

Article

# Sparse Multicategory Generalized Distance Weighted Discrimination in Ultra-High Dimensions

Tong Su <sup>1</sup>, Yafei Wang <sup>2</sup>, Yi Liu <sup>2</sup>, William G. Branton <sup>3</sup>, Eugene Asahchop <sup>3</sup>, Christopher Power <sup>3</sup>, Bei Jiang <sup>2</sup>, Linglong Kong <sup>2,\*</sup>  and Niansheng Tang <sup>1,\*</sup>

<sup>1</sup> Key Lab of Statistical Modeling and Data Analysis of Yunnan Province, Yunnan University, Kunming 650091, China; sutong\_366@sina.com

<sup>2</sup> Department of Mathematical and Statistical Sciences, University of Alberta, Edmonton, AB T6G 2G1, Canada; yafei2@ualberta.ca (Y.W.); yliu16@ualberta.ca (Y.L.); bei1@ualberta.ca (B.J.)

<sup>3</sup> Department of Medicine (Neurology), University of Alberta, Edmonton, AB T6G 2G1, Canada; wbranton@ualberta.ca (W.G.B.); asahchop@ualberta.ca (E.A.); cp9@ualberta.ca (C.P.)

\* Correspondence: lkong@ualberta.ca (L.K.); nstang@ynu.edu.cn (N.T.)

Received: 30 September 2020; Accepted: 2 November 2020; Published: 5 November 2020



**Abstract:** Distance weighted discrimination (DWD) is an appealing classification method that is capable of overcoming data piling problems in high-dimensional settings. Especially when various sparsity structures are assumed in these settings, variable selection in multicategory classification poses great challenges. In this paper, we propose a multicategory generalized DWD (MgDWD) method that maintains intrinsic variable group structures during selection using a sparse group lasso penalty. Theoretically, we derive minimizer uniqueness for the penalized MgDWD loss function and consistency properties for the proposed classifier. We further develop an efficient algorithm based on the proximal operator to solve the optimization problem. The performance of MgDWD is evaluated using finite sample simulations and miRNA data from an HIV study.

**Keywords:** high dimension; multicategory classification; DWD; sparse group lasso;  $L_2$ -consistency; proximal algorithm

## 1. Introduction

Classification problems appear in diverse practical applications, such as spam e-mail classification, disease diagnosis and drug discovery, among many others (e.g., [1–3]). In these classification problems, the goal is to predict class labels based on a given set of variables. Recent research has focused extensively on linear classification: see [4,5] for comprehensive introductions. Among many linear classification methods, support vector machines (SVMs) (see [6,7]) and distance-weighted discrimination (DWD) (see [8–10]) are two commonly used large-margin based classification methods.

Owing to the recent advent of new technologies for data acquisition and storage, classification with high dimensional features, i.e., a large number of variables, has become a ubiquitous problem in both theoretical and applied scientific studies. Typically, only a small number of instances are available, a setting we refer to as high-dimensional, low-sample size (HDLSS), as in [11]. In the HDLSS setting, a so-called “data-piling” phenomenon is observed in [8] for SVMs, occurring when projections of many training instances onto a vector normal to the separating hyperplane are nearly identical, suggesting severe overfitting. DWD was originally proposed to overcome data-piling in the HDLSS setting. In binary classification problems, linear SVMs seek a hyperplane maximizing the smallest margin for all data points, while DWD seeks a hyperplane minimizing the sum of inverse margins over all data points. Reference [8] suggests replacing the inverse margins by the  $q$ -th power of the inverse margins in a generalized DWD method; see [12] for a detailed description. Formally, for a training

data set  $\{(y_i, \mathbf{X}_i)\}_{i=1}^N$  of  $N$  observations, where  $\mathbf{X}_i \in \mathbb{R}^p$  and  $y_i \in \{-1, 1\}$ , binary generalized linear DWD seeks a proper separating hyperplane  $\{\mathbf{X} : a + \mathbf{X}^\top \mathbf{b} = 0\}$  through the optimization problem

$$\begin{aligned} \arg \max_{a, \mathbf{b}} \sum_{i=1}^N \frac{1}{d_i^q} \\ \text{s.t. } d_i = y_i (a + \mathbf{X}_i^\top \mathbf{b}) + \eta_i \geq 0, \forall i, \\ \eta_i \geq 0, \forall i, \sum_i \eta_i \leq c, \\ \|\mathbf{b}\|_2^2 = 1, \end{aligned} \quad (1)$$

where  $a$  and  $\mathbf{b}$  are the intercept and slope parameters, respectively. The slack variable  $\eta_i$  is introduced to ensure that the corresponding margin  $d_i$  is non-negative and the constant  $c > 0$  is a tuning parameter to control the overlap between classes. Problem (1) can also be written in a loss-plus-penalty form (e.g., [12]) as

$$(\hat{a}, \hat{\mathbf{b}}) = \operatorname{argmin}_{a, \mathbf{b}} \left[ \frac{1}{N} \sum_{i=1}^N \phi_q \left\{ y_i (a + \mathbf{X}_i^\top \mathbf{b}) \right\} + \lambda \|\mathbf{b}\|_2^2 \right], \quad (2)$$

where

$$\phi_q(u) = \begin{cases} 1 - u, & \text{if } u \leq Q \\ \varphi_q(u), & \text{if } u > Q, \end{cases} \quad (3)$$

with  $Q = \frac{q}{q+1}$ ,  $q > 0$  and  $\varphi_q(u) = (1 - Q)(Qu^{-1})^q$ . When  $q = 1$ , (1) becomes the standard DWD problem in [8] while problem (2) appears in [9,13].

The binary classification problem (1) is well studied. However, in many applications such as image classification [1], cancer diagnosis [2] and speech recognition [3], to name a few, problems with more than two categories are commonplace. To solve these multicategory problems with the DWD classifier, approaches based on either formulation (1) or (2) are common. One common strategy is to extend problem (1) to multiple classes by solving a series of binary problems in a one-versus-one (OVO) or one-versus-rest (OVR) method (e.g., [14]). Instead of reducing the multicategory problem to a binary one, another strategy based on problem (1) considers all classes at once. As shown in [14], this approach generally works better than the OVO and OVR methods. Based on an extension of problem (2), [15] proposes multicategory DWD, written in a loss-plus-penalty form as

$$\begin{aligned} \min_{a_k, \mathbf{b}_k} \frac{1}{N} \sum_{i=1}^N \phi_q \left( a_{y_i} + \mathbf{X}_i^\top \mathbf{b}_{y_i} \right) + \lambda \sum_{k=1}^K \|\mathbf{b}_k\|_2^2 \\ \text{s.t. } \sum_{k=1}^K a_k = 0; \sum_{k=1}^K b_{jk} = 0, \quad \forall j = 1, \dots, p, \end{aligned} \quad (4)$$

with  $y_i, k \in \{1, \dots, K\}$  and where  $a_k$  and  $\mathbf{b}_k = (b_{1k}, \dots, b_{pk})$  are the intercept and slope parameters for each category  $k$ , respectively. Although these methods can be applied to multicategory classification in the HDLSS setting, both problems (2) and (4) use the  $L_2$  penalty and do not perform feature selection. As discussed in [16], for high dimensional classification, taking all features into consideration does not work well for two reasons. First, based on prior knowledge, only a small number of variables are relevant to the classification problem: a good classifier in high dimensions should have the ability to sparsely select important variables and discard redundant ones. Second, classifiers using all available variables in high-dimensional settings may have poor classification performance.

Much of the SVM literature has considered variable selection in high-dimensional classification problems to improve performance (e.g., [17–19]). Among the DWD literature, to the best of our knowledge, only [16] considered variables selection and classification simultaneously. Wang and Zou [16] considered an  $L_1$  rather than an  $L_2$  penalty in problem (2) to improve interpretability through sparsity in the binary classification. Moreover, [16] made selections based on the strengths of input variables within individual classes but ignored the strengths of input variable groupings, thereby selecting more factors than necessary for each class. To overcome this weakness in this paper, we developed a multicategory generalized DWD method that is capable of performing variable selection and classification simultaneously. Our approach incorporates sparsity and group structure information via the sparse group lasso penalty (see [20–24]).

Although DWD is well studied, it is less popular than the SVM for binary classification, arguably for computational and theoretical reasons. For an up-to-date list of works on DWD mostly focused on the  $q = 1$  case, see [14,15]. Theoretical asymptotic properties of large-margin classifiers in high dimensional settings were studied in [25], and [26] derived an expression for asymptotic generalization error. In terms of computation, [8] solved the standard DWD problem in (1) as a second-order cone programming (SOCP) problem using a primal-dual interior-point method that is computationally expensive when  $N$  or  $p$  is large. To overcome computational bottlenecks, [12] proposed an approach based on a novel formulation of the primal DWD model in (1): this method, proposed in [12], does not scale to large data sets and requires further work. Lam et al. [27] designed a new algorithm for large DWD problems with  $q \geq 2$  and  $K = 2$  based on convergent multi-block ADMM-type methods (see [28]). Wang and Zou [16] solved the lasso-penalized binary DWD problem by combining majorization–minimization and coordinate descent methods: the lasso penalty does not directly permit a SOCP solution. In fact, solution identifiability in the generalized DWD problem with  $q > 1$  requires more constraints and remains an open research problem (see [8]). To the best of our knowledge, no work focusing on computational aspects of lasso penalized multicategory generalized DWD (MgDWD) exists. The same holds for sparse group lasso-penalized MgDWD.

The theoretical and computational contributions of this paper are as follows. First, we establish the uniqueness of the minimizer in the population form of the MgDWD problem. Second, we prove a non-asymptotic  $L_2$  estimation error bound for the sparse group lasso-regularized MgDWD loss function in the ultra-high dimensional setting under mild regularity conditions. Third, we develop a fast, efficient algorithm able to solve the sparse group lasso-penalized MgDWD problem using proximal methods.

The rest of this paper is organized as follows. In Section 2.1, we introduce the MgDWD problem with sparse group lasso penalty. In Sections 2.2 and 2.3, we establish theoretical properties of the population classifier and regularized empirical loss. We propose a computational algorithm in Section 2.4. Section 3 illustrates the finite sample performance of our method through simulation studies and a real data analysis. Proofs for major theorems are given in the Appendix A.

## 2. Methodology

### 2.1. Model Setup

We begin with some basic set-up and notation. Consider the multicategory classification problem for a random sample  $\{(y_i, \mathbf{X}_i)\}_{i=1}^N$  of  $N$  independent and identically distributed (i.i.d.) observations from some distribution  $\mathbb{P}(y, \mathbf{X})$ . Here,  $y$  is the categorical response taking values in  $\mathcal{Y} = \{1, \dots, K\}$ , and  $\mathbf{X} = (x_1, \dots, x_p)^\top \in \mathcal{X} \subset \mathbb{R}^p$  is the covariate vector. We wish to obtain a proper separating hyperplane  $\{\mathbf{X} \in \mathcal{X} | a_k + \mathbf{X}^\top \mathbf{b}_k = 0\}$  for each category  $k \in \mathcal{Y}$ , where  $a_k$  and  $\mathbf{b}_k = (b_{1k}, \dots, b_{pk})^\top$  are intercept and slope parameters, respectively.

In this paper, we consider MgDWD with sparse group lasso regularization. That is, we estimate a classification boundary by solving the constrained optimization problem

$$\begin{aligned} \min_{a_k, \mathbf{b}_k} & \frac{1}{N} \sum_{i=1}^N \phi_q(a_{y_i} + \mathbf{X}_i^\top \mathbf{b}_{y_i}) + \lambda_1 \sum_{k=1}^K \sum_{j=1}^p |b_{jk}| + \lambda_2 \sum_{j=1}^p \sqrt{\sum_{k=1}^K b_{jk}^2} \\ \text{s.t.} & \sum_{k=1}^K a_k = 0; \sum_{k=1}^K b_{jk} = 0, \quad \forall j = 1, \dots, p, \end{aligned} \tag{5}$$

where  $\phi_q$  is as defined in (3).

To approach this problem, we apply the concept of a “margin vector” to extend the definition of a (binary) margin to the multicategory case. Denote the margin vector of an observation  $\mathbf{X}_i$  as  $\mathbf{F}_i = (f_{i1}, \dots, f_{iK})^\top$ , with  $f_{ik} = a_k + \mathbf{X}_i^\top \mathbf{b}_k$  satisfying  $\sum_{k=1}^K f_{ik} = 0$ . Let  $\mathbf{E}_i = (e_{i1}, \dots, e_{iK})^\top$  be the class indicator vector with  $e_{ik} = \mathbb{1}\{y_i = k\}$ . The multicategory margin of the data point  $(y_i, \mathbf{X}_i)$  is then given as  $f_{iy_i} = a_{y_i} + \mathbf{X}_i^\top \mathbf{b}_{y_i} = \mathbf{E}_i^\top \mathbf{F}_i$ . Therefore, the MgDWD loss can be rewritten as

$$\phi_q(a_{y_i} + \mathbf{X}_i^\top \mathbf{b}_{y_i}) = \phi_q(\mathbf{E}_i^\top \mathbf{F}_i) = \mathbf{E}_i^\top \phi_q(\mathbf{F}_i) = \sum_{k=1}^K \mathbb{1}\{y_i = k\} \phi_q(a_k + \mathbf{X}_i^\top \mathbf{b}_k). \tag{6}$$

Based on (6), Lemma 1 describes the Fisher consistency of the MgDWD loss.

