



香港城市大學
City University of Hong Kong

專業 創新 胸懷全球
Professional · Creative
For The World

CityU Scholars

Sentiment analysis with covariate-assisted word embeddings

Xu, Shirong; Dai, Ben; Wang, Junhui

Published in:

Electronic Journal of Statistics

Published: 01/01/2021

Document Version:

Final Published version, also known as Publisher's PDF, Publisher's Final version or Version of Record

License:

CC BY

Publication record in CityU Scholars:

[Go to record](#)

Published version (DOI):

[10.1214/21-EJS1854](https://doi.org/10.1214/21-EJS1854)

Publication details:

Xu, S., Dai, B., & Wang, J. (2021). Sentiment analysis with covariate-assisted word embeddings. *Electronic Journal of Statistics*, 15(1), 3015-3039. <https://doi.org/10.1214/21-EJS1854>

Citing this paper

Please note that where the full-text provided on CityU Scholars is the Post-print version (also known as Accepted Author Manuscript, Peer-reviewed or Author Final version), it may differ from the Final Published version. When citing, ensure that you check and use the publisher's definitive version for pagination and other details.

General rights

Copyright for the publications made accessible via the CityU Scholars portal is retained by the author(s) and/or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights. Users may not further distribute the material or use it for any profit-making activity or commercial gain.

Publisher permission

Permission for previously published items are in accordance with publisher's copyright policies sourced from the SHERPA RoMEO database. Links to full text versions (either Published or Post-print) are only available if corresponding publishers allow open access.

Take down policy

Contact lbscholars@cityu.edu.hk if you believe that this document breaches copyright and provide us with details. We will remove access to the work immediately and investigate your claim.

Sentiment analysis with covariate-assisted word embeddings*

Shirong Xu

*School of Data Science, City University of Hong Kong,
Kowloon Tong, Hong Kong*
e-mail: shirongxu2-c@my.cityu.edu.hk

Ben Dai

*School of Statistics, University of Minnesota,
Minneapolis, MN 55455, USA*
e-mail: bdai@umn.edu

Junhui Wang

*School of Data Science, City University of Hong Kong,
Kowloon Tong, Hong Kong*
e-mail: j.h.wang@cityu.edu.hk

Abstract: Sentiment analysis measures inclination of textual documents, aiming to extract and quantify their subjective sentiment polarity. In literature, most sentiment analysis methods first numericalize textual documents through certain word embeddings framework, and then formulate sentiment analysis as an ordinal regression or classification task. Yet it is often ignored that different people may have different preference of wording, and thus a uniform word embeddings often leads to suboptimal performance. In this article, to accommodate the heterogeneity among individual persons, we propose a covariate-assisted word embeddings in a margin-based ordinal regression framework, where covariates are incorporated through scaling factors to adjust the word embeddings. Moreover, we employ a block-wise coordinate descent scheme to tackle the resultant large-scale optimization task, and establish theoretical results to quantify the asymptotic behavior of the proposed method, guaranteeing its fast convergence rate in terms of prediction accuracy. Finally, we demonstrate the advantages of the proposed method over its competitors in both the Yelp Challenge dataset and synthetic datasets.

MSC2020 subject classifications: Primary 62H30.

Keywords and phrases: Personalized prediction, sentiment analysis, word embeddings, unstructured data, ordinal regression.

Received October 2020.

Contents

1	Introduction	3016
---	------------------------	------

*This work is supported in part by HK RGC grants GRF-11303918, GRF-11300919 and GRF-11304520.

2	Proposed methodology	3017
2.1	Preambles	3017
2.2	Covariate-assisted word embeddings	3019
2.3	Scalable computation	3020
3	Theory	3022
4	Numerical experiments	3024
4.1	Yelp challenge	3024
4.2	Simulations	3027
5	Summary	3029
	Acknowledgments	3029
	Appendix	3029
	Supplementary Material	3037
	References	3037

1. Introduction

Unstructured text data has become increasingly important in recent years, due to the fast advancement of information technology and evolution of information storage. It typically arises from text-heavy documents, including customer reviews, news stories, or online twits. One of the central tasks of text data analysis is to extract subjective sentiment polarity of the textual documents, which has been an essential component in modern business analytics and political surveys [1, 2].

In literature, most sentiment analysis methods first convert textual documents into numerical vectors, and then formulate it as a classification task, where sentiment levels are treated as binary or ordinal responses [3, 4, 5]. The numericalization is often done by using the bag-of-words framework [6], where word presences or frequencies in the textual documents are extracted as the numerical features. The bag-of-words framework is interpretable and easy to implement, but it fails to capture the relationship among meaningful words. Recently, embedding technique [7] has drawn significant interests from both statistics and machine learning communities for its flexibility and interpretability in representing textual documents, including Word2Vec [8] and Global Vectors for Word Representation (GloVe) [9]. The key idea of word embeddings is to embed each word into a low-dimensional vector space so that the corresponding vectors of relevant words are close in the embedded space. A number of embeddings schemes have been proposed from various perspectives, in order to obtain a uniform word embeddings for all individual persons to facilitate the subsequent text data analysis.

A uniform word embeddings is simple in nature but also suffers from some intrinsic limitations, due to the fact that different people may have different preferences of wording. For example, the word “interesting” can be used to express neutral or even negative sentiment level by people who tend not to use negative words to show politeness. Also, it appears very common that people like to use sarcastic expressions on internet. One review in Yelp dataset says

“I feel so excited to check out” to express the extremely negative sentiment without using any negative words. Analyzing such textual statements can be easily misled by the presence of positive words with strong polarity [10, 11]. In literature, difference in wording across genders, ages, educational background and political background has been widely reported [1, 12, 13]. It is thus natural to consider adaptive word embeddings to capture the heterogeneity among individual persons so that their preferences of wording can be incorporated to improve the prediction accuracy. Yet, only a few attempts have been made in literature, including time-varying word embeddings [14] and topic-adaptive word embeddings [15, 16].

In this paper, we propose a sentiment analysis method based on a novel covariate-assisted word embeddings, which integrates covariates into an ordinal regression framework [17] to refine word embeddings for better prediction accuracy. Specifically, a sentiment lexicon and the corresponding word embeddings will be employed to construct covariate-adjusted representation of each textual document. For each covariate level, an adjusting factor is introduced to scale the original word embeddings, which quantifies the deviation of semantics from the pre-trained word embeddings. Furthermore, we also develop a scalable block-wise coordinate descent algorithm to tackle the resultant large-scale optimization task. Theoretically, the asymptotic convergence rate of the proposed method is established in terms of sample size, sentiment levels, covariate levels, lexicon size.

The rest of the paper is organized as follows. Section 2 presents the proposed covariate-assisted word embeddings in an ordinal regression framework for sentiment analysis, as well as the block-wise coordinate descent algorithm. Section 3 establishes the asymptotic results for the proposed method assuring its fast convergence rate under several situations. Section 4 conducts a simulation study to examine the numerical performance of the proposed method in various synthetic datasets and applies the proposed method to analyze the Yelp challenge dataset. A brief summary is given in Section 5, and the Appendix contains the technical proofs.

2. Proposed methodology

2.1. Preambles

In sentiment analysis, a training dataset consists of $\{(\mathbf{t}_{ij}, y_{ij}); i = 1, \dots, u, j = 1, \dots, N_i\}$, where \mathbf{t}_{ij} is the j -th textual document made by the i -th person, and $y_{ij} \in \{1, \dots, K\}$ indicates its sentiment level with ordering $1 \prec 2 \prec \dots \prec K$, where \prec denotes less positive in terms of sentiment level. The primary goal of sentiment analysis is to construct a decision function $\phi(\mathbf{t}_{ij})$ to accurately predict the sentiment level of \mathbf{t}_{ij} , so that the disagreement between $\phi(\mathbf{t}_{ij})$ and y_{ij} can be minimized.

Various disagreement metrics for ordinal regression have been considered in literature [18], including mean absolute error(MAE), mean zero-one error

(MZOE) and mean square error(MSE). However, neither MZOE nor MSE is originally designed for ordinal regression, and both metrics have their own limitations for analyzing ordinal data. Particularly, MZOE fails to take the ordinality into account, which is particularly undesirable in sentiment analysis, where mis-classifying a positive review as neutral is less severe than as negative, whereas these two type of misclassifications are treated equally in MZOE. As for MSE, it regards ordinal response as continuous, leading to unnecessary bias, especially when ordinal responses are only encoded to reflect the ordering but imply no elaboration of the difference among the ordered values. By contrast, MAE appears to be a reasonable choice for ordinal regression, and widely used in literature [17]. It can be written as

$$\text{MAE}(\phi) = E(|y - \phi(\mathbf{t})|) = \sum_{k=1}^{K-1} E\left(I(\text{sgn}(y - k)\text{sgn}(\phi(\mathbf{t}) - k) \leq 0)\right), \quad (2.1)$$

where $I(\cdot)$ is an indicator function and $\text{sgn}(x) = 1$ when $x > 0$, and -1 otherwise. It follows immediately from (2.1) that minimizing MAE is equivalently transformed into solving $K - 1$ binary classification problems, where $\text{sgn}(y - k)$ can be treated as the binary class label and $\text{sgn}(\phi(\mathbf{t}) - k)$ denotes the corresponding classification decision function.

Instead of estimating $\phi(\mathbf{t})$ directly, it is common to introduce $K - 1$ functions with $f_{K-1}(\mathbf{t}) \leq \dots \leq f_1(\mathbf{t})$, and set $\phi(\mathbf{t}) = \min\{k : f_k(\mathbf{t}) \leq 0\}$. Then the estimation of ϕ is converted to estimating $\mathbf{f} = (f_1, \dots, f_{K-1})$, and

$$\text{MAE}(\phi) = \text{MAE}(\mathbf{f}) = \sum_{k=1}^{K-1} E\left(I(\text{sgn}(y - k)f_k(\mathbf{t}) \leq 0)\right). \quad (2.2)$$

The indicator function in (2.2) is computationally intractable for optimization, thus we replace it by some surrogate margin losses. Specifically, an empirical version of (2.2) with a surrogate loss and regularization term can then be constructed to estimate \mathbf{f} ,

$$\min_{\mathbf{f}} \frac{1}{N} \sum_{i=1}^n \sum_{j=1}^{N_i} \sum_{k=1}^{K-1} V(\text{sgn}(y_{ij} - k)f_k(\mathbf{t}_{ij})) + \lambda J(\mathbf{f}) \quad (2.3)$$

subject to $f_{K-1}(\mathbf{t}_{ij}) \leq f_{K-2}(\mathbf{t}_{ij}) \leq \dots \leq f_1(\mathbf{t}_{ij}); i = 1, \dots, n, j = 1, \dots, N_i$,

where $N = \sum_{i=1}^n N_i$, $V(z)$ is a surrogate margin loss function non-decreasing in z , $J(\mathbf{f})$ is a regularization term, and λ is a tuning parameter. Here $V(z)$ can take various forms. For instance, $V(\cdot)$ can be the hinge loss $V(u) = (1 - u)_+$ [19], the ψ -loss $V(u) = \min((1 - u)_+, 1)$ [20], or the logistic loss $V(u) = 1/(1 + \exp(-u))$ [21]. It is interesting to note that setting V as the hinge loss or the logistic loss in (2.3) resembles the ALL-threshold method [22] for ordinal regression.

To facilitate the modelling of \mathbf{f} , textual documents need to be pre-processed into numerical vectors. In literature, primitive approaches extract word presence or frequency in \mathbf{t} as predictors under the bag-of-words framework [6, 23].

Recently, more informative word embeddings frameworks have been developed, such as Word2Vec and GloVe [8, 9]. In particular, Word2Vec learns the word representation via a three-layer neural network, which assumes that words with similar linguistic meaning should be close in textual documents, resulting in frequent co-occurrences within a fixed-sized context window.

