*, 1982.
- [15] E. T. Hale, W. Yin, and Y. Zhang, “Fixed-point continuation for  $\ell_1$ -minimization: Methodology and convergence,” in *SIAM Journal on Optimization*, 2008, vol. 19, pp. 1107–1130.
- [16] Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan, “Thumbs up? sentiment classification using machine learning techniques,” in *Proceedings of EMNLP*, 2002, pp. 79–86.
- [17] V. N. Vapnik, *The nature of statistical learning theory*, Springer-Verlag New York, Inc., New York, NY, USA, 1995.
- [18] Bo Pang and Lillian Lee, “A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts,” in *Proceedings of ACL*, 2004, pp. 271–278.