**Lemma 1.** Given  $\mathbf{X} = \mathbf{u}$ , the minimizer of the conditional expectation of (6) is  $\tilde{\mathbf{F}}(\mathbf{u}) = (\tilde{f}_1(\mathbf{u}), \dots, \tilde{f}_K(\mathbf{u}))^\top$ , satisfying

$$\operatorname{argmax}_{k \in \mathcal{Y}} \tilde{f}_k(\mathbf{u}) = \operatorname{argmax}_{k \in \mathcal{Y}} \Pr\{y = k | \mathbf{X} = \mathbf{u}\},$$

where

$$\tilde{f}_k(\mathbf{u}) = \begin{cases} Q \sqrt{\frac{\Pr\{y = k | \mathbf{X} = \mathbf{u}\}}{\Pr\{y = k_* | \mathbf{X} = \mathbf{u}\}}}, & k \neq k_* \\ -Q \sum_{l \neq k_*} \sqrt{\frac{\Pr\{y = l | \mathbf{X} = \mathbf{u}\}}{\Pr\{y = k_* | \mathbf{X} = \mathbf{u}\}}}, & k = k_*. \end{cases}$$

and  $k_* = \operatorname{argmin}_{k \in \mathcal{Y}} \Pr\{y = k | \mathbf{X} = \mathbf{u}\}$ .

Consequently,  $\tilde{f}_k(\mathbf{u})$  can be treated as an effective proxy of  $\Pr\{y = k | \mathbf{X} = \mathbf{u}\}$  and, for any new observation  $\mathbf{X}_*$ , a reasonable prediction of its label  $y_*$  is

$$\hat{y}_* = \operatorname{argmax}_{k \in \mathcal{Y}} \{a_k + \mathbf{X}_*^\top \mathbf{b}_k\}.$$

Speaking to the sparse group lasso (SGL) regularization in (5), the  $L_1$  penalty encourages an element-wise sparse estimator that selects important variables for each category, indicated by  $\hat{b}_{jk} \neq 0$ . Assuming that parameters in different categories share the same information, we use an  $L_2$  penalty to encourage a group-wise sparsity structure that removes covariates that are irrelevant across all categories, that is, where  $\hat{\boldsymbol{\beta}}_j = (b_{1j}, \dots, b_{Kj})^\top = \mathbf{0}$ . Specifically, let  $x_j = (x_{1j}, \dots, x_{Nj})^\top$  and  $\mathbf{B} = (b_{jk}) \in \mathbb{R}_{jk}^{p \times K}$ , where the  $k$ -th column  $\mathbf{b}_k$  is the slope vector for the category label  $k$  and the  $j$ -th row  $\boldsymbol{\beta}_j^\top$  is the group coefficient for the variable  $x_j$ . If  $x_j$  is noise in the classification problem or is not

relevant to category label  $k$ , then the entry  $b_{jk}$  of  $\mathbf{B}$  should be shrunk to exactly zero. The SGL penalty of (5) can be written as a convex combination of the lasso and group lasso penalties in terms of  $\beta_j$  as

$$\lambda_1 \sum_{k=1}^K \sum_{j=1}^p |b_{jk}| + \lambda_2 \sum_{j=1}^p \sqrt{\sum_{k=1}^K b_{jk}^2} = \lambda \sum_{j=1}^p \{ \tau \|\beta_j\|_1 + (1 - \tau) \|\beta_j\|_2 \}, \tag{7}$$

where  $\lambda > 0$  is the scale of the penalty and  $\tau \in [0, 1]$  tunes the propensity between the element-wise and group-wise sparsity structure.

### 2.2. Population MgDWD

In this subsection, some basic results pertaining to unpenalized population MgDWD are given. These results are necessary for further theoretical analysis.

Denote the marginal probability mass of  $y$  as  $\Pr(y = k) = \pi_k$  with  $\pi_k > 0$  and  $\sum_{k=1}^K \pi_k = 1$ , and the conditional probability density functions of  $\mathbf{X}$  given  $y = k$  by  $g(\mathbf{X} \mid y = k) = g_k(\mathbf{X})$ . Let  $\Theta = (\theta_1, \dots, \theta_K)$  be the collection of coefficient vectors  $\theta_k = (a_k, \mathbf{b}_k^\top)^\top$  for all labels and  $\mathbf{Z} = (\mathbf{1}, \mathbf{X}^\top)^\top$ . The population version of the MgDWD problem in (6) is

$$\mathcal{L}(\boldsymbol{\vartheta}) = \mathbb{E}\{\mathbb{I}(\mathcal{Y})^\top \phi_q(\Theta^\top \mathbf{Z})\} = \sum_{k=1}^K \pi_k \int_{\mathcal{X}} \phi_q(\mathbf{Z}^\top \theta_k) g_k(\mathbf{x}) d\mathbf{x}, \tag{8}$$

where  $\boldsymbol{\vartheta} = \text{vec}\{\Theta\}$  is the vectorization of the matrix  $\Theta$  and  $\mathbb{I}(\mathcal{Y}) = (\mathbf{1}\{y = 1\}, \dots, \mathbf{1}\{y = K\})^\top$  is a random vector. Denote the true parameter value  $\boldsymbol{\vartheta}^*$  as a minimizer of the population MgDWD problem, namely,

$$\boldsymbol{\vartheta}^* \in \underset{\boldsymbol{\vartheta} \in \mathcal{C}}{\text{argmin}} \mathcal{L}(\boldsymbol{\vartheta}),$$

where  $\mathcal{C} = \{\boldsymbol{\vartheta} \in \mathbb{R}^{K(p+1)} \mid \mathbf{C}\boldsymbol{\vartheta} = \mathbf{0}_K\}$  is the set of sum-constrained  $\boldsymbol{\vartheta}$  with  $\mathbf{C} = \mathbf{1}_K^\top \otimes \mathbf{I}_{p+1}$ , where  $\otimes$  denotes the Kronecker product.

To facilitate our theoretical analysis, we first define the gradient vector and Hessian matrix of the population MgDWD loss function. We then introduce some regularity conditions necessary to derive theoretical properties of this problem. Let  $\text{diag}\{v\}$  be a diagonal matrix constructed from the vector  $v$ , and let  $\circ$  and  $\oplus$  be the Hadamard product and the direct matrix sum, respectively. Denote the gradient vector of the population MgDWD loss function (8) as

$$\mathcal{S}(\boldsymbol{\vartheta}) = \mathbb{E}\{\{\mathbb{I}(\mathcal{Y}) \circ \phi'_q(\Theta^\top \mathbf{Z})\} \otimes \mathbf{Z}\} = \text{vec}(\mathcal{S}_1, \dots, \mathcal{S}_K),$$

with

$$\mathcal{S}_k = \mathbb{E}\{\mathbf{1}\{y = k\} \phi'_q(\mathbf{Z}^\top \theta_k) \mathbf{Z}\} = \pi_k \int_{\mathcal{X}} \phi'_q(\mathbf{Z}^\top \theta_k) \mathbf{Z} g_k(\mathbf{X}) d\mathbf{X},$$

and its Hessian matrix as

$$\mathcal{H}(\boldsymbol{\vartheta}) = \mathbb{E}\left\{\text{diag}\{\mathbb{I}(\mathcal{Y}, \mathcal{X}) \circ \phi''_q(\Theta^\top \mathbf{Z})\} \otimes (\mathbf{Z}\mathbf{Z}^\top)\right\} = \bigoplus_{k=1}^K \mathcal{H}_k,$$

where  $\phi''_q$  denotes the second derivative of the function  $\phi_q$ ;  $\mathbb{I}(\mathcal{Y}, \mathcal{X}) = \mathbb{I}(\mathcal{Y}) \circ \mathbb{I}(\mathcal{X})$  is a random vector with  $\mathbb{I}(\mathcal{X}) = (\mathbf{1}\{\mathbf{X} \in \mathcal{X}_1\}, \dots, \mathbf{1}\{\mathbf{X} \in \mathcal{X}_k\})^\top$  and  $\mathcal{X}_k = \{\mathbf{X} \in \mathcal{X} \mid \mathbf{Z}^\top \theta_k > Q\}$ ; and

$$\mathcal{H}_k = \mathbb{E}\{\mathbf{1}\{y = k, \mathbf{X} \in \mathcal{X}_k\} \phi''_q(\mathbf{Z}^\top \theta_k) \mathbf{Z}\mathbf{Z}^\top\} = \pi_k \int_{\mathcal{X}_k} \phi''_q(\mathbf{Z}^\top \theta_k) \mathbf{Z}\mathbf{Z}^\top g_k(\mathbf{X}) d\mathbf{X}.$$

The block structure of  $\mathcal{H}(\boldsymbol{\theta})$  implies a parallel relationship between each category. The relationship between the  $\boldsymbol{\theta}_k$  is reflected by the sum-to-zero constraint in the definition of  $\mathcal{C}$ .

We assume the following regularity conditions.

(C1) The densities of  $\mathbf{X}$  given  $y = k \in \mathcal{Y}$ , i.e., the  $g_k(\mathbf{X})$ , are continuous and have finite second moments.

(C2)  $0 < \Pr\{\mathbf{X} \in \mathcal{X}_k^* | y = k\} < 1$  for all  $k \in \mathcal{Y}$ , where  $\mathcal{X}_k^* = \{\mathbf{X} \in \mathcal{X} | \mathbf{Z}^\top \boldsymbol{\theta}_k^* > Q\}$ .

(C3)  $\text{Var}\{\mathbf{X} | \mathbf{X} \in \mathcal{X}_k^*, y = k\} \succ \mathbf{O}_p$  for all  $k \in \mathcal{Y}$ .

**Remark 1.** Condition (C1) ensures that  $\mathcal{L}$ ,  $\mathcal{S}$  and  $\mathcal{H}$  are well defined and continuous in  $\boldsymbol{\theta}$ . For the theoretically optimal hyperplane  $\{\mathbf{X} \in \mathcal{X} | \mathbf{Z}^\top \boldsymbol{\theta}_k^* = 0\}$ , the case with  $\boldsymbol{\theta}_k^* = \mathbf{0}_{p+1}$  leaves  $\mathcal{X}$  useless for classification. On the other hand, when  $\mathbf{a}_k^* \neq 0$  and  $\mathbf{b}_k^* = \mathbf{0}_p$ , the hyperplane is the empty set and is similarly meaningless. Condition (C2) is proposed to avoid the case where  $\mathbf{b}_k^* = \mathbf{0}_p$  so that  $\boldsymbol{\theta}^*$  always contains information relevant to the classification problem. For bounded random variables, condition (C2) should be assumed with caution. Condition (C3) implies the positive definiteness of  $\mathcal{H}(\boldsymbol{\theta}^*)$ .

By convexity and the second-order Lagrange condition, the following theorem shows that the local minimizer of the population MgDWD problem exists and is unique.

**Theorem 1.** Under the regularity conditions (C1)-(C3), the true parameter  $\boldsymbol{\theta}^* \in \mathcal{C}$  is the unique minimizer of  $\mathcal{L}(\boldsymbol{\theta})$  with  $\mathbf{b}_k^* \neq \mathbf{0}_p$ , and

$$\mathcal{L}(\boldsymbol{\theta}^*) = \sum_{k=1}^K A(k, q) \pi_k,$$

with  $0 \leq u(k, q) \leq A(k, q) \leq v(k, q) \leq 1$ , where

$$\begin{aligned} A(k, q) &= 1 - \mathbb{E}\left\{\mathbb{1}\{\mathbf{X} \in \mathcal{X}_k^*\} \left\{1 - \left(\frac{Q}{\mathbf{Z}^\top \boldsymbol{\theta}_k^*}\right)^q\right\} \mid y = k\right\}, \\ u(k, q) &= \Pr\{\mathbf{X} \notin \mathcal{X}_k^* | y = k\} + Q^{2q} \Pr\{Q < \mathbf{Z}^\top \boldsymbol{\theta}_k^* \leq Q^{-1} | y = k\}, \\ v(k, q) &= \Pr\{\mathbf{Z}^\top \boldsymbol{\theta}_k^* \leq 1 \mid y = k\} + \inf_{\epsilon > 0} \left(\frac{Q}{1 + \epsilon}\right)^q \Pr\{\mathbf{Z}^\top \boldsymbol{\theta}_k^* > 1 + \epsilon \mid y = k\}. \end{aligned}$$

The bounds in Theorem 1 show how  $q$  affects the loss function  $\mathcal{L}(\boldsymbol{\theta}^*)$ . The upper bound  $v(k, q)$  is a decreasing function of  $q$  with

$$\lim_{q \rightarrow 0} v(k, q) = 1 \text{ and } \lim_{q \rightarrow \infty} v(k, q) = \Pr\{\mathbf{Z}^\top \boldsymbol{\theta}_k^* \leq 1 \mid y = k\}.$$

In the lower bound  $u(k, q)$ , the first term  $\Pr\{\mathbf{X} \notin \mathcal{X}_k^* | y = k\}$  is an increasing function of  $q$  and the last term  $Q^{2q} \Pr\{Q < \mathbf{Z}^\top \boldsymbol{\theta}_k^* \leq Q^{-1} | y = k\}$  is a decreasing function of  $q$ , with

$$\lim_{q \rightarrow 0} u(k, q) = 1 \text{ and } \lim_{q \rightarrow \infty} u(k, q) = \Pr\{\mathbf{Z}^\top \boldsymbol{\theta}_k^* \leq 1 \mid y = k\}.$$

Consequently, for the given population  $\mathbb{P}(y, \mathbf{X})$ , a larger  $q$  encourages the population MgDWD estimator to focus more on the regions  $\{\mathbf{X} \notin \mathcal{X}_k, y = k\}$  that correspond to misclassifications. As a result, the estimator’s performance will be similar to the hinge loss as  $q \rightarrow \infty$ . Setting  $q$  too small will lead to an ineffective classifier due to the unreasonable penalty placed on the well classified region  $\{\mathbf{X} \in \mathcal{X}_k, y = k\}$ . This variation in the lower bound with respect to  $q$  provides a necessary condition for the existence of an optimal  $q$ .