Specifically, let $\mathcal{D} = \{\omega_1, \omega_2, \dots, \omega_d\}$ be a lexicon of sentiment words and $\mathbf{E} \in \mathbb{R}^{p \times d}$ be an embedding matrix, where p is the dimension of the embedded space and each column denotes the embedding of the corresponding word in \mathcal{D} . We further define $\mathbf{B}(\mathbf{t}) = (b_1, \dots, b_d)^T$ to be the frequency vector of \mathbf{t} based on \mathcal{D} . Then $\mathbf{EB}(\mathbf{t})$ is the averaged embeddings of words appearing in \mathbf{t} , which can be viewed as the representation of \mathbf{t} in the embedded space, and the sentiment function f_k can be formulated as

$$f_k(\mathbf{t}_{ij}) = \boldsymbol{\beta}^T \mathbf{EB}(\mathbf{t}_{ij}) + \beta_{0,k}, \quad (2.4)$$

where $\boldsymbol{\beta} \in \mathbb{R}^p$, and $\beta_{0,K-1} \leq \dots \leq \beta_{0,1}$. Clearly, f_1, \dots, f_{K-1} are parallel, and their ordering is inherent in $\beta_{0,k}; k = 1, \dots, K-1$. Such structures have been commonly used in literature to enforce ordering among multiple functions [24, 25, 26].

2.2. Covariate-assisted word embeddings

In many scenarios, some covariates $\mathbf{x}_{ij} = (x_{ij1}, x_{ij2}, \dots, x_{ijL})^T$ are also available for each observation $(\mathbf{t}_{ij}, y_{ij})$, where $x_{ijl} \in \{1, \dots, m_l\}$ denotes the l -th covariate and m_l denotes its number of distinct levels. To incorporate available covariates in the word embeddings, the proposed covariate-assisted sentiment function f_k can be formulated as

$$f_k(\mathbf{t}_{ij}, \mathbf{x}_{ij}) = \boldsymbol{\beta}^T \mathbf{E}(\mathbf{B}(\mathbf{t}_{ij}) \circ \mathbf{w}_{x_{ij1}}^{(1)} \circ \dots \circ \mathbf{w}_{x_{ijL}}^{(L)}) + \beta_{0,k}(\mathbf{x}_{ij}), \quad (2.5)$$

where \circ denotes the entry-wise product, $\mathbf{w}_{x_{ijl}}^{(l)} \in \mathbb{R}^d$ denotes the parameter vector corresponding to x_{ijl} , and the intercept $\beta_{0,k}(\mathbf{x})$ is allowed to vary with \mathbf{x} . Note that both \mathbf{E} and \mathbf{B} are pre-trained or pre-specified, and only $\boldsymbol{\beta}$, $\beta_{0,k}$ and $\mathbf{w}_{x_{ijl}}^{(l)}$ are the unknown parameters in (2.5) that need to be estimated. Particularly, $\mathbf{w}_{x_{ijl}}^{(l)}$ serves the purpose of adjusting $\mathbf{B}(\mathbf{t}_{ij})$ in $f_k(\mathbf{t}_{ij}, \mathbf{x}_{ij})$, and the varying intercept can be viewed as baseline in predicting sentiment for each covariate level. Furthermore, although $\boldsymbol{\beta}$ and $\mathbf{w}_{x_{ijl}}^{(l)}$ may not be identifiable, they contribute to $f_k(\mathbf{t}_{ij}, \mathbf{x}_{ij})$ only through their product, and thus does not affect the predictability of the proposed method. If interpretability is also of interest, one may fix $\|\boldsymbol{\beta}\| = 1$ to avoid the non-identifiability, which only requires an additional normalization step of $\boldsymbol{\beta}$. Additionally, the proposed method in (2.5) is mainly designed for binary and categorical covariates, to which a direct extension to continuous covariates is to divide the domain of covariates into exclusive subsets, which are then treated as distinct categorical levels.

Let $\bar{\mathbf{w}}_{\mathbf{x}_{ij}} = \mathbf{w}_{x_{ij1}}^{(1)} \circ \dots \circ \mathbf{w}_{x_{ijL}}^{(L)}$ be the overall adjusting effect, then $f_k(\mathbf{t}_{ij}, \mathbf{x}_{ij})$ can be rewritten as

$$f_k(\mathbf{t}_{ij}, \mathbf{x}_{ij}) = \boldsymbol{\beta}^T [\mathbf{E} \circ (\mathbf{1}_p \otimes \bar{\mathbf{w}}_{\mathbf{x}_{ij}}^T)] \mathbf{B}(\mathbf{t}_{ij}) + \beta_{0,k}(\mathbf{x}_{ij}), \quad (2.6)$$

where \otimes denotes the Kronecker product. This formulation leads to the proposed covariate-assisted word embeddings, where $\bar{\mathbf{w}}_{ij}$ calibrates the embedding matrix \mathbf{E} by multiplying its embedding vector with a scaling factor. It allows for a refined word embedding by incorporating the available covariates, which is in sharp contrast to the uniform word embeddings in literature. For example, a positive word can be used to express negative sentiment when its embedding vector is multiplied by a negative scalar. This flexibility is particularly attractive when analyzing sarcastic statements on internet. Additionally, the overall adjusting effect $\bar{\mathbf{w}}_{ij}$ of documents issued by the same person will be close since they share common covariates, implying similarity among wordings of different textual documents by the same person. It also allows certain similarity among wording of different persons depending on the level of their common covariates.

With the modelling of f_k in (2.6), the proposed method can be organized as

$$\begin{aligned} \min_{\mathbf{W}, \boldsymbol{\beta}, \beta_0} \quad & \frac{1}{N} \sum_{i=1}^u \sum_{j=1}^{N_i} \sum_{k=1}^{K-1} V(\text{sign}(y_{ij} - k)(\boldsymbol{\beta}^T \mathbf{E}(\mathbf{B}(\mathbf{t}_{ij}) \circ \bar{\mathbf{w}}_{\mathbf{x}_{ij}}) + \beta_{0,k}(\mathbf{x}_{ij}))) \\ & + \lambda_1 J(\boldsymbol{\beta}) + \lambda_2 J(\mathbf{W}) \end{aligned} \quad (2.7)$$

$$\text{subject to } \beta_{0,K-1}(\mathbf{x}_{ij}) \leq \beta_{0,K-2}(\mathbf{x}_{ij}) \leq \dots \leq \beta_{0,1}(\mathbf{x}_{ij}) \text{ for all } \mathbf{x}_{ij},$$

where $\mathbf{W} = [\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \dots, \mathbf{W}^{(L)}]$, and $\mathbf{W}^{(l)} = [\mathbf{w}_1^{(l)}, \dots, \mathbf{w}_{m_l}^{(l)}]$ denotes the adjusting matrix of l -th categorical covariate. Here the number of parameters of $\boldsymbol{\beta}$ and \mathbf{W} are p and $d \sum_{l=1}^L m_l$, respectively. To control the complexity of sentiment functions, $J(\boldsymbol{\beta})$ and $J(\mathbf{W})$ can be any regularization term. In the sequel, we illustrate the proposed method by setting $V(\cdot)$ is the hinge loss, $J(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|_2^2$, and $J(\mathbf{W}) = \|\mathbf{W}\|_F^2$.

Note that the proposed method in (2.7) consists of two tuning parameters λ_1 and λ_2 , and Lemma 1 shows they play a similar role in (2.7) and thus significantly simplifies the tuning process for λ_1 and λ_2 .

Lemma 1. *The solution to (2.7) remains the same as long as $\lambda_1 \lambda_2^L$ stays the same.*

Lemma 1 implies that the optimization task in (2.7) with tuning parameters (λ_1, λ_2) has the same solution as that with (λ, λ) satisfying $\lambda^{L+1} = \lambda_1 \lambda_2^L$. Therefore, we simplify the cost function in (2.7) by setting $\lambda_1 = \lambda_2 = \lambda$ in the sequel. The proof of Lemma 1 and all other technical proofs are provided in a supplementary file [27].

2.3. Scalable computation

Note that optimization task in (2.7) is a bi-convex optimization problem with respect to $\boldsymbol{\beta}$ and $\bar{\mathbf{w}}_{\mathbf{x}_{ij}}$, and hence we employ a block-wise coordinate descent

algorithm to update β and $\bar{\mathbf{w}}_{\mathbf{x}_{ij}}$ sequentially. By introducing a slack variable ξ_{ijk} , (2.7) can be reformulated as

$$\begin{aligned} \min_{\mathbf{W}, \beta, \beta_0, \xi_{ijk}} \quad & \sum_{i,j,k} \xi_{ijk} + \lambda \|\beta\|_2^2 + \lambda \|\mathbf{W}\|_F^2 \quad (2.8) \\ \text{subject to} \quad & \text{sgn}(y_{ij} - k) (\beta^T \mathbf{E}(\mathbf{B}(\mathbf{t}_{ij}) \circ \bar{\mathbf{w}}_{\mathbf{x}_{ij}}) + \beta_{0,k}(\mathbf{x}_{ij})) \geq 1 - \xi_{ijk}, \quad \xi_{ijk} \geq 0, \\ & \beta_{0,K-1}(\mathbf{x}_{ij}) \leq \beta_{0,K-2}(\mathbf{x}_{ij}) \leq \dots \leq \beta_{0,1}(\mathbf{x}_{ij}). \end{aligned}$$

We then break (2.8) into multiple sub-tasks, and update β , \mathbf{W} and β_0 alternatively. Specifically, when \mathbf{W} and β_0 are fixed, β can be updated by solving

$$\begin{aligned} \min_{\beta, \xi_{ijk}} \quad & \sum_{i,j,k} \xi_{ijk} + \lambda \|\beta\|_2^2 \quad (2.9) \\ \text{subject to} \quad & \text{sgn}(y_{ij} - k) (\beta^T \mathbf{E}(\mathbf{B}(\mathbf{t}_{ij}) \circ \bar{\mathbf{w}}_{\mathbf{x}_{ij}}) + \beta_{0,k}(\mathbf{x}_{ij})) \geq 1 - \xi_{ijk}, \quad \xi_{ijk} \geq 0. \end{aligned}$$

Note that the optimization task in (2.9) resembles linear support vector machine (SVM) in nature except additional varying intercepts, which can be efficiently solved by `Liblinear` [28]. Therefore, we develop a similar optimization scheme based on dual coordinate descent method as in `Liblinear` to solve (2.9), which is named as driftSVM and available in Python package `VarSVM`.

When β and β_0 are fixed, our strategy is to use the back-fitting scheme to update $\mathbf{W}^{(l)}$; $l = 1, \dots, L$ sequentially. Particularly, $\mathbf{W}^{(l)}$ can be updated by solving

$$\begin{aligned} \min_{\mathbf{W}^{(l)}, \xi_{ijk}} \quad & \sum_{q=1}^{m_l} \sum_{\{\mathbf{x}_{ij}: x_{ijl}=q\}} \sum_k \xi_{ijk} + \lambda \sum_{q=1}^{m_l} \|\mathbf{w}_q^{(l)}\|_2^2 \\ \text{subject to} \quad & \text{sgn}(y_{ij} - k) ((\mathbf{w}_q^{(l)})^T \mathbf{B}_{ij}^{-\mathbf{w}_{x_{ijl}}} \mathbf{E}^T \beta + \beta_{0,k}(\mathbf{x}_{ij})) \geq 1 - \xi_{ijk}, \quad \xi_{ijk} \geq 0, \end{aligned}$$

where $\mathbf{B}_{ij}^{-\mathbf{w}_{x_{ijl}}} = \text{diag}\{\mathbf{B}(\mathbf{t}_{ij}) \circ \mathbf{w}_{x_{ij1}}^{(1)} \circ \dots \circ \mathbf{w}_{x_{ij,l-1}}^{(l-1)} \circ \mathbf{w}_{x_{ij,l+1}}^{(l+1)} \circ \dots \circ \mathbf{w}_{x_{ijL}}^{(L)}\}$ is a diagonal matrix. Furthermore, $\mathbf{w}_q^{(l)}$ can be optimized in a parallel fashion. That is, each $\mathbf{w}_q^{(l)}$, $q = 1, \dots, m_l$ can be updated by solving

$$\begin{aligned} \min_{\mathbf{w}_q^{(l)}, \xi_{ijk}} \quad & \sum_{\{\mathbf{x}_{ij}: x_{ijl}=q\}} \sum_k \xi_{ijk} + \lambda \|\mathbf{w}_q^{(l)}\|_2^2 \quad (2.10) \\ \text{subject to} \quad & \text{sgn}(y_{ij} - k) ((\mathbf{w}_q^{(l)})^T \mathbf{B}_{ij}^{-\mathbf{w}_q^{(l)}} \mathbf{E}^T \beta + \beta_{0,k}(\mathbf{x}_{ij})) \geq 1 - \xi_{ijk}, \quad \xi_{ijk} \geq 0. \end{aligned}$$

The optimization task in (2.10) is exactly the same as (2.9), and hence can be solved by following identical optimization scheme.