**Remark 2.** The explicit relationship between  $q$  and  $\boldsymbol{\theta}^*$  is complicated. While it may be more desirable to prove that a greater value of  $q$  results in a smaller value of the loss function  $\mathcal{L}(\boldsymbol{\theta})$ , there is no explicit formula for the optimal value  $\boldsymbol{\theta}^*$  in terms of  $q$ .

### 2.3. Estimator Consistency

Under the unpenalized framework presented in the previous subsection, all covariates will contribute to the classification task for each category: this scenario may lead to a classifier that overfits to the training data set. In this subsection, we study the consistency of the estimator for (5) in ultra-high dimensional settings.

To achieve structural sparsity in the estimator, the regularization parameter  $\lambda$  in (7) must be large enough to dominate the gradient of the empirical MgDWD loss evaluated at the theoretical minimizer  $\boldsymbol{\theta}^* = \text{vec}\{\boldsymbol{\Theta}^*\}$  with high probability. On the other hand,  $\lambda$  should also be as small as possible to reduce the bias incurred by the SGL regularization term

$$P(\boldsymbol{\beta}) = \sum_{j=1}^p \tau \|\boldsymbol{\beta}_j\|_1 + (1 - \tau) \|\boldsymbol{\beta}_j\|_2.$$

Lemma 2 provides a suitable choice of  $\lambda$  under the following assumption.

(A1) The predictors  $\mathbf{X} = (x_1, \dots, x_p) \in \mathbb{R}^p$  are independent sub-Gaussian random vectors satisfying  $\mathbb{E}\mathbf{X} = \mathbf{0}_p$ , and where  $\text{Var}(\mathbf{X}) = \boldsymbol{\Sigma}$ , there exists a constant  $\kappa > 0$  such that for any  $\gamma \in \mathbb{R}^p$ ,  $\mathbb{E} \exp(\gamma^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{X}) \leq \exp(\|\gamma\|_2^2 \kappa^2 / 2)$ . From here on, we define  $\zeta_1^2$  as the largest eigenvalue of  $\boldsymbol{\Sigma}$ .

**Lemma 2.** Denote  $\mathbf{S}(\boldsymbol{\theta}^*) = (\mathbf{I}_K \otimes \mathbf{Z}^\top) \text{diag}(\text{vec}\{\mathbf{E}\}) \text{vec}\{\phi'_q(\mathbf{Z}\boldsymbol{\theta}^*)\}$ , where  $\mathbf{E} = (\mathbf{E}_1, \dots, \mathbf{E}_N)^\top$ ,  $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_N)^\top$  with  $\mathbf{Z}_i = (1, \mathbf{X}_i^\top)^\top$ , and  $\mathbf{I}_K$  is the identity matrix of size  $K$ . Under condition (A1),

$$\tilde{P}\{\mathbf{P}\mathbf{S}(\boldsymbol{\theta}^*)\} \leq \tau \Lambda_1 + (1 - \tau) \Lambda_2$$

with probability at least  $1 - 2(Kp)^{1-c_1^2} - p^{1-c_2^2}$ , where

$$\begin{aligned} \mathbf{P} &= (\mathbf{I}_K - K^{-1} \mathbf{1}_K \mathbf{1}_K^\top) \otimes \mathbf{I}_{p+1}, \\ \Lambda_1 &= \max\{\zeta_1 \kappa, 1\} \left(1 - \frac{1}{K}\right) \sqrt{\frac{2 \log(pK)}{N}}, \\ \Lambda_2 &= \max\{2\sqrt{2}\zeta_1 \kappa, 1\} \left\{ c_2 \sqrt{\left(1 - \frac{1}{K}\right) \frac{2 \log(p)}{N}} + \sqrt{\frac{K-1}{N}} \right\}, \end{aligned}$$

for constants  $c_1, c_2 > 1$ .

It is difficult to obtain a closed form for the conjugate of the SGL penalty, say,  $\bar{P}(v) = \sup_{\mathbf{u} \in \mathcal{C} \setminus \{0\}} \frac{\langle \mathbf{u}, v \rangle}{P(\mathbf{u})}$ . Instead, we use a regularized upper bound  $\tilde{P}(v) \geq \bar{P}(v)$ . Based on Lemma 2, we propose a theoretical tuning parameter value

$$\lambda = c_0 \sqrt{\frac{\log(pK)}{N}}, \tag{9}$$

where  $c_0$  is some given constant satisfying  $\lambda > \tau \Lambda_1 + (1 - \tau) \Lambda_2$ .

Before we can derive an error bound for the estimator in (5), we impose two additional assumptions.

(A2) For the true parameter value  $\boldsymbol{\theta}^*$ , there is a  $(s_e, s_g)$ -sparse structure in the coefficients  $\mathbf{B}^*$  with element-wise and group-wise support sets

$$\mathcal{E} = \{(j, k) \in \{1, \dots, p\} \times \{1, \dots, K\} | b_{jk}^* \neq 0\} \text{ and } \mathcal{G} = \{j \in \{1, \dots, p\} | \boldsymbol{\beta}_j^* \neq \mathbf{0}_K\}$$

with cardinality  $|\mathcal{E}| = s_e$  and  $|\mathcal{G}| = s_g$ , respectively.

(A3) There exist some positive constants  $\zeta_2$  and  $\zeta_3$  such that

$$\zeta_2^2 = \max_{\gamma \in \mathcal{V}} \frac{\|\text{diag}\{\text{vec}(\mathbf{E}^\top)\}(\mathbf{Z} \otimes \mathbf{I}_K)\gamma\|_2^2}{N\|\gamma\|_2^2} \text{ and } \zeta_3^2 = \min_{\gamma \in \mathcal{W}} \frac{\gamma^\top \mathcal{H}(\boldsymbol{\theta}^*)\gamma}{\gamma^\top \gamma}$$

with  $\mathcal{V} = \{\mathbf{v} \in \mathbb{R}^{K(p+1)} | 0 < \|\mathbf{v}\|_0 \leq s_e + K\}$  and

$$\mathcal{W} = \left\{ \boldsymbol{\delta} \in \mathbb{R}^{K(p+1)} \mid \frac{\tau}{1-\tau} \|\boldsymbol{\delta}_{\mathcal{E}_+}\|_1 + \sum_{j \in \mathcal{G}_+} \|\boldsymbol{\delta}_j\|_2 \geq C_0 \left( \frac{\tau}{1-\tau} \|\boldsymbol{\delta}_{\mathcal{E}^c}\|_1 + \sum_{j \notin \mathcal{G}} \|\boldsymbol{\delta}_j\|_2 \right) \right\},$$

where  $C_0 = \frac{(c_0-1)}{(c_0+1)}$ ,  $\mathcal{E}^c$  is the complement of  $\mathcal{E}$ ,  $\mathcal{E}_+ = \mathcal{E} \cup \{l = 1 + (k-1)(p+1) | k = 1, \dots, K\}$ , and  $\mathcal{G}_+ = \mathcal{G} \cup \{0\}$ .

Under the choice of  $\lambda$  given in (9), we show the  $L_2$ -consistency of the estimator in (5).

**Theorem 2.** Suppose that conditions (A1)-(A3) hold. Then with  $\lambda = c_0 \sqrt{\frac{\log(pK)}{N}}$  in (5), we have that

$$\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 \leq \left\{ C_1 \sqrt{s_e + K} + C_2 \sqrt{s_g + 1} \right\} \sqrt{\frac{\log(pK)}{N}}$$

with probability at least  $1 - 2(Kp)^{2(s_e+K)(1-c_3^2)}$ , where  $C_1 = 2\zeta_3^{-2}\{c_0\tau + (\sqrt{2} + 2c_3)\zeta_2\}$  and  $C_2 = 2\zeta_3^{-2}c_0(1-\tau)$ .

**Remark 3.** The sub-Gaussian distribution assumption (A1) is common in high-dimensional scenarios. This assumption characterizes the tail behavior of a collection of random variables including Gaussian, Bernoulli, and bounded variables as special cases. Assumption (A2) describes structural sparsity at two levels. The element-wise size  $s_e < p$  is the size of the underlying generative model, and the group-wise size  $s_g < pK$  is the size of the signal covariate set. Both  $s_e$  and  $s_g$  are allowed to depend on the sample size  $N$ . As a result, the dimension  $p$  is allowed to increase with the sample size  $N$ . Assumption (A3) guarantees that eigenvalues are positive in this sparse scenario.

**Remark 4.** In practice, the tuning parameters  $\lambda$  and  $\tau$  in (7) are commonly chosen by  $M$ -fold cross validation. That is, we choose the pair  $(\tau, \lambda)$  with the highest prediction accuracy among the sub-data sets  $\mathcal{D}_m$ , specifically,

$$CV(\tau, \lambda) = \sum_{m=1}^M \sum_{i \in \mathcal{D}_m} \mathbb{1}\{y_i \neq \hat{y}_i(\tau, \lambda)\}$$

where  $\hat{y}_i(\tau, \lambda) = \underset{k \in \mathcal{W}}{\text{argmax}} \mathbf{Z}_i^\top \hat{\boldsymbol{\theta}}_k(\tau, \lambda)$ .

### 2.4. Computational Algorithm

In this section, we propose an efficient algorithm to solve problem (5). Our approach uses the proximal algorithm (see [29]) for solving high dimensional regularization problems. In two main steps, this approach obtains a solution to the constrained optimization problem by applying the proximal operator to the solution to the unconstrained problem.

Since regularization is not needed for the intercept terms  $\boldsymbol{\alpha} = (a_1, \dots, a_K)^\top$ , it can be separated from the coefficients in  $\mathbf{B}$ . The empirical MgDWD loss of (8) is given by

$$L(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \mathbf{E}_i^\top \boldsymbol{\phi}_q(\mathbf{F}_i) = \frac{1}{N} \text{tr} \left\{ \mathbf{E} \boldsymbol{\phi}_q(\mathbf{F}^\top) \right\} = \frac{1}{N} \text{vec}\{\mathbf{E}^\top\}^\top \text{vec}\{\boldsymbol{\phi}_q(\mathbf{F}^\top)\}$$



where  $\mathbf{F} = (f_{ik})_{N \times K} = \mathbf{Z}\Theta = \mathbf{1}_N \boldsymbol{\alpha}^\top + \mathbf{X}\mathbf{B}$ . Various properties of the loss function  $L(\boldsymbol{\theta})$  follow below.

**Lemma 3.** *The loss function  $L(\boldsymbol{\theta})$  has Lipschitz continuous partial derivatives. In particular, for  $\mathbf{S}(\boldsymbol{\alpha}) = \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\alpha}} = \frac{1}{N} \{\mathbf{E} \circ \phi'_q(\mathbf{F})\}^\top \mathbf{1}_N$  and any  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^K$ , we have that*

$$\|\mathbf{S}(\mathbf{u}) - \mathbf{S}(\mathbf{v})\|_2 \leq \sqrt{\frac{n_{\max}(q+1)^2}{Nq}} \|\mathbf{u} - \mathbf{v}\|_2,$$

where  $n_{\max}$  is the largest group sample size. For  $\mathbf{S}(\mathbf{B}) = \frac{\partial L(\boldsymbol{\theta})}{\partial \mathbf{B}} = \frac{1}{N} \{\mathbf{E} \circ \phi'_q(\mathbf{F})\}^\top \mathbf{X}$  and any  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{p \times K}$ , we have that

$$\|\text{vec}\{\mathbf{S}(\mathbf{U}) - \mathbf{S}(\mathbf{V})\}\|_2 \leq \frac{\max_k \|\text{diag}(\mathbf{e}_k)\mathbf{X}\|_2^2 (q+1)^2}{Nq} \|\text{vec}\{\mathbf{U} - \mathbf{V}\}\|_2,$$

where  $\mathbf{e}_k$  is the  $k$ -th column of  $\mathbf{E}$  and indicates the observations belonging to the  $k$ -th group.