When β and \mathbf{W} are fixed, β_0 can be updated by solving

$$\begin{aligned} \min_{\beta_0, \xi_{ijk}} \quad & \sum_{i,j,k} \xi_{ijk} \quad (2.11) \\ \text{subject to} \quad & \text{sgn}(y_{ij} - k) \beta_{0,k}(\mathbf{x}_{ij}) + \xi_{ijk} \geq 1 - \text{sgn}(y_{ij} - k) (\beta^T \mathbf{E}(\mathbf{B}(\mathbf{t}_{ij}) \circ \bar{\mathbf{w}}_{\mathbf{x}_{ij}})), \end{aligned}$$

$$\beta_{0,K-1}(\mathbf{x}_{ij}) \leq \beta_{0,K-2}(\mathbf{x}_{ij}) \leq \cdots \leq \beta_{0,1}(\mathbf{x}_{ij}), \xi_{ijk} \geq 0.$$

It is clear that (2.11) is a standard linear programming formulation with respect to $\beta_{\mathbf{x}_{ij}}$ and ξ_{ijk} , and can be efficiently solved by the popular interior-point algorithm which is available in Python package `cvxopt` [29].

The parallel block-wise coordinate descent algorithm for the proposed method is summarized in Algorithm 1. In essence, this is an implementation of the block successive convex minimization, and hence it is guaranteed to converge to a stationary point [30].

Algorithm 1:

(Initialization): Set initial values $(\beta^{(0)}, \beta_0^{(0)}, \mathbf{W}^{(0)})$, tuning parameters λ , the tolerance error ϵ_{tol} , and iteration $t = 1, \epsilon_0 = 1$;

while $\epsilon_{t-1} > \epsilon_{tol}$ **do**

(Update for β):
Estimate $\beta^{(t)}$ by solving (2.9) with $\beta_0^{(t-1)}, \mathbf{W}^{(t-1)}$;

(Update for \mathbf{W}): Set $\epsilon_W = 1, \mathbf{W}^{new} = \mathbf{W}^{(t-1)}$;

while $\epsilon_W < \epsilon_{tol}$ **do**

$\mathbf{W}^{old} \leftarrow \mathbf{W}^{new}$

for $(l = 1, \dots, L)$ **do**

Estimate $(\mathbf{w}_m^{(l)})^{new}$ by solving (2.10) in a parallel fashion with $\beta^{(t)}, \beta_0^{(t-1)}$
and $\mathbf{W}^{(l')}$ as $(\mathbf{W}^{(l')})^{(t)}$ for $l' = 1, \dots, l$;

end

$\epsilon_W = \sum_{l=1}^L \|(\mathbf{W}^{new})^{(l)} - (\mathbf{W}^{old})^{(l)}\|_F^2 / \sum_{l=1}^L \|(\mathbf{W}^{old})^{(l)}\|_F^2$;

end

$\mathbf{W}^{(t)} = \mathbf{W}^{new}$;

(Update for β_0):
Estimate $\beta_0^{(t)}$ by solving (2.11) with $\beta^{(t)}$ and $\mathbf{W}^{(t)}$;

Set $\epsilon_t = \sum_{l=1}^c \|(\mathbf{W}^{(l)})^{(t)} - (\mathbf{W}^{(l)})^{(t-1)}\|_F^2 / \sum_{l=1}^c \|(\mathbf{W}^{(l)})^{(t-1)}\|_F^2$
+ $\|\beta^{(t)} - \beta^{(t-1)}\|_2^2 / \|\beta^{(t-1)}\|_2^2 + \|\beta_0^{(t)} - \beta_0^{(t-1)}\|_2^2 / \|\beta_0^{(t-1)}\|_2^2$

end

3. Theory

This section establishes the asymptotic convergence of the proposed method in estimating the ideal sentiment function $\mathbf{f}^0 = (f_1^0, \dots, f_{K-1}^0)$, which is showed to satisfy that $\text{sgn}(f_k^0) = \text{sgn}(P(y \geq k | \mathbf{x}, \mathbf{t}) - \frac{1}{2})$ for $k = 1, \dots, K-1$ [17]. Then the regret of \mathbf{f} is defined as

$$e(\mathbf{f}, \mathbf{f}^0) = (K-1)^{-1} \sum_{k=1}^{K-1} E(I(\text{sgn}(f_k(\mathbf{x}, \mathbf{t})) \neq \text{sgn}(f_k^0(\mathbf{x}, \mathbf{t}))).$$

Furthermore, denote $\bar{V}(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) = (K-1)^{-1} \sum_{k=1}^{K-1} V(\text{sgn}(y-k)f_k(\mathbf{x}, \mathbf{t}))$, and $e_V(\mathbf{f}, \mathbf{f}^0) = E\bar{V}(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - E\bar{V}(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t}))$. We further denote $\mathcal{F} = \{\mathbf{f} = (f_1, \dots, f_{K-1}) : f_k(\mathbf{x}, \mathbf{t}) = \beta^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} + \beta_{0,k}(\mathbf{x}), \beta_{0,K-1}(\mathbf{x}) \leq \beta_{0,K-2}(\mathbf{x}) \leq \dots \leq \beta_{0,1}(\mathbf{x})\}$, where $\mathbf{E}_t = \mathbf{E} \circ (\mathbf{1}_p \otimes \mathbf{B}(\mathbf{t}))$. Let $n = \sum_{i=1}^u N_i$ denote the total number of observations, and two technical assumptions are made to quantify the asymptotic behavior of the proposed method.

Assumption A. For any $\xi_n > 0$, there exists $\mathbf{f}^* \in \mathcal{F}$ such that $e_V(\mathbf{f}^*, \mathbf{f}^0) \leq \xi_n$.

Assumption A assures that the approximation error of \mathcal{F} in approximating \mathbf{f}^0 is governed by ξ_n , which eventually will impact the asymptotic behavior of the proposed sentiment function $\hat{\mathbf{f}}$.

Next, let $V^T(u) = \min(V(u), T)$, where T is chosen so that $V(\text{sgn}(y-k)f_k^0(\mathbf{x}, \mathbf{t})) \leq T$ almost surely. In addition, we let $\bar{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) = (K-1)^{-1} \sum_{k=1}^{K-1} V^T(\text{sgn}(y-k)f_k(\mathbf{x}, \mathbf{t}))$ and define $e_{VT}(\mathbf{f}, \mathbf{f}^0) = E(\bar{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - E(\bar{V}^T(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t})))$.

Assumption B. There exist constants $\alpha > 0, 1 \geq \gamma \geq 0$ and $a_1, a_2 > 0$ such that for any sufficiently small $\delta_n > 0$,

$$\begin{aligned} \sup_{\mathbf{f} \in \mathcal{F}_{\delta_n}} e(\mathbf{f}, \mathbf{f}^0) &\leq a_1 \delta_n^\alpha, \\ \sup_{\mathbf{f} \in \mathcal{F}_{\delta_n}} \text{Var}(\bar{V}^T(\text{sgn}(y-k), f_k(\mathbf{x}, \mathbf{t})) - \bar{V}^T(\text{sgn}(y-k), f_k^0(\mathbf{x}, \mathbf{t}))) &\leq a_2 \delta_n^\gamma, \end{aligned} \quad (3.1)$$

where $\mathcal{F}_{\delta_n} = \{\mathbf{f} \in \mathcal{F} : e_{VT}(\mathbf{f}, \mathbf{f}^0) \leq \delta_n\}$.

Assumption B implies the local smoothness of $e(\mathbf{f}, \mathbf{f}^0)$ and $\text{Var}(\bar{V}^T(\text{sgn}(y-k)f_k(\mathbf{x}, \mathbf{t})) - \bar{V}^T(\text{sgn}(y-k)f_k^0(\mathbf{x}, \mathbf{t})))$ within a neighborhood of \mathbf{f}^0 . Here α and γ are determined by the joint distribution of (\mathbf{x}, \mathbf{t}) and the loss function V . Additionally, (3.1) provides a connection between the first and second moments of $\bar{V}^T(\text{sgn}(y-k)f_k(\mathbf{x}, \mathbf{t})) - \bar{V}^T(\text{sgn}(y-k)f_k^0(\mathbf{x}, \mathbf{t}))$, which is essential for establishing the subsequent large deviation inequalities. In fact, Assumption B is a mild assumption and has been verified for various losses and distributions in literature [31, 32]. For example, Assumption B holds true for the hinge loss and any distribution $P(\mathbf{x}, \mathbf{t})$ with $\alpha = 1$ and $\gamma = 1$ [32].

Theorem 1. Suppose that Assumptions A-B are met. Then for the proposed sentiment function $\hat{\mathbf{f}}$ defined in (2.7), there exists a constant $a_3 > 0$ such that

$$P(e(\hat{\mathbf{f}}, \mathbf{f}^0) \geq a_1 \delta_n^{2\alpha}) \leq 3.5 \exp(-a_3 n (\lambda J^*)^{2-\min(\gamma, 1)}),$$

provided that $\delta_n^2 \geq 4\lambda J^*$, where $J^* = \min\{\|\beta\|_2^2 + \|\mathbf{W}\|_F^2 : \text{sgn}(f_k) = \text{sgn}(f_k^0); k = 1, \dots, K-1, \mathbf{f} \in \mathcal{F}\}$, $\delta_n^2 = \min(\epsilon_n^2 + 2\xi_n, 1)$, and $\epsilon_n^2 = (D_1 n^{-1} \log(n/D_1))^{1/(2-\gamma)}$, where $D_1 = \max\{p + d \sum_{l=1}^c m_l, \overline{MK}\}$, $\overline{M} = \prod_{l=1}^L m_l$, and a_1, γ, ξ_n are defined as in Assumptions A-B.

Corollary 1. *Suppose the assumptions of Theorem 1 are met, and that λ satisfies $n^{-1}(\lambda J^*)^{\min(\gamma, 1)-2} = o(1)$. Then it holds true that*

$$e(\hat{\mathbf{f}}, \mathbf{f}^0) = O_p(\delta_n^{2\alpha}) \text{ and } E|e(\hat{\mathbf{f}}, \mathbf{f}^0)| = O(\delta_n^{2\alpha}).$$

Clearly, the rate δ_n^2 is governed by both ϵ_n^2 and ξ_n , where ϵ_n^2 is determined by the complexity of \mathcal{F}^V depending on the dimension of the embedded space p , the size of sentiment lexicon d , and the number of levels of all covariates $m_l; l = 1, \dots, L$. Usually, there is a trade-off between the approximation error ξ_n and the complexity of \mathcal{F}^V over the choice of \mathbf{f}^0 , so as to attain the optimal convergence rate $\delta_n^{2\alpha}$.

4. Numerical experiments

In this section, we conduct a series of numerical experiments on simulated datasets and the Yelp challenge dataset to examine the performance of the proposed method. We compare it against various baseline word embedding methods in literature, including word embeddings based on Google news trained by word2vec technique [8], word embeddings based on Wikipedia trained by GloVe [9], and random word embeddings generated by multivariate normal distribution. Here the random word embeddings are used as baseline to verify the effectiveness of the other two pre-trained embeddings. Moreover, we let Google_p , Wiki_p , and Random_p denote the corresponding covariate-assisted word embeddings, whereas Google , Wiki and Random denote their corresponding baseline embeddings in (2.4), respectively.

For each method, the tuning parameters are selected via grid search over $[10^{-6}, 10^3]$, and their numerical performance is measured by MAE evaluated on a test set,

$$\text{TE}(\mathbf{f}) = \frac{1}{n_{\text{test}}(K-1)} \sum_{k=1}^{K-1} \sum_{i=1}^{n_{\text{test}}} I(\text{sgn}(y_i - k) \neq \text{sgn}(f_k(\mathbf{x}_i, \mathbf{t}_i))), \quad (4.1)$$

where n_{test} is the size of the test set.

4.1. Yelp challenge

The Yelp challenge dataset consists of four parts, including “business”, “review”, “user” and “check-in”, and is publicly available at <https://www.yelp.com/dataset/challenge>. In “business”, it contains location, latitude-longitude, averaged stars, opening hours, review counts and business categories. In “review”, a specific review is composed of textual comment, stars, business, user and corresponding feedback given by other users. In “user”, personal information associated with each user is given, including users’ social network, starting time and elite-experience in the Yelp community. Additionally, users’ behavior like votes and stars are also provided. In “checking”, the counts of check-ins at each

business are provided. In fact, all capitalized words in the reviews are converted into lower case by using the `nltk` package in Python. Other pre-processing steps, including removing spaces, stop words and punctuation, are also conducted.

We implement the proposed method based on the “review” and “user” parts. The “review” part contains reviews with star ratings from 1 to 5. Due to the imbalance of classes in “stars” of reviews, we encode “1” and “2” as 1, “3” as 2, and “4” and “5” as 3. In the pre-processing step, all capitalized words are converted into lower case, stop words and punctuation are removed, and frequency vectors are constructed for each review under the bag-of-words framework against a sentiment lexicon consisting of about 6,800 positive and negative words [33] combined with 1,000 1-gram features extracted based on term frequency–inverse document frequency (TF-IDF). The “user” part provides a personal social network, number of fans, counts of “useful”, counts of “cool”, counts of “funny” and elite-experience. In particular, elite-experience indicates the years when the user was selected as elite for well-written reviews, high-quality tips, or a detailed personal profile. Furthermore, elite users are characterized by pertinent comments and useful tips, with which other users have resonated to cast “useful”, “funny” and “cool” votes. About 3.25% of the users in the Yelp community have elite experience.

One salient difference between elite and non-elite users is their preference of wording, in particular the frequencies of sentiment words in reviews. For instances, as showed in Figure 1, “reputable”, “diligence” and “abnormal” are used much more frequently by non-elite users, whereas “slut” and “catchy” are much more popular among elite users. Also, as surprising as it appears, non-elite users tend to use “reputable” in an ironic way to show that the service they receive does not live up to their expectations, such as “A reputable apartment would try to fix their errors”, “I’m spending my money on an experience from

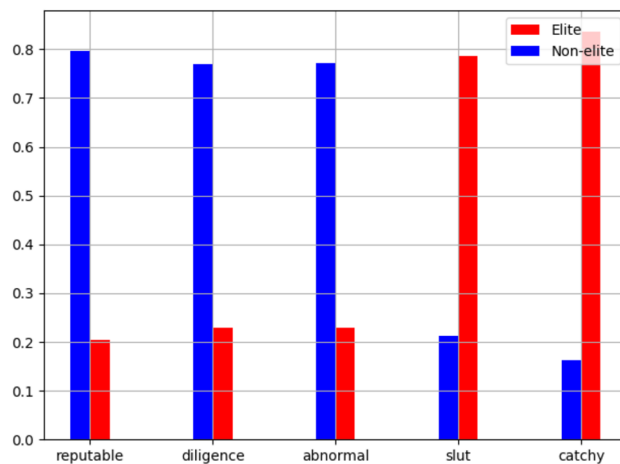


Fig 1: Relative frequency of words in average by elite and non-elite users.

a reputable Salon that has been nothing but rude and unhelpful”, and so on. In sharp contrast, elite users use “reputable” as a positive comment, leading to a 3.34 stars in average among those reviews containing “reputable”, while only 2.00 for by non-elite reviews.

Another interesting difference between elite and non-elite users is their preferences in giving “stars”. As seen in Figure 2, the distribution of averaged stars given by elite users appears to be a normal distribution, whereas the stars given by non-elite users appear to be more disordered in that they tend to give 1-star or 5-star reviews.

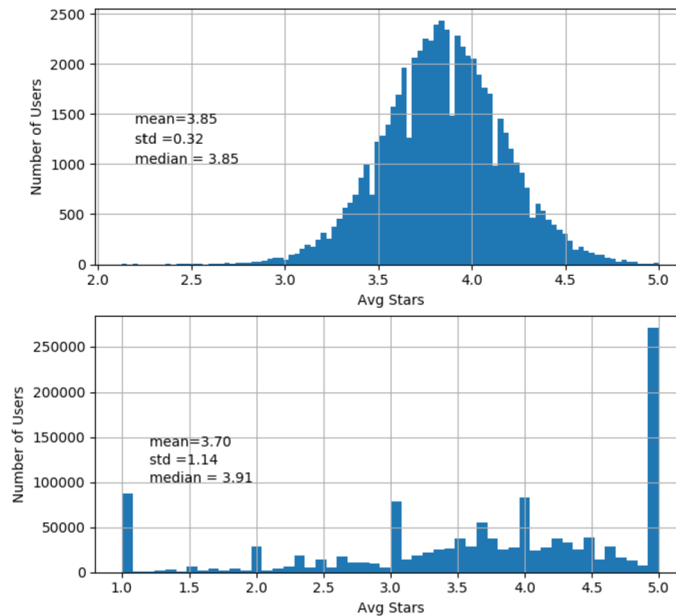


Fig 2: Histograms of averaged stars for elite and non-elite users, respectively.

Furthermore, number of feedbacks including “helpful”, “cool”, and “funny” tend to provide useful information about users. Users with a large number of feedbacks are popular for their objective comments, interesting expressions or humorous reviews. In fact, these three covariates are proportional to each other, and hence we only include “helpful” for application, which is converted to a binary covariate indicating whether “useful” count of the user is in the top 10 percent.

In this numerical experiment, 100,000 reviews are sampled from the Yelp challenge dataset, which are then split into three equal-sized sets used for training, validation and testing. The averaged test errors and their standard errors over 100 replications are reported in Table 1. To further evaluate the effectiveness of the proposed method in sentiment analysis, we also include the word embeddings trained by Embeddings from Language Models (ELMo) for comparison

[34], which generate word embeddings with information of context. Specifically, each review is converted to averaged word embeddings of length 1,024 obtained from ELMo based on 1 Billion Word Benchmark [35].

TABLE 1
Averaged test errors of various methods as well as their standard errors (in parentheses) over 100 replications in the Yelp dataset.

Google	Google _p	Wiki	Wiki _p	Random	Random _p	ELMo
0.1492 (0.0002)	0.1407 (0.0002)	0.1534 (0.0002)	0.1410 (0.0002)	0.1833 (0.0003)	0.1443 (0.0001)	0.1423 (0.0003)

Table 1 shows that the proposed method is able to improve the performance of Google, Wiki and Random by incorporating covariates, with improvement ranging from 5.7% to 21.3%. Even though word embeddings by ELMo appears to be more accurate than those trained by Word2Vec and GloVe, the proposed method yields the best performance on the word embeddings on Google news, showing that the proposed method is capable of learning covariate-varying word embeddings to improve prediction accuracy. Interestingly, when using covariates to adjust random embeddings which does not make use of semantic information at all, the improvement is much more significant and the performance of Random_p is almost comparable to Google_p and Wiki_p, showing that the proposed method is capable of training a word embeddings adaptive to prediction task. To verify the significance of the improvement, we further conduct t-tests between the proposed method on three word embeddings with their corresponding baselines, and the best performer with ELMo.

TABLE 2
P-values for various pair-wise t-tests.

Google vs Google _p	Wiki vs Wiki _p	Random vs Random _p	Google _p vs ELMo
0.0000	0.0000	0.0000	0.0003

As shown in Table 2, the improvements on three baseline word embeddings are statistically significant, showing that a uniform word embeddings in sentiment analysis may lead to sub-optimal performance, and hence entails adjusting effect from covariates. Additionally, Google_p also outperforms ELMo with statistically significant improvement, suggesting the proposed framework is competitive in sentiment analysis.

4.2. Simulations

We further verify the effectiveness of our proposed method under the assumed model in 2.4 with various numbers of covariates and degrees of covariate effect. The simulated examples are generated as follows. We first choose 100 words from the sentiment lexicon [33] and obtain corresponding word embeddings \mathbf{E} based on GoogleNews. Then we generate $\mathbf{w}_j = (\mathbf{1}_{100-r}, \tilde{\mathbf{w}}_j)$; $j = 1, \dots, m$, where $\tilde{\mathbf{w}}_j \in$

\mathbb{R}^r with each component generated from $unif(-1, 1)$ and r adjusts the degree of covariate effect. Then we generate $(x_i, \mathbf{b}_i); i = 1, \dots, n$, where $\mathbf{b}_i \in \mathbb{R}^{100}$ is a sequence of integers denoting word frequencies with each element of \mathbf{b}_i generated independently from $Pois(1)$, and x_i is uniformly chosen from $\{1, \dots, m\}$. The sentiment level y_i is generated via $y_i = \max\{k : \boldsymbol{\beta}^T \mathbf{E}(\mathbf{b}_i \circ \mathbf{w}_{x_i}) + \beta_{0,k}(\mathbf{x}_i) \leq 0\}$, where $\boldsymbol{\beta}$ is generated from $N(\mathbf{0}, \mathbf{I}_{300})$ and β_0 is set to generate K equal-sized classes.

Under this data generation scheme, we consider various cases with $(n, m, K) = (2000, 2, 5), (2000, 4, 5), (4000, 2, 5)$ and $(4000, 4, 5)$ and $r = 20, 40, 60, 80$, respectively. In each case, we split the dataset into training, validation and test sets with ratio 1:1:1. The averaged test errors of various methods over 100 replications are summarized in Table 3.

TABLE 3
Averaged test errors of various methods as well as their standard errors (in parentheses) over 100 replications.

	(n, m, K)	Google _p	Google
$r = 20$	(2000,2,5)	0.0679(0.0030)	0.1122(0.0030)
	(2000,4,5)	0.0794(0.0036)	0.1261(0.0027)
	(4000,2,5)	0.0659(0.0035)	0.1135(0.0034)
	(4000,4,5)	0.0664(0.0031)	0.1246(0.0024)
$r = 40$	(2000,2,5)	0.0632(0.0028)	0.1446(0.0024)
	(2000,4,5)	0.0697(0.0024)	0.1687(0.0021)
	(4000,2,5)	0.0503(0.0026)	0.1340(0.0026)
	(4000,4,5)	0.0583(0.0024)	0.1575(0.0018)
$r = 60$	(2000,2,5)	0.0497(0.0008)	0.1593(0.0024)
	(2000,4,5)	0.0612(0.0008)	0.1936(0.0024)
	(4000,2,5)	0.0443(0.0015)	0.1559(0.0023)
	(4000,4,5)	0.0495(0.0013)	0.1957(0.0015)
$r = 80$	(2000,2,5)	0.0420(0.0008)	0.1830(0.0024)
	(2000,4,5)	0.0539(0.0008)	0.2311(0.0013)
	(4000,2,5)	0.0380(0.0012)	0.1766(0.0020)
	(4000,4,5)	0.0427(0.0007)	0.2273(0.0015)

As shown in Table 3, the proposed method outperforms its baseline embedding method under all settings, showing that our proposed method is capable of enhancing the sentiment performance by incorporating covariate effect into word embeddings. The advantage becomes more substantial when r gets larger and \mathbf{w}_j 's become more different, showing that employing a homogeneous word embeddings may yield poor performance when embeddings varies with some covariates.

To verify the efficiency of the proposed method, we examine the computing time of three sub-optimization tasks, where the sample size increases from 1,000 to 10,000, or the dictionary size increases from 100 to 500. The averaged computing time over 50 replications of three sub-optimization tasks under all settings are reported in Figure 3.