Hence, following the majorization–minimization scheme, we can majorize the empirical MgDWD loss  $L(\boldsymbol{\theta})$  by a quadratic function, that is,

$$\begin{aligned} L(\boldsymbol{\theta}) \leq & L(\boldsymbol{\theta}_*) + \mathbf{S}(\boldsymbol{\alpha}_*)^\top (\boldsymbol{\alpha} - \boldsymbol{\alpha}_*) + \frac{L_\alpha}{2} \|\boldsymbol{\alpha} - \boldsymbol{\alpha}_*\|_2^2 \\ & + \text{vec}\{\mathbf{S}(\mathbf{B})\}^\top \text{vec}\{\mathbf{B} - \mathbf{B}_*\} + \frac{L_{\mathbf{B}_*}}{2} \|\text{vec}\{\mathbf{B} - \mathbf{B}_*\}\|_2^2, \end{aligned}$$

for some  $\boldsymbol{\theta}_* = \text{vec}\{(\boldsymbol{\alpha}_*, \mathbf{B}_*)^\top\}$ , where  $L_\alpha$  and  $L_{\mathbf{B}}$  denote the Lipschitz constants in Lemma 3. Instead of minimizing  $L(\boldsymbol{\theta})$  directly, we apply gradient descent to minimize its surrogate upper bound function. The gradient descent updates are given by

$$\boldsymbol{\alpha}_* = \boldsymbol{\alpha} - \frac{q(q+1)^{-2}}{\sqrt{n_{\max}N}} \{\mathbf{E} \circ \phi'_q(\mathbf{F})\}^\top \mathbf{1}_N, \tag{10}$$

$$\mathbf{B}_* = \mathbf{B} - \frac{q(q+1)^{-2}}{\max_k \|\text{diag}(\mathbf{e}_k)\mathbf{X}\|_2^2} \{\mathbf{E} \circ \phi'_q(\mathbf{F})\}^\top \mathbf{X}. \tag{11}$$

Next, we address the problem’s constraints and regularization simultaneously by applying the proximal operator. For  $\boldsymbol{\alpha}_*$ , it is clear that

$$\boldsymbol{\alpha}_{\text{new}} = \underset{\boldsymbol{\alpha}^\top \mathbf{1}_K = 0}{\text{argmin}} \|\boldsymbol{\alpha} - \boldsymbol{\alpha}_*\|_2^2 = \mathbf{P}_K \boldsymbol{\alpha}_*, \tag{12}$$

where  $\mathbf{P}_K = \mathbf{I}_K - k^{-1} \mathbf{1}_K \mathbf{1}_K^\top$ . For  $\mathbf{B}_* = (\boldsymbol{\beta}_{1*}, \dots, \boldsymbol{\beta}_{p*})^\top$ , the minimization problem can be expressed as

$$\begin{aligned} \mathbf{B}_{\text{new}} &= \underset{\mathbf{B} \mathbf{1}_K = \mathbf{0}_p}{\text{argmin}} \frac{1}{2} \|\text{vec}\{\mathbf{B} - \mathbf{B}_*\}\|_2^2 + \frac{\lambda_1}{L_{\mathbf{B}}} \|\text{vec}\{\mathbf{B}\}\|_1 + \frac{\lambda_2}{L_{\mathbf{B}}} \|\text{vec}\{\mathbf{B}\}\|_{1,2} \\ &= \underset{\mathbf{B} \mathbf{1}_K = \mathbf{0}_p}{\text{argmin}} \sum_{j=1}^p \frac{1}{2} \|\boldsymbol{\beta}_j - \boldsymbol{\beta}_{j*}\|_2^2 + \frac{\lambda_1}{L_{\mathbf{B}}} \|\boldsymbol{\beta}_j\|_1 + \frac{\lambda_2}{L_{\mathbf{B}}} \|\boldsymbol{\beta}_j\|_2, \end{aligned} \tag{13}$$

which implies that we can implement minimization for  $p$  groups in parallel. The following theorem provides the solution to (13).

**Theorem 3.** Let  $\rho_1, \rho_2 \geq 0$  and  $\beta_* \in \mathbb{R}^K$ . Then the constrained regularization problem

$$\min_{\beta \in \mathbb{R}^K} \frac{1}{2} \|\beta - \beta_*\|_2^2 + \rho_1 \|\beta\|_1 + \rho_2 \|\beta\|_2$$

$$\text{s.t. } \beta^\top \mathbf{1}_K = 0$$

has a solution of the form

$$\beta^* = \left\{ 1 - \frac{\rho_2}{\|\mathbf{P}_K(\beta_* - \rho_1 \mathbf{u})\|_2} \right\}_+ \mathbf{P}_K(\beta_* - \rho_1 \mathbf{u}) \tag{14}$$

for some  $\mathbf{u} \in \partial \|\beta\|_1$ .

In the special case with  $\rho_2 = 0$ , the constrained regularization problem in Theorem 3 reduces to the constrained lasso problem with solution  $\tilde{\beta}^* = \mathbf{P}_K(\beta_* - \rho_1 \mathbf{u})$ . Combined with (14), the proximal operator  $\mathcal{U}$ , given by

$$\beta^* = \mathcal{U}(\tilde{\beta}^*, \rho_2) = \left\{ 1 - \frac{\rho_2}{\|\tilde{\beta}^*\|_2} \right\}_+ \tilde{\beta}^*, \tag{15}$$

can be introduced to realize the group sparsity of  $\tilde{\beta}^*$ .

For the standard lasso problem, the subgradient  $\mathbf{u}$  has a closed form given by  $\tilde{\beta}^* = \beta_* - \rho_1 \mathbf{u} = \mathcal{S}(\beta_*, \rho_1)$ , with  $\mathcal{S}(u, v) = \text{sign}(u)(|u| - v)_+$ . However, under the constraint on  $\tilde{\beta}^*$ , the naive solution  $\mathbf{P}_K \mathcal{S}(\beta_*, \rho_1)$  is misleading in that it satisfies the constraint but does not achieve shrinkage, let alone loss function minimization. The term  $\mathbf{P}_K \mathbf{u}$  is suggestive of the intersection between the subdifferential set  $\partial \|\beta\|_1$  and the constraint set  $\{\beta \in \mathbb{R}^K \mid \beta^\top \mathbf{1}_K = 0\}$ ; in this sense,  $\tilde{\beta}^*$  might not have a closed form. Here we consider using coordinate descent to solve the constrained lasso problem. For some fixed coordinate  $m$ , since  $\beta^\top \mathbf{1}_K = 0$ , we have that  $b_m = -\sum_{l \neq m} b_l$ . Rewriting the objective function of the lasso-constrained problem in a coordinate-wise form, we obtain

$$\sum_{l=1}^K \frac{1}{2} (b_l - b_{l*})^2 + \rho_1 |b_l| = \left( b_k - \frac{(b_{k*} - b_{m*})}{2} + \frac{1}{2} \sum_{l \neq k, m} b_l \right)^2 + \rho_1 \left\{ |b_k| + \left| b_k + \sum_{l \neq k, m} b_l \right| \right\}$$

$$+ \frac{1}{4} \left( b_{k*} + b_{m*} + \sum_{l \neq k, m} b_l \right)^2 + \sum_{l \neq k, m} \frac{1}{2} (b_l - b_{l*})^2 + \rho_1 |b_l|. \tag{16}$$

Next, Theorem 4 provides the solution to the optimization problem (16).

**Theorem 4.** Suppose that  $t, s \in \mathbb{R}$  and  $\rho \geq 0$ . Then the regularization problem

$$\min_{b \in \mathbb{R}} \frac{1}{2} (b - t)^2 + \rho \{|b| + |b + s|\}$$

has solution

$$b^* = \begin{cases} t, & |t| < C(s, t) \\ -C(s, t), & C(s, t) \leq |t| \leq C(s, t) + 2\rho \\ \text{sign}(t)(|t| - 2\rho), & |t| > C(s, t) + 2\rho \end{cases}$$

$$= t - \mathcal{S}(t, C(s, t)) + \mathcal{S}\{\mathcal{S}(t, C(s, t)), 2\rho\},$$

where  $C(s, t) = \frac{1 - \text{sign}(s)\text{sign}(t)}{2} |s|$ .

By Theorem 4, given some  $m \in \{1, \dots, K\}$ , the coordinate-wise minimizer for any  $k \neq m$  can be expressed as the proximal operator

$$b_k = \mathcal{T}(t, s, \rho_1) = t - \mathcal{S}(t, C(s, t)) + \mathcal{S}\{\mathcal{S}(t, C(s, t)), \rho_1\}, \tag{17}$$

with  $s = \sum_{l \neq k, m} b_l$  and  $t = (b_{k^*} - b_{m^*} - s)/2$ . If we fix  $m$  during iteration, then the shrinkage of  $b_m$  will be indirectly reflected in the other  $b_k$ . We propose that  $m$  change with  $k$  in the coordinate-wise minimization process to ensure that every coordinate can be equally shrunk. We summarize our proposed algorithm in Algorithm 1.

**Algorithm 1** Proximal gradient descent algorithm for SGL-MgDWD.

**Input:**  $\lambda_1, \lambda_2$ .

**Initialization:**  $\alpha^{(0)} = \mathbf{0}_K, \mathbf{B}^{(0)} = \mathbf{O}_{p \times K}, l = 0$ .

- 1: **repeat**
- 2:     Update  $\alpha$  according to (10) and (12):

$$\alpha^{(l+1)} = \mathbf{P}_K\{\alpha^{(l)} - L_\alpha^{-1}\mathcal{S}(\alpha^{(l)})\}.$$

- 3:     Update  $\tilde{\mathbf{B}}$  according to (11):

$$\tilde{\mathbf{B}} = \mathbf{B}^{(l)} - L_{\mathbf{B}}^{-1}\mathcal{S}(\mathbf{B}^{(l)}).$$

- 4:     Set  $\mathbf{B}^{(l+1)} \leftarrow \tilde{\mathbf{B}}$ .
- 5:     **repeat**
- 6:         **for**  $m = 1$  to  $K$  **do**
- 7:             **for**  $k$  in  $\{1, \dots, K\} \setminus m$  **do**
- 8:                 Update  $(t, s)$ :

$$t = \tilde{b}_k - \tilde{b}_m, \quad s = \sum_{r=1}^K \tilde{b}_r^{(l+1)} - \tilde{b}_k^{(l+1)} - \tilde{b}_m^{(l+1)}.$$

- 9:             Update  $b_k^{(l+1)}$  according to (17) and  $b_m^{(l+1)}$  by constraint:

$$b_k^{(l+1)} = \mathcal{T}(t, s, L_{\mathbf{B}}^{-1}\lambda_1), \quad b_m^{(l+1)} = -s - b_k^{(l+1)}.$$

- 10:         **end for**
- 11:         **end for**
- 12:         **until**  $\mathbf{B}^{(l+1)}$  convergence.
- 13:         Update  $\mathbf{B}^{(l+1)}$  according to (15):

$$\mathbf{B}^{(l+1)} = \mathcal{U}(\mathbf{B}^{(l+1)}, L_{\mathbf{B}}^{-1}\lambda_2).$$

- 14:     Set  $l \leftarrow l + 1$ .
  - 15: **until** some condition is met.
- Output:**  $\alpha^{(l)}$  and  $\mathbf{B}^{(l)}$ .

### 3. Numerical Analysis

In the following section, we use both simulated and real data sets to evaluate the finite sample properties of our proposed method. We compare the finite sample performance of SGL-MgDWD with  $L_1$ -regularized multinomial logistic regression ( $L_1$ -logistic).

#### 3.1. Simulation Studies

The data is generated from the following model. Consider the  $K$ -category classification problem where  $\pi_k = K^{-1}$  and  $g_k(\mathbf{X})$  is the density function of a normal distribution with mean vector  $\boldsymbol{\mu}_k = (\mu_{1k}, \mu_{2k}, \mathbf{0}_{p-2}^\top)^\top$  and covariance matrix  $\mathbf{I}_p$ , where  $(\mu_{1k}, \mu_{2k}) = (2 \cos(\pi r_k), 2 \sin(\pi r_k))$  with  $r_k = \frac{2(k-1)}{K}$ , for  $k = 1, \dots, K$ . In this model, only the first two variables contribute to the classification and their corresponding parameter vectors  $\boldsymbol{\beta}_1$  and  $\boldsymbol{\beta}_2$  form two groups of coefficients. The true model has the sparsity structure  $(s_e, s_g) = (2K, 2)$  for a total of  $K(p+1)$  coefficients. We set the sample size for each category to  $n_k = 50, 100, 200$  and  $400$ , and the number of classes to  $K = 5$  and  $11$ . We consider dimensionality  $p = 100$  and  $1000$ .

In what follows, we compare the proposed SGL-MgDWD method with the OVR method based on SGL-MgDWD with  $K = 2$  (OVR-SGL-gDWD). For SGL-MgDWD, logistic regression and OVR, the tuning parameter  $\lambda$  is optimized over a discrete set by minimizing the prediction error using 5-fold cross validation. In each simulation, we conduct 100 runs and use a testing set of equal size to evaluate each method's performance using the following criteria:

- Testing set accuracy, measuring the rate of correct classification;
- Signal, as the average number of correctly-selected element-wise and group-wise signals, that is, with  $\hat{b}_{jk} \neq 0$  and  $\hat{\boldsymbol{\beta}}_j \neq \mathbf{0}$ , respectively, denoted by the pair  $(s_e^+, s_g^+)$ ;
- Noise, as the average number of incorrectly-selected element-wise and group-wise components, that is, with  $\hat{b}_{jk} = 0$  and  $\hat{\boldsymbol{\beta}}_j = \mathbf{0}$ , respectively, denoted by the pair  $(n_e^+, n_g^+)$ .

Simulation results are summarized in Tables 1 and 2.

As shown in Tables 1 and 2, the proposed SGL-MgDWD method performs better than the  $L_1$ -logistic and OVR methods. Specifically, in each scenario, predictions from the SGL-MgDWD method had higher accuracy relative to the other two methods. Similarly, the SGL-MgDWD method correctly selected the signal components of the model with fewer incorrectly-selected noise components, again relative to the  $L_1$ -logistic and OVR methods. These simulation results also demonstrate that test accuracy increases with increasing sample size  $n_k$  and that test accuracy decreases with higher dimension  $p$  at fixed  $n_k$ . This is consistent with the derived theoretical properties. All computations were performed on a Tensorflow 2.3 CPU on Threadripper 2950X at 4.1 Ghz.

#### 3.2. HIV Data Analysis

Symptomatic distal sensory polyneuropathy (sDSP) is a common debilitating condition among people with HIV. This condition leads to neuropathic pain and is associated with substantial comorbidities and increased health care costs. Plasma miRNA profiles show differences between HIV patients with and without sDSP, and several miRNA biomarkers are reported to be associated with the presence of sDSP in HIV patients (see [30]). The corresponding binary classification problem was analyzed in [30] using random forest classifiers. However, the HIV data set can be further classified into four classes. The HIV data set has 1715 miRNA measures for 40 patients and is partitioned into four groups ( $K = 4$ ) with  $n_k = 10$  patients each category: non-HIV, HIV with no brain damage (HIVNBD), HIV with brain damage but stable (HIVBDS) and HIV with brain damage and unstable (HIVBDU). In the following analysis, we apply our proposed method to this classification problem. The primary aim was to identify critical miRNA biomarkers for each of the four groups. Beyond achieving a finer classification, this analysis is helpful in assessing related pathogenic effects for each patient group.