As shown in Figure 3, the averaged computing time of three optimizations tasks are all linearly proportional to sample size. This is due to the fact that the

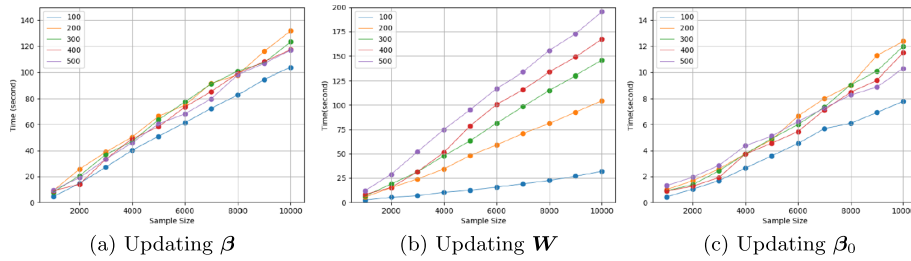


Fig 3: The averaged computing time of 50 replications of three sub-optimization tasks.

optimization tasks for β and \mathbf{W} resemble linear SVM in nature except varying intercepts, and hence can be solved efficiently by a dual coordinate descent algorithm as in `Liblinear` [28]. Moreover, the computational time for updating \mathbf{W} also depends on the size of dictionary, whereas those for updating β and β_0 appear to be less affected.

5. Summary

This article proposes a flexible framework for covariate-assisted sentiment analysis by incorporating covariates into word embeddings to improve prediction accuracy. Specifically, the proposed method admits varying document representations over covariate information, such as gender, education level and so on. This is also equivalent to admitting varying sentiment functions over levels of covariates, and then endows the proposed method with the ability to capture distinctions in wording and sentiment derived from covariate. Additionally, we propose an scalable block-wise coordinate descent algorithm to solve the resultant optimization task. We also establish the asymptotic properties of the proposed method, which provides a theoretical guarantee of its convergence to the ideal sentiment function. Note that even though our proposed method is formulated under the ordinal regression framework, the key idea of integrating covariates can be employed in other models.

Acknowledgments

The authors are also grateful to the editor, associate editor, and anonymous reviewers for their constructive comments and suggestions, which have significantly improved the manuscript.

Appendix

Lemma 2. *The solution to (7) remains the same as long as $\lambda_1\lambda_2^L$ stays the same.*

Proof of Lemma 2. Let

$$G(\boldsymbol{\beta}, \mathbf{W}, \boldsymbol{\beta}_0; \lambda_1, \lambda_2) = \frac{1}{N} \sum_{i=1}^n \sum_{j=1}^{N_i} \sum_{k=1}^{K-1} V(\text{sgn}(y_{ij} - k) \mathbf{f}_k(\mathbf{x}_{ij}, t_{ij})) + \lambda_1 \|\boldsymbol{\beta}\|^2 + \lambda_2 \|\mathbf{W}\|_F^2,$$

where $\|\mathbf{W}\|_F = \sum_{l=1}^L \|\mathbf{W}^{(l)}\|_F^2$, and then

$$(\hat{\boldsymbol{\beta}}, \hat{\mathbf{W}}, \hat{\boldsymbol{\beta}}_0) = \underset{\boldsymbol{\beta}, \mathbf{W}, \boldsymbol{\beta}_0}{\text{argmin}} G(\boldsymbol{\beta}, \mathbf{W}, \boldsymbol{\beta}_0, \lambda_1, \lambda_2).$$

Consider another tuning parameter pair $(\tilde{\lambda}_1, \tilde{\lambda}_2)$ satisfying $\lambda_1 \lambda_2^L = \tilde{\lambda}_1 \tilde{\lambda}_2^L$, and denote $(\tilde{\boldsymbol{\beta}}, \tilde{\mathbf{W}}, \tilde{\boldsymbol{\beta}}_0) = \underset{\boldsymbol{\beta}, \mathbf{W}, \boldsymbol{\beta}_0}{\text{argmin}} G(\boldsymbol{\beta}, \mathbf{W}, \boldsymbol{\beta}_0; \tilde{\lambda}_1, \tilde{\lambda}_2)$. Then we have

$$G(\hat{\boldsymbol{\beta}}, \hat{\mathbf{W}}, \hat{\boldsymbol{\beta}}_0; \lambda_1, \lambda_2) = G(\sqrt{\lambda_1/\tilde{\lambda}_1} \hat{\boldsymbol{\beta}}, \sqrt{\lambda_2/\tilde{\lambda}_2} \hat{\mathbf{W}}, \hat{\boldsymbol{\beta}}_0, \tilde{\lambda}_1, \tilde{\lambda}_2).$$

It hence follows that $(\tilde{\boldsymbol{\beta}}, \tilde{\mathbf{W}}, \tilde{\boldsymbol{\beta}}_0) = (\sqrt{\lambda_1/\tilde{\lambda}_1} \hat{\boldsymbol{\beta}}, \sqrt{\lambda_2/\tilde{\lambda}_2} \hat{\mathbf{W}}, \hat{\boldsymbol{\beta}}_0)$, which leads to

$$\begin{aligned} & \tilde{\boldsymbol{\beta}}^T \mathbf{E}(\mathbf{B}(\mathbf{t}) \circ \bar{\mathbf{w}}_{\mathbf{x}}) \\ &= \sqrt{\lambda_1/\tilde{\lambda}_1} \hat{\boldsymbol{\beta}}^T \mathbf{E}(\mathbf{B}(\mathbf{t}) \circ \sqrt{\lambda_2/\tilde{\lambda}_2} \hat{\mathbf{w}}_{x_1}^{(1)} \circ \sqrt{\lambda_2/\tilde{\lambda}_2} \circ \cdots \circ \sqrt{\lambda_2/\tilde{\lambda}_2} \hat{\mathbf{w}}_{x_L}^{(L)}) \\ &= \sqrt{\lambda_1 \lambda_2^L / (\tilde{\lambda}_1 \tilde{\lambda}_2^L)} \hat{\boldsymbol{\beta}}^T \mathbf{E}(\mathbf{B}(\mathbf{t}) \circ \bar{\mathbf{w}}_{\mathbf{x}}). \end{aligned}$$

The desired result then follows immediately after the fact that $\lambda_1 \lambda_2^L = \tilde{\lambda}_1 \tilde{\lambda}_2^L$. \square

Lemma 3. Let $C_1(\tau) = d^{\frac{1}{2}}(\tau J^*)^{\frac{L}{2}+1} / L^{\frac{L}{2}} \|\mathbf{E}_{\mathbf{t}}\|_F$, then for any $\beta_{0,k}(\mathbf{x})$ with $|\beta_{0,k}(\mathbf{x})| \geq C_1(\tau) + T + 1$, we have

$$\begin{aligned} & V^T(\text{sgn}(y - k)(\boldsymbol{\beta}^T \mathbf{E}_{\mathbf{t}} \bar{\mathbf{w}}_{\mathbf{x}} + \beta_{0,k}(\mathbf{x}))) \\ &= \begin{cases} TI(\text{sgn}(y - k) \neq 1), & \text{if } \beta_{0,k}(\mathbf{x}) \geq C_1(\tau) + T + 1; \\ TI(\text{sgn}(y - k) = 1), & \text{if } \beta_{0,k}(\mathbf{x}) \leq -C_1(\tau) - T - 1, \end{cases} \end{aligned} \quad (5.1)$$

for any $(\boldsymbol{\beta}, \mathbf{W}) \in \{(\boldsymbol{\beta}, \mathbf{W}) : \|\boldsymbol{\beta}\|^2 + \|\mathbf{W}\|_F^2 \leq J^* \tau\}$.

Proof of Lemma 3. We first show that $|\boldsymbol{\beta}^T \mathbf{E}_{\mathbf{t}} \bar{\mathbf{w}}_{\mathbf{x}}| \leq C_1(\tau)$, for any $(\boldsymbol{\beta}, \mathbf{W}) \in \{(\boldsymbol{\beta}, \mathbf{W}) : \|\boldsymbol{\beta}\|^2 + \|\mathbf{W}\|_F^2 \leq J^* \tau\}$. Particularly,

$$\begin{aligned} & |\boldsymbol{\beta}^T \mathbf{E}_{\mathbf{t}} \bar{\mathbf{w}}_{\mathbf{x}}| = |(\boldsymbol{\beta}^T \mathbf{E}_{t,1}, \boldsymbol{\beta}^T \mathbf{E}_{t,2}, \dots, \boldsymbol{\beta}^T \mathbf{E}_{t,d}) \bar{\mathbf{w}}_{\mathbf{x}}| \\ & \leq \|\boldsymbol{\beta}^T \mathbf{E}_{\mathbf{t}}\|_2 \|\bar{\mathbf{w}}_{\mathbf{x}}\|_2 = \left(\sum_{j=1}^d (\boldsymbol{\beta}^T \mathbf{E}_{t,j})^2 \right)^{1/2} \|\bar{\mathbf{w}}_{\mathbf{x}}\|_2 \\ & \leq \left(\sum_{j=1}^d \|\boldsymbol{\beta}\|_2^2 \|\mathbf{E}_{t,j}\|_2^2 \right)^{1/2} \|\bar{\mathbf{w}}_{\mathbf{x}}\|_2 = \|\boldsymbol{\beta}\|_2 \|\mathbf{E}_{\mathbf{t}}\|_F \|\mathbf{w}_{x_1}^{(1)} \circ \mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_L}^{(L)}\|_2 \end{aligned}$$

$$\begin{aligned}
&= \|\boldsymbol{\beta}\|_2 \|\mathbf{E}_t\|_F \left(\sum_{j=1}^d \left(\prod_{i=1}^L \mathbf{w}_{x_i,j}^{(i)} \right)^2 \right)^{\frac{1}{2}} \leq \|\boldsymbol{\beta}\|_2 \|\mathbf{E}_t\|_F \left(\sum_{j=1}^d \left(\sum_{i=1}^L \frac{(\mathbf{w}_{x_i,j}^{(i)})^2}{L} \right)^L \right)^{\frac{1}{2}} \\
&\leq (\|\boldsymbol{\beta}\|_2^2 + \|\mathbf{W}\|_F^2)^{1/2} \|\mathbf{E}_t\|_F d^{\frac{1}{2}} \left(\frac{\|\mathbf{W}\|_F^2}{L} \right)^{\frac{L}{2}} \\
&\leq (\|\boldsymbol{\beta}\|_2^2 + \|\mathbf{W}\|_F^2)^{1/2} \|\mathbf{E}_t\|_F d^{\frac{1}{2}} \left(\frac{\tau J^*}{L} \right)^{\frac{L}{2}} \leq d^{\frac{1}{2}} \frac{(\tau J^*)^{\frac{L}{2}+1}}{L^{\frac{L}{2}}} \|\mathbf{E}_t\|_F = C_1(\tau),
\end{aligned}$$

where $\mathbf{E}_{t,j}$ is the j -th column of \mathbf{E}_t , the first inequality follows from the Cauchy-Schwarz inequality, the third inequality follows from the inequality of arithmetic and geometric means.

Next we verify only (5.1) for $\beta_{0,k}(\mathbf{x}) \geq C_1(\tau) + T + 1$, and it can be verified similarly for $\beta_{0,k}(\mathbf{x}) \leq -C_1(\tau) - T - 1$. For any $\beta_{0,k}(\mathbf{x}) \geq C_1(\tau) + T + 1$, when $\text{sgn}(y - k) = 1$ we have

$$0 \leq V^T(\boldsymbol{\beta}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} + \beta_{0,k}(\mathbf{x})) \leq V^T(-|\boldsymbol{\beta}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}}| + C_1(\tau) + T + 1) \leq 0,$$

where both inequalities follow from the non-increasing property of $V^T(\cdot)$. When $\text{sgn}(y - k) = -1$, we have

$$T \geq V^T(-\boldsymbol{\beta}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} - \beta_{0,k}(\mathbf{x})) \geq V^T(|\boldsymbol{\beta}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}}| - C_1(\tau) - T - 1) \geq T.$$

The desirable result then follows. \square

Lemma 4. Let $C_2(\tau) = \|\mathbf{E}_t\|_F \sum_{k=0}^L M(\sqrt{J^* \tau})^k$. For any $f_k, \tilde{f}_k \in \mathcal{F}_k(\tau)$, we have

$$|f_k(\mathbf{x}, \mathbf{t}) - \tilde{f}_k(\mathbf{x}, \mathbf{t})| \leq C_2(\tau) \left(\|\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}\|_2 + \|\mathbf{W} - \tilde{\mathbf{W}}\|_F \right) + |\beta_{0,k}(\mathbf{x}) - \tilde{\beta}_{0,k}(\mathbf{x})|.$$

Proof of Lemma 4. Let $f_k, \tilde{f}_k \in \mathcal{F}_k(\tau)$, then we have

$$\begin{aligned}
&|f_k(\mathbf{x}, \mathbf{t}) - \tilde{f}_k(\mathbf{x}, \mathbf{t})| \\
&= |\boldsymbol{\beta}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} + \beta_{0,k}(\mathbf{x}) - \tilde{\boldsymbol{\beta}}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} - \tilde{\beta}_{0,k}(\mathbf{x})| \\
&= |\boldsymbol{\beta}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} + \beta_{0,k}(\mathbf{x}) - \tilde{\boldsymbol{\beta}}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} + \tilde{\boldsymbol{\beta}}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} - \tilde{\boldsymbol{\beta}}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} - \tilde{\beta}_{0,k}(\mathbf{x})| \\
&\leq |(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} + \tilde{\boldsymbol{\beta}}^T \mathbf{E}_t (\bar{\mathbf{w}}_{\mathbf{x}} - \tilde{\bar{\mathbf{w}}}_{\mathbf{x}})| + |\beta_{0,k}(\mathbf{x}) - \tilde{\beta}_{0,k}(\mathbf{x})| \\
&\leq \|\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}\|_2 \|\mathbf{E}_t\|_F \|\bar{\mathbf{w}}_{\mathbf{x}}\|_2 + \|\tilde{\boldsymbol{\beta}}\|_2 \|\mathbf{E}_t\|_F \|\bar{\mathbf{w}}_{\mathbf{x}} - \tilde{\bar{\mathbf{w}}}_{\mathbf{x}}\|_2 + |\beta_{0,k}(\mathbf{x}) - \tilde{\beta}_{0,k}(\mathbf{x})|
\end{aligned}$$

By the definition of $\mathcal{F}_k(\tau)$, $\boldsymbol{\beta}$ is bounded by $J^* \tau$. It suffices to bound $\|\bar{\mathbf{w}}_{\mathbf{x}}\|_2$ and $\|\bar{\mathbf{w}}_{\mathbf{x}} - \tilde{\bar{\mathbf{w}}}_{\mathbf{x}}\|_2$ respectively. For $\bar{\mathbf{w}}_{\mathbf{x}}$, by inequality of arithmetic and geometric means, we have

$$\|\mathbf{w}_{\mathbf{x}_1}^{(1)} \circ \mathbf{w}_{\mathbf{x}_2}^{(2)} \cdots \circ \mathbf{w}_{\mathbf{x}_L}^{(L)}\|_2^2 = \sum_{j=1}^d \left(\prod_{i=1}^L w_{x_i,j}^{(i)} \right)^2 \leq \sum_{j=1}^d \left(\sum_{i=1}^L \frac{(w_{x_i,j}^{(i)})^2}{L} \right)^L$$

$$\leq d \left(\frac{J(\mathbf{W})}{L} \right)^L \leq d \left(\frac{J^*\tau}{L} \right)^L \equiv M$$

Similarly for $\|\bar{\mathbf{w}}_{\mathbf{x}} - \tilde{\mathbf{w}}_{\mathbf{x}}\|_2$, for any integer $1 \leq m \leq L$, we have

$$\begin{aligned} & \|\mathbf{w}_{x_1}^{(1)} \circ \mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_m}^{(m)} - \tilde{\mathbf{w}}_{x_1}^{(1)} \circ \tilde{\mathbf{w}}_{x_2}^{(2)} \cdots \circ \tilde{\mathbf{w}}_{x_m}^{(m)}\|_2 \\ & \leq \|\mathbf{w}_{x_1}^{(1)} \circ \mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_m}^{(m)} - \tilde{\mathbf{w}}_{x_1}^{(1)} \circ \mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_m}^{(m)} + \\ & \quad \tilde{\mathbf{w}}_{x_1}^{(1)} \circ \mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_m}^{(m)} - \tilde{\mathbf{w}}_{x_1}^{(1)} \circ \tilde{\mathbf{w}}_{x_2}^{(2)} \cdots \circ \tilde{\mathbf{w}}_{x_m}^{(m)}\| \\ & = \|(\mathbf{w}_{x_1}^{(1)} - \tilde{\mathbf{w}}_{x_1}^{(1)}) \circ \mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_m}^{(m)} + \\ & \quad \tilde{\mathbf{w}}_{x_1}^{(1)} \circ (\mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_m}^{(m)} - \tilde{\mathbf{w}}_{x_2}^{(2)} \cdots \circ \tilde{\mathbf{w}}_{x_m}^{(m)})\| \\ & \leq \|\mathbf{w}_{x_1}^{(1)} - \tilde{\mathbf{w}}_{x_1}^{(1)}\|_2 \|\mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_m}^{(m)}\|_2 + \\ & \quad \|\tilde{\mathbf{w}}_{x_1}^{(1)} \circ (\mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_m}^{(m)} - \tilde{\mathbf{w}}_{x_2}^{(2)} \cdots \circ \tilde{\mathbf{w}}_{x_m}^{(m)})\| \\ & \leq d \|\mathbf{W} - \tilde{\mathbf{W}}\|_F \left(\frac{J^*\tau}{L} \right)^L + \|\tilde{\mathbf{w}}_{x_1}^{(1)}\|_2 \|\mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_m}^{(m)} - \tilde{\mathbf{w}}_{x_2}^{(2)} \cdots \circ \tilde{\mathbf{w}}_{x_m}^{(m)}\|_2 \\ & \leq d \|\mathbf{W} - \tilde{\mathbf{W}}\|_F \left(\frac{J^*\tau}{L} \right)^L + \sqrt{J^*\tau} \|\mathbf{w}_{x_2}^{(2)} \cdots \circ \mathbf{w}_{x_m}^{(m)} - \tilde{\mathbf{w}}_{x_2}^{(2)} \cdots \circ \tilde{\mathbf{w}}_{x_m}^{(m)}\|_2 \\ & \leq \sum_{k=1}^L M (\sqrt{J^*\tau})^k \|\mathbf{W} - \tilde{\mathbf{W}}\|_F \end{aligned}$$

where the last inequality follows by applying similar steps iteratively. This completes the proof. \square

Proof of Theorem 1. By Assumption B, we have $\{e(\hat{\mathbf{f}}, \mathbf{f}^0) \geq a_1 \delta_n^{2\alpha}\} \subset \{e_{VT}(\hat{\mathbf{f}}, \mathbf{f}^0) \geq \delta_n^2\}$, and thus it suffices to bound $P(e_{VT}(\hat{\mathbf{f}}, \mathbf{f}^0) \geq \delta_n^2)$. Since $\hat{\mathbf{f}}$ is a global minimizer of (7) of the manuscript, it yields that

$$\begin{aligned} & P(e(\hat{\mathbf{f}}, \mathbf{f}^0) \geq a_1 \delta_n^{2\alpha}) \\ & \leq P \left(\sup_{e_{VT}(\mathbf{f}, \mathbf{f}^0) \geq \delta_n^2} \sum_{i,j} \tilde{V}^T(y_{ij}, \mathbf{f}^*(\mathbf{x}_{ij}, \mathbf{t}_{ij})) - \sum_{i,j} \tilde{V}^T(y_{ij}, \mathbf{f}(\mathbf{x}_{ij}, \mathbf{t}_{ij})) \geq 0 \right) \equiv I, \end{aligned} \quad (5.2)$$

where $\tilde{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) = \bar{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) + \lambda J(\boldsymbol{\beta}) + \lambda J(\mathbf{W})$. Next we define a scaled empirical process as

$$\begin{aligned} & E_n(\tilde{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t}))) - \tilde{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) \\ & = \frac{1}{n} \sum_{i,j} (\tilde{V}^T(y_{ij}, \mathbf{f}^*(\mathbf{x}_{ij}, \mathbf{t}_{ij}))) - \tilde{V}^T(y_{ij}, \mathbf{f}(\mathbf{x}_{ij}, \mathbf{t}_{ij})) \\ & \quad - E(\tilde{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t}))) - \tilde{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})). \end{aligned}$$

Let $A_{ij} = \{\mathbf{f} \in \mathcal{F} : 2^{i-1} \delta_n^2 \leq e_{VT}(\mathbf{f}, \mathbf{f}^0) \leq 2^i \delta_n^2, 2^{j-1} J^* \leq J(\boldsymbol{\beta}) + J(\mathbf{W}) \leq 2^j J^*\}$, and $A_{i0} = \{\mathbf{f} \in \mathcal{F} : 2^{i-1} \delta_n^2 \leq e_{VT}(\mathbf{f}, \mathbf{f}^0) \leq 2^i \delta_n^2, J(\boldsymbol{\beta}) + J(\mathbf{W}) \leq J^*\}$

for $i, j \geq 1$. Then we have

$$\begin{aligned}
I &\leq \sum_{i,j \geq 1} P\left(\sup_{\mathbf{f} \in A_{ij}} \sum_{i,j} (\tilde{V}^T(y_{ij}, \mathbf{f}^*(\mathbf{x}_{ij}, \mathbf{t}_{ij})) - \tilde{V}(y_{ij}, \mathbf{f}(\mathbf{x}_{ij}, \mathbf{t}_{ij}))) \geq 0\right) \\
&\quad + \sum_{i=1}^{\infty} P\left(\sup_{\mathbf{f} \in A_{i0}} \sum_{i,j} (\tilde{V}^T(y_{ij}, \mathbf{f}^*(\mathbf{x}_{ij}, \mathbf{t}_{ij})) - \tilde{V}^T(y_{ij}, \mathbf{f}(\mathbf{x}_{ij}, \mathbf{t}_{ij}))) \geq 0\right) \\
&= \sum_{i,j \geq 1} P\left(\sup_{\mathbf{f} \in A_{ij}} \frac{1}{n} \sum_{i,j} (\tilde{V}^T(y_{ij}, \mathbf{f}^*(\mathbf{x}_{ij}, \mathbf{t}_{ij})) - \tilde{V}(y_{ij}, \mathbf{f}(\mathbf{x}_{ij}, \mathbf{t}_{ij}))) + \right. \\
&\quad \left. M(i, j) \geq M(i, j)\right) \\
&\quad + \sum_{i=1}^{\infty} P\left(\sup_{\mathbf{f} \in A_{i0}} \frac{1}{n} \sum_{i,j} (\tilde{V}^T(y_{ij}, \mathbf{f}^*(\mathbf{x}_{ij}, \mathbf{t}_{ij})) - \tilde{V}^T(y_{ij}, \mathbf{f}(\mathbf{x}_{ij}, \mathbf{t}_{ij}))) + \right. \\
&\quad \left. M(i, 0) \geq M(i, 0)\right) \\
&\leq \sum_{i,j \geq 1} P\left(\sup_{\mathbf{f} \in A_{ij}} E_n(\bar{V}^T(y, \mathbf{f}^*(\mathbf{x})) - \bar{V}^T(y, \mathbf{f}(\mathbf{x}))) \geq M(i, j)\right) \\
&\quad + \sum_{i=1}^{\infty} P\left(\sup_{\mathbf{f} \in A_{i0}} E_n(\bar{V}^T(y, \mathbf{f}^*(\mathbf{x})) - \bar{V}^T(y, \mathbf{f}(\mathbf{x}))) \geq M(i, 0)\right) \equiv I_1 + I_2.
\end{aligned}$$