Given the small sample size of  $N = 40$ , we chose the tuning parameter  $\lambda$  by maximizing leave-one-out cross validation accuracy. We fixed  $(q, \tau) = (1, 0.1)$ . Table 3 shows the signal for

coefficient estimates obtained from the SGL-MgDWD method using the selected  $\lambda$ . We conclude that there are 22 critical miRNA biomarkers important to the classification problem. In particular, the biomarkers miR-25-star, miR-3171, miR-3924 and miR-4307 are not relevant to the non-HIV group; miR-4641, miR-4655-3p and miR-660 are not relevant to the HIVNBD group; miR-217 and miR-4683 are not relevant to the HIVBDS group; and miR-217 and miR-4307 are not relevant to the HIVBDU group.

**Table 1.** Simulation results for the SGL-MgDWD,  $L_1$ -logistic, and OVR methods with  $K = 5$ . Time is measured relative to a baseline logistic regression model with  $K = 5, p = 100$ , and  $N = 50$ . Numbers in parentheses denote standard deviations.

$n_k$	$p$	Method	Test Accuracy	Signal ( $s_e^+, s_g^+$ )	Noise ( $n_e^+, n_g^+$ )	Time (SD)
50	100	SGL-MgDWD	0.980	(9.99, 2)	(0, 0)	1.150 (0.173)
		$L_1$ -logistic	0.979	(9.00, 2)	(116.98, 26.17)	1.000 (0.153)
		OVR-SGL-gDWD	0.912	-	-	-
	1000	SGL-MgDWD	0.979	(10, 2)	(6.96, 1.94)	5.290 (0.166)
		$L_1$ -logistic	0.966	(10, 2)	(2793.65, 722.38)	5.130 (0.063)
		OVR-SGL-gDWD	0.740	-	-	-
100	100	SGL-MgDWD	0.981	(10, 2)	(0.07, 0.03)	1.453 (0.155)
		$L_1$ -logistic	0.980	(8.82, 2)	(35.18, 3.98)	1.258 (0.127)
		OVR-SGL-gDWD	0.828	-	-	-
	1000	SGL-MgDWD	0.980	(10, 2)	(1.01, 0.25)	4.863 (0.150)
		$L_1$ -logistic	0.978	(9.93, 2)	(1380.38, 192.37)	4.703 (0.061)
		OVR-SGL-gDWD	0.546	-	-	-
200	100	SGL-MgDWD	0.980	(10, 2)	(7.67, 2.08)	1.776 (0.164)
		$L_1$ -logistic	0.980	(9.39, 2)	(13.1, 0.72)	1.709 (0.175)
		OVR-SGL-gDWD	0.934	-	-	-
	1000	SGL-MgDWD	0.982	(10, 2)	(1.09, 0.29)	8.641 (0.186)
		$L_1$ -logistic	0.981	(9.79, 2)	(199.02, 2.51)	2.505 (0.121)
		OVR-SGL-gDWD	0.950	-	-	-
400	100	SGL-MgDWD	0.981	(10, 2)	(0.02, 0)	2.792 (0.159)
		$L_1$ -logistic	0.981	(10, 2)	(4.72, 3.95)	2.828 (0.115)
		OVR-SGL-gDWD	0.979	-	-	-
	1000	SGL-MgDWD	0.981	(10, 2)	(4.72, 3.95)	15.800 (0.221)
		$L_1$ -logistic	0.981	(9.6, 2)	(16.17, 0.02)	17.915 (1.585)
		OVR-SGL-gDWD	0.964	-	-	-

**Table 2.** Simulation results for the SGL-MgDWD,  $L_1$ -logistic, and OVR methods with  $K = 11$ . Time is measured relative to a baseline logistic regression model with  $K = 5$ ,  $p = 100$ , and  $N = 50$ . Numbers in parentheses denote standard deviations.

$n_k$	$p$	Method	Test Accuracy	Signal ( $s_e^+, s_g^+$ )	Noise ( $n_e^+, n_g^+$ )	Time (SD)
50	100	SGL-MgDWD	0.735	(21.41, 2)	(0.14, 0.02)	1.661 (0.143)
		$L_1$ -logistic	0.735	(20.13, 2)	(337.77, 22.07)	1.610 (0.110)
		OVR-SGL-gDWD	0.647	-	-	-
	1000	SGL-MgDWD	0.733	(21.25, 2)	(0, 0)	7.105 (0.205)
		$L_1$ -logistic	0.566	(20.67, 2)	(3805.97, 265.82)	6.933 (0.205)
		OVR-SGL-gDWD	0.382	-	-	-
100	100	SGL-MgDWD	0.737	(21.82, 2)	(0.06, 0.01)	2.518 (0.099)
		$L_1$ -logistic	0.721	(20, 2)	(173.17, 5.81)	2.418 (0.103)
		OVR-SGL-gDWD	0.609	-	-	-
	1000	SGL-MgDWD	0.737	(21.88, 2)	(5.4, 0.77)	12.371 (0.109)
		$L_1$ -logistic	0.697	(20.15, 2)	(1859.51, 9.04)	12.279 (0.114)
		OVR-SGL-gDWD	0.214	-	-	-
200	100	SGL-MgDWD	0.738	(22, 2)	(0, 0)	5.191 (0.079)
		$L_1$ -logistic	0.730	(20, 2)	(50.7, 0.08)	4.246 (0.100)
		OVR-SGL-gDWD	0.609	-	-	-
	1000	SGL-MgDWD	0.738	(21.98, 2)	(0.23, 0.04)	21.950 (0.241)
		$L_1$ -logistic	0.730	(20, 2)	(523.08, 1.07)	22.158 (0.163)
		OVR-SGL-gDWD	0.490	-	-	-
400	100	SGL-MgDWD	0.740	(22, 2)	(0, 0)	7.025 (0.172)
		$L_1$ -logistic	0.738	(20, 2)	(3.71, 3.48)	7.997 (0.122)
		OVR-SGL-gDWD	0.709	-	-	-
	1000	SGL-MgDWD	0.738	(22, 2)	(0.68, 0.11)	38.301 (0.200)
		$L_1$ -logistic	0.734	(20, 2)	(38.84, 35.37)	41.059 (2.064)
		OVR-SGL-gDWD	0.556	-	-	-

**Table 3.** Signal for the coefficient estimates obtained from the SGL-MgDWD method with  $(q, \tau) = (1, 0.1)$  for the HIV data set. The symbols “+” and “-” denote positive and negative coefficient estimates, respectively, while “0” denotes a zero coefficient (i.e., an irrelevant variable).

	Non-HIV	HIVNBD	HIVBDS	HIVBDU
interception	+	+	-	+
miR-255b	-	+	-	+
miR-217	+	-	0	0
miR-25-star	0	+	+	-
miR-3136-5p	-	-	+	-
miR-3152-3p	+	-	-	+
miR-3159	-	-	-	+
miR-3171	0	+	-	-

Table 3. Cont.

	Non-HIV	HIVNBD	HIVBDS	HIVBDU
miR-33b	-	-	-	+
miR-34c-3p	-	-	+	+
miR-3545-5p	-	+	-	+
miR-3654	-	-	-	+
miR-3924	0	-	+	-
miR-4307	0	-	+	0
miR-4474-5p	-	+	+	+
miR-4526	+	-	-	-
miR-4641	+	0	-	-
miR-4655-3p	+	0	-	-
miR-4680-5p	-	-	+	-
miR-4683	-	-	0	+
miR-589	-	+	+	-
miR-619	+	-	-	+
miR-660	+	0	-	+

**Author Contributions:** Conceptualization, L.K. and N.T.; Methodology, T.S., L.K. and N.T.; Formal Analysis, Y.L.; Data Curation, W.G.B., E.A. and C.P.; Writing—Review & Editing, Y.W., B.J. and L.K.; Supervision, B.J., L.K. and N.T. All authors have read and agreed to the published version of the manuscript.

**Funding:** A Canadian Institutes of Health Research Team Grant and Canadian HIV-Ageing Multidisciplinary Programmatic Strategy (CHAMPS) in NeuroHIV (Christopher Power) supported these studies. Bei Jiang and Linglong Kong were supported by the Natural Sciences and Engineering Research Council of Canada (NSERC). Christopher Power and Linglong Kong were supported by Canada Research Chairs in Neurological Infection and Immunity and Statistical Learning, respectively. Niansheng Tang was supported by grants from the National Natural Science Foundation of China (grant number: 11671349) and the Key Projects of the National Natural Science Foundation of China (grant number: 11731011).

**Acknowledgments:** The authors are thankful for the invitation of the two guest editors, Farouk Nathoo and Ejaz Ahmed. This work has also benefited from two anonymous reviewers' constructive comments and valuable feedback. The authors also thank the great help of Matthew Pietrosanu with editing.

**Conflicts of Interest:** The authors declare no conflict of interest.

## Appendix A. Proofs

### Appendix A.1. Proof of Lemma 1

**Proof.** For simplicity, we write  $p_j = P(y = j|X)$  and  $f_k = f_k(X)$ . Using the Lagrange multiplier method, we define

$$L(\mathbf{F}) = \mathbb{E} \left\{ \sum_{k=1}^K \mathbb{1}\{y = k\} \phi_q \{F(\mathbf{X})\} \middle| \mathbf{X} = \mathbf{u} \right\} + \mu \mathbf{1}_K^\top \mathbf{F}(\mathbf{X}) = \sum_{k=1}^K p_k \phi_q(f_k) + \mu f_k.$$

Then for each  $k$ ,

$$\frac{\partial L(\mathbf{F})}{\partial f_k} = \phi'_q(f_k) p_k + \mu = 0 \quad (\text{A1})$$

with

$$\phi'_q(f_j) = \begin{cases} -1, & f_k \leq Q \\ -(Q f_k^{-1})^q, & f_k > Q. \end{cases}$$

Without loss of generality, assume that  $p_1 > p_2 \geq p_3 \geq \dots \geq p_{K-1} > p_K$ . Note that  $-1 \leq \phi'_q < 0$ , and so  $p_j \geq -\phi'_q(f_k) p_k = \mu > 0$  and  $\mu = p_k$  if and only if  $f_k \leq Q$ .

If  $\mu < p_K < p_k$ , then  $p_K \neq \mu$  when  $f_K > Q$ , which implies that  $f_k > f_K > Q$  for all  $1 \leq k \leq K$ . Hence, substituting  $\phi'_q(f_k) = -(Qf_k^{-1})^q$  into (A1) yields

$$f_k = Q\sqrt[q]{p_k\mu^{-1}} > Q > 0.$$

However,  $\sum_{k=1}^K f_k > 0$ , contradicting the sum-to-zero constraint. Therefore,  $\mu = p_K < p_k$  for  $k < K$  and the result follows.  $\square$

Appendix A.2. Proof of Theorem 1

**Lemma A1.** Under (C1),  $\mathcal{L}(\boldsymbol{\theta})$  exists, and it is convex on  $\boldsymbol{\theta}$ .

**Proof.** The existence of  $\mathcal{L}(\boldsymbol{\theta})$  will be satisfied if

$$\mathbb{E}_{\mathbf{X}|y} \{ |\phi_q(\mathbf{Z}^\top \boldsymbol{\theta}_k)| \mid y = k \} = \int_{\mathcal{X}} |\phi_q(\mathbf{Z}^\top \boldsymbol{\theta}_k)| g_k(\mathbf{X}) d\mathbf{X} < \infty.$$

We divide  $\mathcal{X}$  into two disjoint subsets. Defining  $\mathcal{X}_k = \{ \mathbf{X} \in \mathcal{X} \mid \mathbf{Z}^\top \boldsymbol{\theta}_k > Q \}$ , it is clear that

$$\int_{\mathcal{X}_k} |\phi_q(\mathbf{Z}^\top \boldsymbol{\theta}_k)| g_k(\mathbf{X}) d\mathbf{X} \leq (q + 1)^{-1} \int_{\mathcal{X}_k} g_k(\mathbf{X}) d\mathbf{X} < \infty.$$

Note that  $0 < \phi_q(u) < (1 + q)^{-1} < 1$  when  $u > Q$ . On the other hand, for  $\mathcal{X}_k^c = \{ \mathbf{X} \in \mathcal{X} \mid \mathbf{Z}^\top \boldsymbol{\theta}_k \leq Q \}$ ,

$$\int_{\mathcal{X}_k^c} |\phi_q(\mathbf{Z}^\top \boldsymbol{\theta}_k)| g_k(\mathbf{X}) d\mathbf{X} \leq |1 - a_k| + \sum_{j=1}^p b_{jk} \int_{\mathcal{X}} |x_j| g_k(\mathbf{X}) d\mathbf{X} < \infty,$$

if  $\mathbb{E}_{\mathbf{X}|y} \{ |x_j| \mid y = k \} < \infty$  for all  $k \in \mathcal{Y}$ . This completes the proof of the existence of  $\mathcal{L}(\boldsymbol{\theta})$ .