To bound I , we proceed to bound I_1 and I_2 respectively. By Assumption A, we have

$$\begin{aligned}
&E(\bar{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - E\bar{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t}))) \\
&= E(\bar{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - E\bar{V}^T(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t}))) \\
&\quad + E(\bar{V}^T(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t})) - E\bar{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t}))) \\
&\geq e_{V^T}(\mathbf{f}, \mathbf{f}^0) - \delta_n^2/2,
\end{aligned}$$

where the last inequality follows from Assumption A. Then for any $\mathbf{f} \in A_{ij}; i \geq 1, j \geq 1$,

$$E(\tilde{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - \tilde{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t}))) \geq (2^{i-1} - \frac{1}{2})\delta_n^2 + \lambda(2^{j-1} - 1)J^* \equiv M(i, j),$$

and for any $\mathbf{f} \in A_{i0}; i \geq 1$,

$$E(\tilde{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - \tilde{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t}))) \geq (2^{i-1} - \frac{1}{2})\delta_n^2 - \lambda J^* \geq 2^{i-3}\delta_n^2 \equiv M(i, 0).$$

For the variance, by Assumptions A and B,

$$\begin{aligned}
&\sup_{\mathbf{f} \in A_{ij}} \text{Var}\left(\tilde{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - \tilde{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t}))\right) \\
&= \sup_{\mathbf{f} \in A_{ij}} \text{Var}\left(\bar{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - \bar{V}(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t}))\right)
\end{aligned}$$

$$\begin{aligned}
 & + \bar{V}(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t})) - \bar{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t})) \\
 = & \sup_{\mathbf{f} \in A_{ij}} \text{Var} \left(\bar{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - \bar{V}(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t})) \right) \\
 & + \text{Var} \left(\bar{V}(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t})) - \bar{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t})) \right) \\
 & + 2\text{Cov}(\bar{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - \bar{V}(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t})), \bar{V}(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t})) - \bar{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t}))) \\
 \leq & 2 \sup_{\mathbf{f} \in A_{ij}} \text{Var} \left(\bar{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - \bar{V}(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t})) \right) + \\
 & 2\text{Var} \left(\bar{V}(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t})) - \bar{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t})) \right) \\
 \leq & 2a_2(M(i, j))^\gamma + 2a_2\xi_n^\gamma \leq 4a_2(M(i, j))^\gamma \equiv v(i, j).
 \end{aligned}$$

To bound I_1 and I_2 , we first denote $J^* = J(\boldsymbol{\beta}^*) + J(\mathbf{W}^*)$, and $\mathcal{F}(\tau) = \{\mathbf{f} = (f_1, f_2, \dots, f_{K-1}) : f_k \in \mathcal{F}_k(\tau)\}$ with

$$\mathcal{F}_k(\tau) = \{f_k(\mathbf{x}, \mathbf{t}) = \boldsymbol{\beta}^T \mathbf{E}_t \bar{\mathbf{w}}_{\mathbf{x}} + \beta_{0,k}(\mathbf{x}) : J(\boldsymbol{\beta}) + J(\mathbf{W}) \leq \tau J^*\}.$$

Then the associated spaces of the loss function are

$$\begin{aligned}
 \mathcal{F}^V(\tau) & = \{(K-1)^{-1} \sum_{k=1}^{K-1} V^T(\text{sgn}(y-k)f_k(\mathbf{x}, \mathbf{t})) : \mathbf{f} \in \mathcal{F}(\tau)\}, \\
 \mathcal{F}_k^V(\tau) & = \{V^T(\text{sgn}(y-k)f_k(\mathbf{x}, \mathbf{t})) : f_k \in \mathcal{F}_k(\tau)\}.
 \end{aligned}$$

Next, we employ the L_2 -metric entropy to measure the cardinality of $\mathcal{F}^V(\tau)$, denoted as $H(\epsilon, \mathcal{F}^V(\tau))$. For $\bar{V}(y, f(\mathbf{x}, \mathbf{t})), \tilde{\bar{V}}(y, f(\mathbf{x}, \mathbf{t})) \in \mathcal{F}^V(\tau)$, we have

$$\begin{aligned}
 & \|\bar{V}(y, f(\mathbf{x}, \mathbf{t})) - \tilde{\bar{V}}(y, f(\mathbf{x}, \mathbf{t}))\|_2 \\
 \leq & \frac{1}{K-1} \sum_{k=1}^{K-1} \|V^T(\text{sgn}(y-k)f_k(\mathbf{x}, \mathbf{t})) - V^T(\text{sgn}(y-k)\tilde{f}_k(\mathbf{x}, \mathbf{t}))\|_2. \quad (5.3)
 \end{aligned}$$

Therefore, it suffices to compute the metric entropy of $\mathcal{F}_k^V(\tau)$.

Let $\bar{M} = \prod_{k=1}^L m_k$, and then $\mathcal{F}_k^V(\tau)$ can be expressed as the union of $4^{\bar{M}}$ subspaces

$$\mathcal{F}_k^V(\tau) = \bigcup_{l,r \leq 2^{\bar{M}}} \mathcal{F}_k^V(\tau, \mathbf{p}_l, \mathbf{q}_r), \quad (5.4)$$

where $\mathcal{F}_k^V(\tau, \mathbf{p}_l, \mathbf{q}_r)$ is defined as

$$\begin{aligned}
 \mathcal{F}_k^V(\tau, \mathbf{p}_l, \mathbf{q}_r) & = \{V^T(\text{sgn}(y-k)f_k(\mathbf{x}, \mathbf{t})), f_k \in \mathcal{F}_k(\tau), \\
 & p_{lu}|\beta_{0,k}(\mathbf{x}^{(u)})| \leq p_{lu}(C_1(\tau) + T + 1), q_{ru}\beta_{0,k}(\mathbf{x}^{(u)}) \leq 0; u = 1, \dots, \bar{M}\},
 \end{aligned}$$

where $\mathbf{p}_l, \mathbf{q}_r; l, r = 1, \dots, 2^{\bar{M}}$ taking all possible values in $\{-1, 1\}^{\bar{M}}$. Next, we proceed to compute the entropy of $\mathcal{F}_k^V(\tau, \mathbf{p}_l, \mathbf{q}_r)$. For any two functions $V^T(\text{sgn}(y-k)f_k(\mathbf{x}, \mathbf{t})), V^T(\text{sgn}(y-k)\tilde{f}_k(\mathbf{x}, \mathbf{t})) \in \mathcal{F}_k^V(\tau, \mathbf{p}_l, \mathbf{q}_r)$, we have

$$|V^T(\text{sgn}(y-k)f_k(\mathbf{x}, \mathbf{t})) - V^T(\text{sgn}(y-k)\tilde{f}_k(\mathbf{x}, \mathbf{t}))|$$

$$\begin{aligned}
&\leq \sum_{u=1}^{\overline{M}} I(\mathbf{x} = \mathbf{x}^{(u)}) (I(p_{lu} = -1) \cdot (I(q_{ru} = 1) \cdot 0 + I(q_{ru} = -1) \cdot 0) \\
&\quad + I(p_{lu} = 1) |f_k(\mathbf{x}, \mathbf{t}) - \tilde{f}_k(\mathbf{x}, \mathbf{t})|) \\
&\leq \sum_{u=1}^{\overline{M}} I(\mathbf{x} = \mathbf{x}^{(u)}) I(p_{lu} = 1) \left(C_2(\tau) (\|\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}\|_2 + \|\mathbf{W} - \tilde{\mathbf{W}}\|_F) \right. \\
&\quad \left. + |\beta_{0,k}(\mathbf{x}) - \tilde{\beta}_{0,k}(\mathbf{x})| \right) \\
&\leq C_2(\tau) \sqrt{2(\|\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}\|_2^2 + \|\mathbf{W} - \tilde{\mathbf{W}}\|_F^2)} \\
&\quad + \sum_{u=1}^{\overline{M}} I(\mathbf{x} = \mathbf{x}^{(u)}) I(p_{lu} = 1) |\beta_{0,k}(\mathbf{x}) - \tilde{\beta}_{0,k}(\mathbf{x})|,
\end{aligned}$$

where the first and second inequalities follow from Lemma 3 and Lemma 4 respectively. Combined with (5.3), we have

$$\begin{aligned}
&\|\bar{V}(y, f(\mathbf{x}, \mathbf{t})) - \tilde{V}(y, f(\mathbf{x}, \mathbf{t}))\|_2 \\
&\leq C_2(\tau) \sqrt{2(\|\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}\|_2^2 + \|\mathbf{W} - \tilde{\mathbf{W}}\|_F^2)} \\
&\quad + \frac{1}{K-1} \sum_{k=1}^{K-1} \sum_{u=1}^{\overline{M}} I_{\{\mathbf{x}=\mathbf{x}^{(u)}\}} I_{\{p_{lu}=1\}} |\beta_{0,k}(\mathbf{x}) - \tilde{\beta}_{0,k}(\mathbf{x})|.
\end{aligned}$$

Let $\mathcal{B}(\tau J^*) = \{(\boldsymbol{\beta}, \mathbf{W}) : \|\boldsymbol{\beta}\|_2^2 + \|\mathbf{W}\|_F^2 \leq \tau J^*\}$, and the corresponding $\epsilon/(2\sqrt{2}C_2(\tau))$ -covering number with L_2 -norm is $O\left(\left(\frac{2C_2(\tau)\sqrt{2\tau J^*}}{\epsilon}\right)^{(p+d\sum_{l=1}^L m_l)}\right)$. For $u = 1, \dots, \overline{M}$ we define the set

$$\begin{aligned}
&\{(\beta_{0,1}(\mathbf{x}^{(u)}), \beta_{0,2}(\mathbf{x}^{(u)}), \dots, \beta_{0,K-1}(\mathbf{x}^{(u)})) : \\
&\quad - (C_1(\tau) + T + 1) \leq \beta_{0,1}(\mathbf{x}^{(u)}) \leq \dots \leq \beta_{0,K-1}(\mathbf{x}^{(u)}) \leq C_1(\tau) + T + 1\},
\end{aligned}$$

and its corresponding number of $\epsilon/2$ -covering balls with L_1 -norm is $C_{\tilde{N}+K-2}^{K-1}$ where $\tilde{N} = \lceil \frac{4(C_1(\tau)+T+1)}{\epsilon} \rceil$ with $\lceil \cdot \rceil$ being the ceiling function.

Finally, the metric entropy of $\mathcal{F}^V(\tau)$ can be upper-bounded as follows:

$$\begin{aligned}
H(\epsilon, \mathcal{F}^V(\tau)) &\leq \max \left\{ O\left((p + d \sum_{l=1}^L m_l) \log \frac{2C_2(\tau)(2\tau J^*)^{1/2}}{\epsilon} \right. \right. \\
&\quad \left. \left. + \overline{M}K \log\left(\frac{4(C_1(\tau) + T + 1)}{\epsilon} + K \right) - \overline{M} \log K! \right), 1 \right\} \\
&\leq \max \left\{ O\left((p + d \sum_{l=1}^L m_l) \log \frac{2C_2(\tau)(2\tau J^*)^{1/2}}{\epsilon} \right. \right. \\
&\quad \left. \left. + \overline{M}K \log\left(\frac{4(C_1(\tau) + T + 1)}{K\epsilon} + 1 \right) \right), 1 \right\}
\end{aligned}$$

$$\leq \max\{O(D_1 \log \frac{D_2(\tau)}{\epsilon} + 1), 1\},$$

where $D_1 = \max\{p+d\sum_{l=1}^L m_l, \overline{MK}\}$, $D_2(\tau) = \max\{2C_2(\tau)(2\tau J^*)^{1/2}, 4(C_1(\tau)+T+1)K^{-1}\}$, and the second inequality follows from the fact that $\log K! \geq K \log K - K - \log K$.

In the following, we proceed to verify (4.5) – (4.7) in [36]. We first notice that $\int_{\frac{\epsilon}{32}M(i,j)}^{v^{\frac{1}{2}}(i,j)} H^{\frac{1}{2}}(u, \mathcal{F}^V(\tau))du/M(i, j)$ is non-increasing in i and $M(i, j)$, then we have

$$\begin{aligned} & \int_{\frac{\epsilon}{32}M(i,j)}^{v^{\frac{1}{2}}(i,j)} H^{\frac{1}{2}}(u, \mathcal{F}^V(2^j))du/M(i, j) \\ & \leq \int_{\frac{\epsilon}{32}M(1,j)}^{v^{\frac{1}{2}}(1,j)} H^{\frac{1}{2}}(u, \mathcal{F}^V(2^j))du/M(1, j). \end{aligned} \quad (5.5)$$

Rearranging the right-hand side of (5.5) and let $\delta_n^2 = (D_1 n^{-1} \log(n/D_1))^{1/(2-\gamma)}$, yielding

$$\int_{\frac{\epsilon}{32}M(1,j)}^{v^{\frac{1}{2}}(1,j)} H^{\frac{1}{2}}(u, \mathcal{F}^V(2^j))du/M(1, j) \leq n^{\frac{1}{2}}.$$

Then, it is easy to see that $\inf_{\mathbf{f} \in A_{ij}} E(\tilde{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - \tilde{V}^T(y, \mathbf{f}^*(\mathbf{x}, \mathbf{t})))$ is lower bounded by

$$\begin{aligned} & \inf_{\mathbf{f} \in A_{ij}} E(\bar{V}^T(y, \mathbf{f}(\mathbf{x}, \mathbf{t})) - \bar{V}^T(y, \mathbf{f}^0(\mathbf{x}, \mathbf{t}))) + \lambda(J(\mathbf{f}) - J(\mathbf{f}^*)) \\ & - e_{VT}(\mathbf{f}^0 - \mathbf{f}^*) \geq M(i, j), \end{aligned}$$

it then follows that $M(i, j)/v(i, j) \leq 1/(8 \max\{T, 1\})$ and (4.7) directly implies (4.5).

According to Theorem 3 in [36] with $M = n^{1/2}M(i, j)$ and $v = v(i, j)$ yields that

$$\begin{aligned} I_1 & \leq \sum_{i,j:M(i,j) \leq T} 3 \exp\left(-\frac{(1-\epsilon)nM(i,j)^2}{2(4v(i,j) + M(1,j)T/3)}\right) \\ & \leq \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} 3 \exp(-C_5 n M^{2-\min(\gamma,1)}(i, j)) \\ & = \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} 3 \exp\left((-C_5 n ((2^{i-1} - \frac{1}{2})\delta_n^2 + \lambda(2^{j-1} - 1)J^*)^{2-\min(\gamma,1)})\right) \\ & \leq \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} 3 \exp\left(-C_5 n ((i - \frac{1}{2})\delta_n^{4-2\min(\gamma,1)} + (j-1)(\lambda J^*)^{2-\min(\gamma,1)})\right) \\ & \leq \frac{3 \exp(-4C_5 n (\lambda J^*)^{2-\min(\gamma,1)})}{(1 - \exp(-4C_5 n (\lambda J^*)^{2-\min(\gamma,1)}))^2}. \end{aligned} \quad (5.6)$$

Similarly, I_2 can be bounded by

$$\begin{aligned}
 I_2 &\leq \sum_{i=1}^{\infty} 3 \exp\left(-\frac{(1-\epsilon)n(M(i,0))^2}{2(4v(i,0)) + M(i,0)T/3}\right) \\
 &\leq \sum_{i=1}^{\infty} 3 \exp\left(-C_5 n (M(i,0))^{2-\min(\gamma,1)}\right) \\
 &\leq \sum_{i=1}^{\infty} 3 \exp\left(-C_5 n (2^{i-3} \delta_n^2)^{2-\min(\gamma,1)}\right) \\
 &\leq \frac{3}{1 - \exp(-4C_5 n (\lambda J^*)^{2-\min(\gamma,1)})}. \tag{5.7}
 \end{aligned}$$

Combining (5.6) and (5.7), we have $I \leq \frac{6 \exp(-4C_5 n (\lambda J^*)^{2-\min(\gamma,1)})}{(1 - \exp(-4C_5 n (\lambda J^*)^{2-\min(\gamma,1)}))^2}$. To simplify, let $Q = \exp(-4C_5 n (\lambda J^*)^{2-\min(\gamma,1)})$, by the fact that $I \leq I^{1/2} \leq 1$ and $I^{1/2}(1 - \sqrt{Q}) \leq I^{1/2}(1 - Q) \leq \sqrt{6Q} \leq 2.5\sqrt{Q}$, we have $I \leq (I^{1/2} + 2.5)\sqrt{Q} \leq 3.5\sqrt{Q}$. The desired result then follows immediately. \square

Supplementary Material

Supplement to “Sentiment analysis with covariate-assisted word embeddings”

(doi: [10.1214/21-EJS1854SUPP](https://doi.org/10.1214/21-EJS1854SUPP); .zip).

References

- [1] GENTZKOW, M., AND SHAPIRO, J. M. (2010). What drives media slant? Evidence from US daily newspapers. *Econometrica*, 78(1), 35-71. [MR2642856](#)
- [2] TUMASJAN, A., SPRENGER, T., SANDNER, P., AND WELPE, I. (2010, May). Predicting elections with twitter: What 140 characters reveal about political sentiment. *In Proceedings of the International AAAI Conference on Web and Social Media* (Vol. 4, No. 1).
- [3] YAN, X., AND BIEN, J. (2020). Rare feature selection in high dimensions. *Journal of the American Statistical Association*, 1-14.
- [4] GENKIN, A., LEWIS, D. D., AND MADIGAN, D. (2007). Large-scale Bayesian logistic regression for text categorization. *Technometrics*, 49(3), 291-304. [MR2408634](#)
- [5] TADDY, M. (2013). Multinomial inverse regression for text analysis. *Journal of the American Statistical Association*, 108(503), 755-770. [MR3174658](#)
- [6] PANG, B., LEE, L., AND VAITHYANATHAN, S. (2002). Thumbs up? Sentiment classification using machine learning techniques. *arXiv preprint cs/0205070*.
- [7] DAI, B., SHEN, X., AND WANG, J. (2020). Embedding learning. *Journal of the American Statistical Association*, 1-13.

- [8] LE, Q., AND MIKOLOV, T. (2014, June). Distributed representations of sentences and documents. *In International conference on machine learning* (pp. 1188-1196). PMLR.
- [9] PENNINGTON, J., SOCHER, R., AND MANNING, C. D. (2014, October). Glove: Global vectors for word representation. *In Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (pp. 1532-1543).
- [10] RILOFF, E., QADIR, A., SURVE, P., DE SILVA, L., GILBERT, N., AND HUANG, R. (2013, October). Sarcasm as contrast between a positive sentiment and negative situation. *In Proceedings of the 2013 conference on empirical methods in natural language processing*, (pp. 704-714).
- [11] CHEN, T., XU, R., HE, Y., AND WANG, X. (2017). Improving sentiment analysis via sentence type classification using BiLSTM-CRF and CNN. *Expert Systems with Applications*, 72, 221-230.
- [12] THELWALL, M., WILKINSON, D., AND UPPAL, S. (2010). Data mining emotion in social network communication: Gender differences in MySpace. *Journal of the American Society for Information Science and Technology*, 61(1), 190-199.
- [13] KUCUKTUNC, O., CAMBAZOGLU, B. B., WEBER, I., AND FERHATOSMANOGLU, H. (2012, February). A large-scale sentiment analysis for Yahoo! answers. *In Proceedings of the fifth ACM international conference on Web search and data mining*, (pp. 633-642).
- [14] BAMLER, R., AND MANDT, S. (2017, July). Dynamic word embeddings. *In International conference on Machine learning*, (pp. 380-389). PMLR.
- [15] DAS, R., ZAHEER, M., AND DYER, C. (2015, July). Gaussian lda for topic models with word embeddings. *In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)* (pp. 795-804).
- [16] REN, Y., WANG, R., AND JI, D. (2016). A topic-enhanced word embedding for Twitter sentiment classification. *Information Sciences*, 369, 188-198.
- [17] WANG, J., SHEN, X., SUN, Y., AND QU, A. (2016). Classification with unstructured predictors and an application to sentiment analysis. *Journal of the American Statistical Association*, 111(515), 1242-1253. [MR3561946](#)
- [18] BACCIANELLA, S., ESULI, A., AND SEBASTIANI, F. (2009, November). Evaluation measures for ordinal regression. *In 2009 Ninth international conference on intelligent systems design and applications* (pp. 283-287). IEEE.
- [19] CORTES, C., AND VAPNIK, V. (1995). Support-vector networks. *Machine learning*, 20(3), 273-297.
- [20] SHEN, X., TSENG, G. C., ZHANG, X., AND WONG, W. H. (2003). On ψ -learning. *Journal of the American Statistical Association*, 98(463), 724-734. [MR2011686](#)
- [21] ZHU, J., AND HASTIE, T. (2005). Kernel logistic regression and the import vector machine. *Journal of Computational and Graphical Statistics*, 14(1),

- 185-205. [MR2137897](#)
- [22] PEDREGOSA, F., BACH, F., AND GRAMFORT, A. (2017). On the consistency of ordinal regression methods. *Journal of Machine Learning Research*, 18, 1-35. [MR3687598](#)
- [23] WANG, S. I., AND MANNING, C. D. (2012, July). Baselines and bigrams: Simple, good sentiment and topic classification. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)* (pp. 90-94).
- [24] ZOU, H., AND YUAN, M. (2008). Composite quantile regression and the oracle model selection theory. *The Annals of Statistics*, 36(3), 1108-1126. [MR2418651](#)
- [25] WU, Y., AND LIU, Y. (2009). Variable selection in quantile regression. *Statistica Sinica*, 801-817. [MR2514189](#)
- [26] LIU, Y., AND WU, Y. (2011). Simultaneous multiple non-crossing quantile regression estimation using kernel constraints. *Journal of Nonparametric Statistics*, 23(2), 415-437. [MR2801302](#)
- [27] XU, S., DAI, B., AND WANG, J. (2021). Supplement to "Sentiment analysis with covariate-assisted word embeddings". DOI: [10.1214/21-EJS1854SUPP](#)
- [28] HSIEH, C. J., CHANG, K. W., LIN, C. J., KEERTHI, S. S., AND SUNDARARAJAN, S. (2008, July). A dual coordinate descent method for large-scale linear SVM. In *Proceedings of the 25th international conference on Machine learning* (pp. 408-415).
- [29] ANDERSEN, M., DAHL, J., LIU, Z., VANDENBERGHE, L., SRA, S., NOWOZIN, S., AND WRIGHT, S. J. (2011). Interior-point methods for large-scale cone programming. *Optimization for Machine Learning* 5583.
- [30] RAZAVIYAYN, M., HONG, M., LUO, Z. Q., AND PANG, J. S. (2014). Parallel successive convex approximation for nonsmooth nonconvex optimization. *arXiv preprint [arXiv:1406.3665](#)*.
- [31] BARTLETT, P. L., JORDAN, M. I., AND MCAULIFFE, J. D. (2006). Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473), 138-156. [MR2268032](#)
- [32] SHEN, X., AND WANG, L. (2007). Generalization error for multi-class margin classification. *Electronic Journal of Statistics*, 1, 307-330. [MR2336036](#)
- [33] HU, M., AND LIU, B. (2004, August). Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining* (pp. 168-177).
- [34] PETERS, M. E., NEUMANN, M., IYYER, M., GARDNER, M., CLARK, C., LEE, K., AND ZETTLEMOYER, L. (2018). Deep contextualized word representations. *arXiv preprint [arXiv:1802.05365](#)*.
- [35] ABADI, M., BARHAM, P., CHEN, J., CHEN, Z., DAVIS, A., DEAN, J., ... AND ZHENG, X. (2016). Tensorflow: A system for large-scale machine learning. In *12th USENIX symposium on operating systems design and implementation (OSDI 16)* (pp. 265-283).
- [36] SHEN, X., AND WONG, W. H. (1994). Convergence rate of sieve estimates. *The Annals of Statistics*, 580-615. [MR1292531](#)