Recall that

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \int_{\mathcal{X}} \phi_q(\mathbf{Z}^\top \boldsymbol{\theta}_k) g_k(\mathbf{X}) d\mathbf{X},$$

where  $\phi_q(u)$  is a convex function of  $u$ , so its composition with the affine mapping  $u = \mathbf{Z}^\top \boldsymbol{\theta}_k$  is still convex in  $\boldsymbol{\theta}_k$ . Clearly,  $g_k(\mathbf{X})$ ,  $\pi_k > 0$ , so the non-negatively-weighted integral and sum both preserve convexity.  $\square$

**Lemma A2.** Existence of minimizers of  $\mathcal{L}(\boldsymbol{\theta})$  on  $\mathcal{C} = \{ \boldsymbol{\theta} \in \mathbb{R}^{K(p+1)} \mid \mathbf{C}\boldsymbol{\theta} = \mathbf{0}_K \}$ , where  $\mathbf{C} = \mathbf{1}_K^\top \otimes \mathbf{I}_{p+1}$ .

**Proof.** By Jensen’s inequality, for any  $\boldsymbol{\theta} \in \mathcal{C}$ , we have that

$$\mathcal{L}(\boldsymbol{\theta}) \geq \phi_q \left( \sum_{k=1}^K \pi_k \mathbb{E} \{ \mathbf{Z}^\top \boldsymbol{\theta}_k \mid y = k \} \right).$$

Let  $\boldsymbol{\mu} = \text{vec} \{ (\pi_k \mathbb{E} \{ z_j \mid y = k \})_{jk} \}$ , where  $\|\boldsymbol{\mu}\|_2 \geq (\sum_{k=1}^K \pi_k^2)^{\frac{1}{2}} \geq K^{-\frac{1}{2}} > 0$ . For some  $C > 0$ , we have that



$$\begin{aligned}
 \mathcal{L}(\boldsymbol{\theta}) &\geq \phi_q(\boldsymbol{\mu}^\top \boldsymbol{\theta}) = \mathbb{1}\{\boldsymbol{\mu}^\top \boldsymbol{\theta} < Q\}(1 - \boldsymbol{\mu}^\top \boldsymbol{\theta}) + \mathbb{1}\{\boldsymbol{\mu}^\top \boldsymbol{\theta} \geq Q\}\phi_q(\boldsymbol{\mu}^\top \boldsymbol{\theta}) \\
 &\geq \mathbb{1}\{\boldsymbol{\mu}^\top \boldsymbol{\theta} < Q\}|1 - |\boldsymbol{\mu}^\top \boldsymbol{\theta}|| \\
 &= \mathbb{1}\{\boldsymbol{\mu}^\top \boldsymbol{\theta} < -(C + 1)\}(|\boldsymbol{\mu}^\top \boldsymbol{\theta}| - 1) + \mathbb{1}\{-(C + 1) < \boldsymbol{\mu}^\top \boldsymbol{\theta} < -1\}(|\boldsymbol{\mu}^\top \boldsymbol{\theta}| - 1) \\
 &\quad + \mathbb{1}\{-1 < \boldsymbol{\mu}^\top \boldsymbol{\theta} < Q\}(1 - |\boldsymbol{\mu}^\top \boldsymbol{\theta}|) \\
 &> \mathbb{1}\{\|\boldsymbol{\mu}\|_2 \|\boldsymbol{\theta}\|_2 > C + 1\}C \\
 &= \mathbb{1}\left\{\|\boldsymbol{\theta}\|_2 > \frac{C + 1}{\|\boldsymbol{\mu}\|_2}\right\}C.
 \end{aligned}$$

Note that  $1 - \boldsymbol{\mu}^\top \boldsymbol{\theta} > 1 - Q > 0$  when  $\boldsymbol{\mu}^\top \boldsymbol{\theta} < Q$ . By the Cauchy–Schwarz inequality,  $-\boldsymbol{\mu}^\top \boldsymbol{\theta} = |\boldsymbol{\mu}^\top \boldsymbol{\theta}| \leq \|\boldsymbol{\mu}\|_2 \|\boldsymbol{\theta}\|_2$ .

Hence, if  $\|\boldsymbol{\theta}\|_2 > \frac{C + 1}{\|\boldsymbol{\mu}\|_2} > 0$ , then  $\mathcal{L}(\boldsymbol{\theta}) > C > 0$ . The contrapositive of this result implies the existence of a minimizer in the unconstrained problem. That is, the closed set  $\{\boldsymbol{\theta} \in \mathcal{C} \mid \mathcal{L}(\boldsymbol{\theta}) \leq C\}$  is bounded for some large enough  $C$ . This guarantees the existence of a solution, as desired.  $\square$

**Lemma A3.** Under (C1),  $\mathcal{S}(\boldsymbol{\theta})$  exists and

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathcal{S}(\boldsymbol{\theta}).$$

**Proof.** The existence of  $\mathcal{S}(\boldsymbol{\theta})$  will follow if

$$\int_{\mathcal{X}} |\phi'_q(\mathbf{Z}^\top \theta_k) z_j| \pi_k g_k(\mathbf{X}) d\mathbf{X} \leq \pi_k \int_{\mathcal{X}} |z_j| g_k(\mathbf{X}) d\mathbf{X} < \infty$$

for  $j = 1, \dots, p + 1$ . Note that  $|\phi'_q(u)| \leq 1$  when  $u > Q$ .

For every  $\theta_{kj} \in \mathbb{R}$ ,  $\phi_q(\mathbf{Z}^\top \theta_k)$  is a Lebesgue integrable function of  $\mathbf{X}$ . For any  $u \in \mathbb{R}$ ,  $\phi'_q(u)$  exists and  $|\phi'_q(u)| \leq 1$ . Hence, by the Leibniz integral rule, we have that

$$\begin{aligned}
 \frac{\partial}{\partial \theta_{jk}} \int_{\mathcal{X}} \phi_q(\mathbf{Z}^\top \theta_k) \pi_k g_k(\mathbf{X}) d\mathbf{X} &= \int_{\mathcal{X}} \frac{\partial \phi_q(\mathbf{Z}^\top \theta_k)}{\partial \theta_{jk}} \pi_k g_k(\mathbf{X}) d\mathbf{X} \\
 &= \int_{\mathcal{X}} \phi'_q(\mathbf{Z}^\top \theta_k) z_j \pi_k g_k(\mathbf{X}) d\mathbf{X}
 \end{aligned}$$

and for any  $l \neq k$ ,

$$\frac{\partial}{\partial \theta_{jl}} \int_{\mathcal{X}} \phi_q(\mathbf{Z}^\top \theta_k) \pi_k g_k(\mathbf{X}) d\mathbf{X} = 0,$$

which is sufficient to show that

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathcal{S}(\boldsymbol{\theta}).$$

$\square$

**Lemma A4.** Suppose (C1) is satisfied. Then (C2) implies that  $\mathbf{b}_k^* \neq \mathbf{0}$ .

**Proof.** We can rewrite  $\phi_q(u)$  as

$$\begin{aligned} \phi_q(u) &= \mathbb{1}\{u \leq Q\}(1-u) + \mathbb{1}\{u > Q\}(1-Q)\left(\frac{Q}{u}\right)^q \\ &= \left\{ -\mathbb{1}\{u \leq Q\} - \mathbb{1}\{u > Q\}\left(\frac{Q}{u}\right)^{q+1} \right\}u + \mathbb{1}\{u \leq Q\} + \mathbb{1}\{u > Q\}\left(\frac{Q}{u}\right)^q \\ &= \phi'_q(u)u + \mathbb{1}\{u \leq Q\} + \mathbb{1}\{u > Q\}\left(\frac{Q}{u}\right)^q. \end{aligned}$$

Then for any  $\gamma \in \mathbb{R}^{p+1}$  and its corresponding  $\mathcal{X}_k = \{\mathbf{X} \in \mathcal{X} | \mathbf{Z}^\top \gamma > Q\}$ , we have that

$$\begin{aligned} &\mathbb{E}\{\mathbb{1}\{y = k\}\phi_q(\mathbf{Z}^\top \gamma)\} \\ &= \mathbb{E}\{\mathbb{1}\{y = k\}\phi'_q(\mathbf{Z}^\top \gamma)\mathbf{Z}^\top \gamma\} + \mathbb{E}\{\mathbb{1}\{y = k, \mathbf{Z}^\top \gamma \leq Q\}\} \\ &\quad + \mathbb{E}\left\{\mathbb{1}\{y = k, \mathbf{Z}^\top \gamma > Q\}\left(\frac{Q}{\mathbf{Z}^\top \gamma}\right)^q\right\} \\ &= \mathbf{S}_k^\top(\gamma)\gamma + \Pr\{y = k, \mathbf{X} \notin \mathcal{X}_k\} + \mathbb{E}\left\{\mathbb{1}\{y = k, \mathbf{X} \in \mathcal{X}_k\}\left(\frac{Q}{\mathbf{Z}^\top \gamma}\right)^q\right\} \\ &= \mathbf{S}_k^\top(\gamma)\gamma + \pi_k\left(1 - \mathbb{E}\left\{\mathbb{1}\{\mathbf{X} \in \mathcal{X}_k\}\left\{1 - \left(\frac{Q}{\mathbf{Z}^\top \gamma}\right)^q\right\} \mid y = k\right\}\right). \end{aligned}$$

Let  $\boldsymbol{\theta}^* \in \mathcal{C}$  be a local minimizer. It follows that  $\mathbf{PS}(\boldsymbol{\theta}^*) = \mathbf{0}$  and  $\sum_{k=1}^K \mathbf{S}_k^\top(\boldsymbol{\theta}_k^*)\boldsymbol{\theta}_k^* = \mathbf{S}^\top(\boldsymbol{\theta}^*)\boldsymbol{\theta}^* = \mathbf{0}$  since  $\boldsymbol{\theta}^* = \mathbf{P}\boldsymbol{\theta}^*$  and  $\mathbf{P} = (\mathbf{I}_K - K^{-1}\mathbf{1}_K\mathbf{1}_K^\top) \otimes \mathbf{I}_{p+1}$ . Therefore,

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}^*) &= \mathbb{E}\{\mathbb{1}\{y = k\}\phi_q(\mathbf{Z}^\top \boldsymbol{\theta}_k^*)\} \\ &= \sum_{k=1}^K \pi_k\left(1 - \mathbb{E}\left\{\mathbb{1}\{\mathbf{X} \in \mathcal{X}_k^*\}\left\{1 - \left(\frac{Q}{\mathbf{Z}^\top \boldsymbol{\theta}_k^*}\right)^q\right\} \mid y = k\right\}\right) \\ &= \sum_{k=1}^K \pi_k\left(1 - \Pr\{\mathbf{X} \in \mathcal{X}_k^* | y = k\}\mathbb{E}\left\{1 - \left(\frac{Q}{\mathbf{Z}^\top \boldsymbol{\theta}_k^*}\right)^q \mid y = k, \mathbf{X} \in \mathcal{X}_k^*\right\}\right). \end{aligned} \tag{A2}$$

For any  $\gamma \in \mathbb{R}^{p+1}$  and its corresponding  $\mathcal{X}_k = \{\mathbf{X} \in \mathcal{X} | \mathbf{Z}^\top \gamma > Q\}$ , we always have that

$$0 < \mathbb{E}\left\{\left(\frac{Q}{\mathbf{Z}^\top \gamma}\right)^q \mid y = k, \mathbf{X} \in \mathcal{X}_k\right\} < 1.$$

If  $\gamma = \mathbf{0}_{p+1}$ , then  $\mathcal{X}_k = \emptyset$  so that  $\Pr\{y = k, \mathbf{X} \notin \mathcal{X}_k\} = \pi_k$  and  $\Pr\{y = k, \mathbf{X} \in \mathcal{X}_k\} = 0$ . If  $\gamma_1 \leq Q$  and  $\gamma_{/1} = \mathbf{0}_p$ , then  $\mathcal{X}_k = \emptyset$ , giving the same conclusions as the previous case. If  $\gamma_1 > Q$  and  $\gamma_{/1} = \mathbf{0}_p$ , then  $\mathcal{X}_k = \mathcal{X}$  so that  $\Pr\{y = k, \mathbf{X} \notin \mathcal{X}_k\} = 0$  and  $\Pr\{y = k, \mathbf{X} \in \mathcal{X}_k\} = \pi_k$ . Consequently, when  $0 < \Pr\{\mathbf{X} \in \mathcal{X}_k | y = k\} < 1$ , then neither  $\mathcal{X}_k$  nor  $\mathcal{X}$  equal  $\emptyset$ , so  $\mathbf{b}_k \neq 0$  follows.

Note that  $\Pr\{\mathbf{X} \notin \mathcal{X}_k | y = k\} > 0$  implies that  $\Pr\{0 < \mathbf{Z}^\top \gamma \leq Q | y = k\} > 0$  or  $\Pr\{\mathbf{Z}^\top \gamma \leq 0 | y = k\} > 0$ , and so special attention should be paid to bounded random variables.  $\square$

**Lemma A5.** Under (C1),  $\mathcal{H}(\boldsymbol{\theta})$  exists and

$$\frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} = \mathcal{H}(\boldsymbol{\theta}).$$

Furthermore,  $\mathcal{H}(\boldsymbol{\theta}^*) \succ \mathbf{O}_{K(p+1)}$  when (C2) and (C3) hold.

**Proof.** The existence of  $\mathcal{H}(\boldsymbol{\theta})$  follows if its all entries are absolutely integrable, that is, for any  $j, k = 1, \dots, p + 1$ ,

$$\begin{aligned} & \int_{\mathcal{X}} |\mathbb{1}\{\mathbf{Z}^\top \boldsymbol{\theta}_k > Q\} \phi_q''(\mathbf{Z}^\top \boldsymbol{\theta}_k) z_j z_l | \pi_k g_k(\mathbf{X}) d\mathbf{X} \\ & \leq (q + q^{-1} + 2) \int_{\mathcal{X}_k^c} |z_j z_l| g_k(\mathbf{X}) d\mathbf{X} \\ & < \infty. \end{aligned}$$

Equivalently, the result follows if  $\mathbb{E}_{\mathbf{X}|y}\{|z_j z_l| \mid y = k\} < \infty$  for all  $k \in \mathcal{Y}$ . Note that  $0 < \phi_q''(u) \leq q + q^{-1} + 2$  when  $u > Q$ .

Let  $\eta$  be a test function belonging to the Schwartz space  $\mathcal{D}$ . Then  $\eta' \in \mathcal{D}$  with some support denoted by  $\text{supp}(\eta')$ .

Clearly,  $\phi_q'(u)$  is not differentiable at  $Q$  but is Lipschitz continuous. Therefore, the measurable function  $S_k(\boldsymbol{\theta}_k)$  is a locally integrable function of  $\boldsymbol{\theta}_k$ . Then the (regular) generalized functions  $S_k(\boldsymbol{\theta}_k)$  belong to the dual space of  $\mathcal{D}$ .

For the distributional derivative of  $S_k(\boldsymbol{\theta}_k)$  with respect to  $\theta_{jk}$ , we have that

$$\begin{aligned} \left| \left\langle \frac{\partial S_k(\boldsymbol{\theta}_k)}{\partial \theta_{jk}}, \eta(\theta_{jk}) \right\rangle \right| &= \left| - \left\langle S_k(\boldsymbol{\theta}_k), \frac{d\eta(\theta_{jk})}{d\theta_{jk}} \right\rangle \right| \\ &\leq \int_{\mathbb{R}} |S_k(\boldsymbol{\theta}_k) \eta'(\theta_{jk})| d\theta_{jk} \\ &\leq \max_{\theta_{jk} \in \text{supp}(\eta')} |\eta'(\theta_{jk})| \int_{\text{supp}(\eta')} |S_k(\boldsymbol{\theta}_k)| d\theta_{jk} \\ &< \infty \end{aligned}$$

implying that the function  $f(\theta_{jk}, \mathbf{X}) = \phi_q'(\mathbf{Z}^\top \boldsymbol{\theta}_k) \mathbf{Z} \pi_k g_k(\mathbf{X}) \eta'(\theta_{jk})$  is integrable on  $\mathbb{R} \times \mathcal{X}$ . Therefore, by Fubini's Theorem,

$$\begin{aligned} \left\langle \frac{\partial S_k(\boldsymbol{\theta}_k)}{\partial \theta_{jk}}, \eta(\theta_{jk}) \right\rangle &= - \left\langle S_k(\boldsymbol{\theta}_k), \frac{d\eta(\theta_{jk})}{d\theta_{jk}} \right\rangle \\ &= \int_{\mathcal{X}} - \left\langle \phi_q'(\mathbf{Z}^\top \boldsymbol{\theta}_k) \mathbf{Z} \pi_k g_k(\mathbf{X}), \frac{d\eta(\theta_{jk})}{d\theta_{jk}} \right\rangle d\mathbf{X} \\ &= \int_{\mathcal{X}} \left\langle \frac{\partial \phi_q'(\mathbf{Z}^\top \boldsymbol{\theta}_k)}{\partial \theta_{jk}} \mathbf{Z} \pi_k g_k(\mathbf{X}), \eta(\theta_{jk}) \right\rangle d\mathbf{X} \\ &= \left\langle \mathbb{E} \left\{ \frac{\partial \phi_q'(\mathbf{Z}^\top \boldsymbol{\theta}_k)}{\partial \theta_{jk}} \mathbf{Z} \mathbb{1}\{y = k\} \right\}, \eta(\theta_{jk}) \right\rangle, \end{aligned}$$

which implies that

$$\frac{\partial S_k(\boldsymbol{\theta}_k)}{\partial \theta_{jk}} = \mathbb{E} \left\{ \frac{\partial \phi_q'(\mathbf{Z}^\top \boldsymbol{\theta}_k)}{\partial \theta_{jk}} \mathbf{Z} \mathbb{1}\{y = k\} \right\}.$$

Recall that  $\phi_q'$  can be written as

$$\phi_q'(u) = \phi_q'(u) \mathbb{1}\{u > Q\} + (-1) \mathbb{1}\{u \leq Q\} = (\phi_q'(u) + 1) \mathbb{1}\{u > Q\} - 1,$$

which contains a Schwartz product between the differentiable function  $\phi'_q(u)$  and the generalized function  $\mathbb{1}\{u > Q\}$ . Note that

$$\begin{aligned} \mathbb{1}\{\mathbf{Z}^\top \boldsymbol{\theta}_k > Q\} &= \mathbb{1}\{z_j > 0, \theta_{jk} > c_{jk}\} + \mathbb{1}\{z_j \leq 0, \theta_{jk} \leq c_{jk}\} \\ &= (2\mathbb{1}\{z_j > 0\} - 1)\mathbb{1}\{\theta_{jk} > c_{jk}\} + (1 - \mathbb{1}\{z_j > 0\}) \\ &= \text{sign}(z_j)\mathbb{1}\{\theta_{jk} > c_{jk}\} + \mathbb{1}\{z_j \leq 0\}, \end{aligned}$$

where  $c_{jk} = (Q - \sum_{l \neq j} z_l \theta_{lk}) / z_j$  and

$$\begin{aligned} \frac{\partial \mathbb{1}\{\mathbf{Z}^\top \boldsymbol{\theta}_k > Q\}}{\partial \theta_{jk}} + 0 &= \text{sign}(z_j)\delta(\theta_{jk} - c_{jk}) \\ &= \text{sign}(z_j)|z_j|\delta(\mathbf{Z}^\top \boldsymbol{\theta}_k - Q) \\ &= z_j\delta(\mathbf{Z}^\top \boldsymbol{\theta}_k - Q), \end{aligned}$$

where  $\delta(x)$  is the Dirac delta function and the distributional derivative of  $\mathbb{1}\{x > 0\}$ . Recall that  $\delta(cx) = \delta(x)/|c|$  and  $f(x)\delta(x - c) = f(c)\delta(x - c)$  for some constant  $c$  and function  $f$ .

Thus, by the product rule for the distributional derivative of the Schwartz product,

$$\begin{aligned} \frac{\partial \phi'_q(\mathbf{Z}^\top \boldsymbol{\theta}_k)}{\partial \theta_{jk}} &= \frac{\partial(\phi'_q(\mathbf{Z}^\top \boldsymbol{\theta}_k) + 1)}{\partial \theta_{jk}} \mathbb{1}\{\mathbf{Z}^\top \boldsymbol{\theta}_k > Q\} + (\phi'_q(\mathbf{Z}^\top \boldsymbol{\theta}_k) + 1) \frac{\partial \mathbb{1}\{\mathbf{Z}^\top \boldsymbol{\theta}_k > Q\}}{\partial \theta_{jk}} \\ &= \phi''_q(\mathbf{Z}^\top \boldsymbol{\theta}_k) z_j \mathbb{1}\{\mathbf{Z}^\top \boldsymbol{\theta}_k > Q\} + (\phi'_q(\mathbf{Z}^\top \boldsymbol{\theta}_k) + 1) z_j \delta(\mathbf{Z}^\top \boldsymbol{\theta}_k - Q) \\ &= \phi''_q(\mathbf{Z}^\top \boldsymbol{\theta}_k) z_j \mathbb{1}\{\mathbf{Z}^\top \boldsymbol{\theta}_k > Q\}. \end{aligned}$$

Substituting the above expression, we obtain

$$\frac{\partial \mathcal{S}_k(\boldsymbol{\theta}_k)}{\partial \theta_{jk}} = \mathbb{E} \left\{ \phi''_q(\mathbf{Z}^\top \boldsymbol{\theta}_k) \mathbf{Z} z_j \mathbb{1}\{\mathbf{Z}^\top \boldsymbol{\theta}_k > Q\} \mathbb{1}\{y = k\} \right\}.$$

Similarly, for  $l \neq k$ , we have the distributional derivative

$$\frac{\partial \mathcal{S}_k(\boldsymbol{\theta}_k)}{\partial \theta_{jl}} = 0.$$

Recall that the distributional derivative does not depend on the order of differentiation and agrees with the classical derivative whenever the latter exists. To summarize, we have that

$$\mathbf{H}_k(\boldsymbol{\theta}_k) = \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_k \partial \boldsymbol{\theta}_k^\top} = \frac{\partial \mathcal{S}_k(\boldsymbol{\theta}_k)}{\partial \boldsymbol{\theta}_k^\top}, \quad \mathcal{H}(\boldsymbol{\theta}) = \bigoplus_{k=1}^K \mathbf{H}_k(\boldsymbol{\theta}_k).$$

The  $\mathbf{H}_k(\boldsymbol{\theta}_k)$  are symmetric matrices, so  $\mathcal{H}(\boldsymbol{\theta})$  is also symmetric.

In the sense of generalized functions, differentiation is a continuous operation with respect to convergence in  $\mathcal{D}'$ . Therefore,  $\phi'_0 = \lim_{q \rightarrow 0} \phi'_q = -\mathbb{1}\{u \leq 0\}$  and  $\phi''_0 = \lim_{q \rightarrow 0} \phi''_q = \delta(u)$ ;  $\phi'_\infty = \lim_{q \rightarrow \infty} \phi'_q = -\mathbb{1}\{u \leq 1\}$  and  $\phi''_\infty = \lim_{q \rightarrow \infty} \phi''_q = \delta(u - 1)$ , which coincides with results from the hinge loss.

Next,  $\mathcal{H}(\boldsymbol{\theta}) \succ \mathbf{O}_{K(p+1)}$  if and only if both  $\mathbf{H}_1(\boldsymbol{\theta}_1)$  and its Schur complement  $\bigoplus_{k=2}^K \mathbf{H}_k(\boldsymbol{\theta}_k)$  are both symmetric and positive definite. We can deduce that  $\mathcal{H}(\boldsymbol{\theta}) \succ \mathbf{O}_{K(p+1)}$  if and only if  $\mathbf{H}_k(\boldsymbol{\theta}_k) \succ \mathbf{O}_{p+1}$  for all  $k$ .

Note that there exists  $c > 0$  such that  $\varphi''_q(\mathbf{Z}^\top \boldsymbol{\theta}_k) \geq c$  on  $\mathcal{X}_k$ . Then for any  $\boldsymbol{\gamma} \in \mathbb{R}^{p+1}$ ,

$$\begin{aligned} \boldsymbol{\gamma}^\top \mathbf{H}_k(\boldsymbol{\theta}_k) \boldsymbol{\gamma} &= \pi_k \int_{\mathcal{X}_k} \varphi''_q(\mathbf{Z}^\top \boldsymbol{\theta}_k) (\mathbf{Z}^\top \boldsymbol{\gamma})^2 g_k(\mathbf{X}) d\mathbf{X} \\ &\geq c \Pr\{\mathbf{X} \in \mathcal{X}_k, y = k\} \mathbb{E}\{(\mathbf{Z}^\top \boldsymbol{\gamma})^2 | \mathbf{X} \in \mathcal{X}_k, y = k\} \\ &\geq c \Pr\{\mathbf{X} \in \mathcal{X}_k, y = k\} (\gamma_0^2 + \boldsymbol{\gamma}_1^\top \text{Var}\{\mathbf{X} | \mathbf{X} \in \mathcal{X}_k, y = k\} \boldsymbol{\gamma}_1), \end{aligned}$$

which implies that  $\boldsymbol{\gamma}^\top \mathbf{H}_k(\boldsymbol{\theta}_k) \boldsymbol{\gamma} = 0$  if and only if  $\boldsymbol{\gamma} = \mathbf{0}_{p+1}$  when  $\text{Var}\{\mathbf{X} | \mathbf{X} \in \mathcal{X}_k, y = k\}$  is assumed to be non-singular. Assuming that  $\text{Var}\{\mathbf{X} | y = k\} \succ \mathbf{0}$  implies that  $\text{Var}\{\mathbf{X} | \mathbf{X} \in \mathcal{X}_k, y = k\} \succeq \mathbf{0}$ .  $\square$

**Proof of Theorem 1.** By Lemma A2, a minimizer  $\boldsymbol{\theta}^* \in \mathcal{C}$  exists with  $\mathbf{b}_k^* \neq \mathbf{0}_p$  (by Lemma A4) and  $\mathcal{H}(\boldsymbol{\theta}^*) \succ \mathbf{0}_{K(p+1)}$  (by Lemma A5). By the second-order Lagrange condition and the convexity of  $\mathcal{L}(\boldsymbol{\theta})$  (by Lemma A1), a minimizer of the population MgDWD loss is unique.

Recall from (A2) that

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}^*) &= \mathbb{E}\{\mathbb{1}\{y = k\} \phi_q(\mathbf{Z}^\top \boldsymbol{\theta}_k^*)\} \\ &= \sum_{k=1}^K \pi_k \left(1 - \mathbb{E}\left\{\mathbb{1}\{\mathbf{X} \in \mathcal{X}_k^*\} \left\{1 - \left(\frac{Q}{\mathbf{Z}^\top \boldsymbol{\theta}_k^*}\right)^q\right\} \middle| y = k\right\}\right) \\ &= \sum_{k=1}^K A(k, q) \pi_k. \end{aligned}$$

It follows that

$$\begin{aligned} 0 &\leq \mathbb{E}\left\{\mathbb{1}\{\mathbf{X} \in \mathcal{X}_k^*\} \left\{1 - \left(\frac{Q}{\mathbf{Z}^\top \boldsymbol{\theta}_k^*}\right)^q\right\} \middle| y = k\right\} \\ &< \mathbb{E}\left\{\mathbb{1}\{\mathbf{Z}^\top \boldsymbol{\gamma} > 1 + q^{-1}\} + \mathbb{1}\{Q < \mathbf{Z}^\top \boldsymbol{\gamma} \leq 1 + q^{-1}\} \left\{1 - \left(\frac{Q}{1 + q^{-1}}\right)^q\right\} \middle| y = m\right\} \\ &= \Pr\{\mathbf{Z}^\top \boldsymbol{\gamma} > Q | y = m\} - \Pr\{Q < \mathbf{Z}^\top \boldsymbol{\gamma} \leq Q^{-1} | y = m\} Q^{2q} \\ &\leq 1 \end{aligned}$$

and

$$\begin{aligned} 1 &\geq \mathbb{E}\left\{\mathbb{1}\{\mathbf{X} \in \mathcal{X}_k^*\} \left\{1 - \left(\frac{Q}{\mathbf{Z}^\top \boldsymbol{\theta}_k^*}\right)^q\right\} \middle| y = k\right\} \\ &> \mathbb{E}\left\{\mathbb{1}\{\mathbf{Z}^\top \boldsymbol{\theta}_k^* > 1 + \epsilon\} \left\{1 - \left(\frac{Q}{1 + \epsilon}\right)^q\right\} \middle| y = k\right\} \\ &\geq \sup_{\epsilon > 0} \left\{1 - \left(\frac{Q}{1 + \epsilon}\right)^q\right\} \Pr\{\mathbf{Z}^\top \boldsymbol{\theta}_k^* > 1 + \epsilon \mid y = m\} \\ &\geq 0. \end{aligned}$$

Consequently,  $0 \leq u(k, q) \leq A(k, q) \leq v(k, q) \leq 1$ .

Note that  $\lim_{q \rightarrow \infty} (1 + \epsilon)^{-q} Q^q = e^{-1}$  when  $\epsilon = 0$  and  $\lim_{q \rightarrow \infty} (1 + \epsilon)^{-q} Q^q = 0$  when  $\epsilon > 0$ .

The difference between these two results is attributed to pointwise convergence.

Let  $f_m = 1 - A(k, m) \in \mathcal{D}'$  with  $m = 1, 2, \dots$  and  $\eta \in \mathcal{D}$ . By Fubini's theorem and the dominated convergence theorem,

$$\begin{aligned} \lim_{m \rightarrow \infty} \langle f_m, \eta \rangle &= \lim_{m \rightarrow \infty} \left\langle \mathbb{E} \left\{ \mathbb{1} \{ \mathbf{X} \in \mathcal{X}_k^* \} \left( \frac{Q}{\mathbf{Z}^\top \boldsymbol{\theta}_k^*} \right)^q \middle| y = k \right\}, \eta(\gamma) \right\rangle \\ &= \lim_{m \rightarrow \infty} \mathbb{E} \left\{ \left\langle \mathbb{1} \{ \mathbf{X} \in \mathcal{X}_k^* \} \left( \frac{Q}{\mathbf{Z}^\top \boldsymbol{\theta}_k^*} \right)^q, \eta(\boldsymbol{\theta}_k^*) \right\rangle \middle| y = k \right\} \\ &= \mathbb{E} \left\{ \lim_{m \rightarrow \infty} \left\langle \mathbb{1} \{ \mathbf{X} \in \mathcal{X}_k^* \} \left( \frac{Q}{\mathbf{Z}^\top \boldsymbol{\theta}_k^*} \right)^q, \eta(\boldsymbol{\theta}_k^*) \right\rangle \middle| y = k \right\} \\ &= 0 = \langle 0, \eta(\boldsymbol{\theta}_k^*) \rangle. \end{aligned}$$

Similarly,

$$\begin{aligned} \lim_{m \rightarrow 0} \langle f_m, \eta \rangle &= \mathbb{E} \left\{ \lim_{m \rightarrow 0} \left\langle \mathbb{1} \{ \mathbf{X} \in \mathcal{X}_k^* \} \left( \frac{Q}{\mathbf{Z}^\top \boldsymbol{\theta}_k^*} \right)^q, \eta(\gamma) \right\rangle \middle| y = k \right\} \\ &= \mathbb{E} \left\{ \left\langle \mathbb{1} \{ \mathbf{Z}^\top \boldsymbol{\theta}_k^* > 0 \}, \eta(\gamma) \right\rangle \middle| y = k \right\} \\ &= \left\langle \mathbb{E} \{ \mathbb{1} \{ \mathbf{Z}^\top \boldsymbol{\theta}_k^* > 0 \} \middle| y = k \}, \eta(\boldsymbol{\theta}_k^*) \right\rangle \\ &= \left\langle \Pr \{ \mathbf{Z}^\top \boldsymbol{\theta}_k^* > 0 \mid y = k \}, \eta(\boldsymbol{\theta}_k^*) \right\rangle, \end{aligned}$$

hence

$$A(k, \infty) = \lim_{q \rightarrow \infty} A(k, q) = \Pr \{ \mathbf{X} \notin \mathcal{X}_k^* \mid y = k \}, \text{ and } A(k, 0) = \lim_{q \rightarrow 0} A(k, q) = 1.$$

As a result,  $A(k, \infty)$  coincides with the population hinge/SVM loss and  $A(k, 0)$  is independent of  $\boldsymbol{\theta}_k^*$ .  $\square$

Appendix A.3. Proof of Lemma 2

**Proof.** By the definition of  $\tilde{P}$ ,

$$\tilde{P} \{ \mathbf{PS}(\boldsymbol{\theta}^*) \} = \tau \| \mathbf{PS}(\boldsymbol{\theta}^*) \|_\infty + (1 - \tau) \max_j \left\{ \| \mathbf{P}_K \mathbf{S}(\boldsymbol{\alpha}^*) \|_2, \| \mathbf{P}_K \mathbf{S}(\boldsymbol{\beta}_j^*) \|_2 \right\},$$

where

$$\begin{aligned} \mathbf{P}_K \mathbf{S}(\boldsymbol{\alpha}^*) &= \mathbf{P}_K (\mathbf{E} \circ \phi'_q \{ \mathbf{F}(\boldsymbol{\theta}^*) \})^\top \mathbf{1}_K = \frac{1}{N} \sum_{i=1}^N \mathbf{P}_K \text{diag} \{ \mathbf{E}_i \} \phi'_q(\mathbf{F}_i^*), \\ \mathbf{P}_K \mathbf{S}(\boldsymbol{\beta}_j^*) &= \mathbf{P}_K (\mathbf{E} \circ \phi'_q \{ \mathbf{F}(\boldsymbol{\theta}^*) \})^\top \mathbf{x}_j = \frac{1}{N} \sum_{i=1}^N x_{ij} \mathbf{P}_K \text{diag} \{ \mathbf{E}_i \} \phi'_q(\mathbf{F}_i^*), \end{aligned}$$

$\mathbf{P}_K = (\mathbf{p}_1, \dots, \mathbf{p}_K)$  with  $\mathbf{p}_k = (p_{lk}) = \mathbb{1} \{ l = k \} - K^{-1}$ , and

$$\mathbb{E} \{ \mathbf{P}_K \mathbf{S}(\boldsymbol{\alpha}^*) \} = \mathbf{P}_K \mathbf{S}(\boldsymbol{\alpha}^*) = \mathbf{0}_K, \quad \mathbb{E} \{ \mathbf{P}_K \mathbf{S}(\boldsymbol{\beta}_j^*) \} = \mathbf{P}_K \mathbf{S}(\boldsymbol{\beta}_j^*) = \mathbf{0}_K.$$

Denoting

$$d_{ik} = \{ \mathbf{p}_k^\top \text{diag} \{ \mathbf{E}_i \} \phi'_q(\mathbf{F}_i^*) \} = \sum_{l=1}^K \left( \mathbb{1} \{ y_i = k \} - \frac{1}{K} \right) e_{il} \phi'_q(f_{il}^*),$$

we have that  $|d_{ik}| \leq 1 - K^{-1}$ . Note that the  $d_{ik}$  are  $N$  i.i.d. random variables with

$$\frac{1}{N} \sum_{i=1}^N \mathbb{E}(d_{ik}) = \mathbf{p}_k^\top \mathbf{S}(\boldsymbol{\alpha}^*) = 0 \text{ and } \frac{1}{N} \sum_{i=1}^N \mathbb{E}(d_{ik} x_{ij}) = \mathbf{p}_k^\top \mathbf{S}(\boldsymbol{\beta}_j^*) = 0.$$

By Hoeffding’s inequality, we have that

$$\Pr \left\{ \left| \mathbf{p}_k^\top \mathbf{S}(\boldsymbol{\alpha}^*) \right| > c_1 \left( 1 - \frac{1}{K} \right) \sqrt{\frac{2 \log(pK)}{N}} \right\} \leq 2(pK)^{-c_1^2}, \tag{A3}$$

where  $c_1 > 1$ .

Regarding the  $d_{ik}x_{ij}$ , we have that

$$\mathbb{E} \exp\{d_{ik}x_{ij}\} \leq \mathbb{E} \exp\{(1 - K^{-1})|x_{ij}|\} \leq \exp\{4(1 - K^{-1})^2 \zeta_1^2 \kappa^2\},$$

which implies that the  $d_{ik}x_{ij}$  are  $N$  independent sub-Gaussian random variables with variance proxy  $(1 - K^{-1})^2 \zeta_1^2 \kappa^2$ . Taking  $c_1 > 1$ , we have that

$$\Pr \left\{ \left| \mathbf{p}_k^\top \mathbf{S}(\boldsymbol{\beta}_j^*) \right| > c_1 \zeta_1 \kappa \left( 1 - \frac{1}{K} \right) \sqrt{\frac{2 \log(pK)}{N}} \right\} \leq 2(pK)^{-c_1^2}. \tag{A4}$$

Then by (A3) and (A4),

$$\Pr \left\{ \max_j \left\{ \left| \mathbf{p}_k^\top \mathbf{S}(\boldsymbol{\alpha}^*) \right|, \left| \mathbf{p}_k^\top \mathbf{S}(\boldsymbol{\beta}_j^*) \right| \right\} > \Lambda_1 \right\} \leq 2(pK)^{-c_1^2} \tag{A5}$$

with

$$\Lambda_1 = \max\{\zeta_1 \kappa, 1\} c_1 \left( 1 - \frac{1}{K} \right) \sqrt{\frac{2 \log(pK)}{N}}.$$

Taking a union bound over the  $Kp$  entries of  $\mathbf{PS}(\boldsymbol{\beta}^*)$  yields that

$$\begin{aligned} \Pr\{\|\mathbf{PS}(\boldsymbol{\beta}^*)\|_\infty \geq \Lambda_1\} &= \Pr \left\{ \max_{j,k} \left\{ \left| \frac{1}{N} \sum_{i=1}^N \mathbf{p}_k^\top \mathbf{S}(\boldsymbol{\alpha}^*) \right|, \left| \frac{1}{N} \sum_{i=1}^N \mathbf{p}_k^\top \mathbf{S}(\boldsymbol{\beta}_j^*) \right| \right\} \geq \Lambda_1 \right\} \\ &\leq 2K(p+1)(Kp)^{-c_1^2}. \end{aligned}$$

On one hand,

$$\|\mathbf{Pdiag}\{\mathbf{E}_i\} \phi'_q(\mathbf{F}_i^*)\|_2^2 = \|(\mathbf{E}_i - K^{-1}) \circ \phi'_q(\mathbf{F}_i^*)\|_2^2 \leq \sum_{l=1}^K (e_{il} - K^{-1})^2 \cdot 1 = 1 - K^{-1},$$

so for any  $\boldsymbol{\gamma} \in \mathbb{R}^K$ ,

$$|\boldsymbol{\gamma}^\top \mathbf{Pdiag}\{\mathbf{E}_i\} \phi'_q(\mathbf{F}_i^*)| \leq \|\boldsymbol{\gamma}\|_2 \sqrt{1 - \frac{1}{K}}$$

and  $\mathbb{E}\{\boldsymbol{\gamma}^\top \mathbf{Pdiag}\{\mathbf{E}_i\} \phi'_q(\mathbf{F}_i^*)\} = 0$ . Applying Hoeffding’s lemma,

$$\mathbb{E} \exp\{\boldsymbol{\gamma}^\top \mathbf{P}_K \mathbf{S}(\boldsymbol{\alpha}^*)\} = \prod_{i=1}^N \mathbb{E} \exp \left\{ \frac{1}{N} \boldsymbol{\gamma}^\top \mathbf{P}_K \text{diag}\{\mathbf{E}_i\} \phi'_q(\mathbf{F}_i^*) \right\} \leq \exp \left\{ \frac{\|\boldsymbol{\gamma}\|_2^2}{2N} \left( 1 - \frac{1}{K} \right) \right\}.$$

Applying a square root to Theorem 2.1 of [31] with  $c_2 > 1$ , we have that

$$\Pr \left\{ \|\mathbf{PS}(\boldsymbol{\alpha}^*)\|_2 \geq \sqrt{\frac{K-1}{N}} + c_2 \sqrt{\left( 1 - \frac{1}{K} \right) \frac{2 \log(p)}{N}} \right\} \leq p^{-c_2^2}. \tag{A6}$$

On the other hand, since the  $x_{ij}$  are  $N$  independent sub-Gaussian random variables with variance proxy  $\zeta_1^2 \kappa^2$ ,

$$\begin{aligned} \mathbb{E} \exp\{\gamma^\top \mathbf{PS}(\beta_j^*)\} &= \prod_{i=1}^N \mathbb{E} \exp\left\{\frac{x_{ij}}{N} \{\gamma^\top \mathbf{P} \text{diag}\{\mathbf{E}_i\} \phi'_q(\mathbf{F}_i^*)\}\right\} \\ &\leq \prod_{i=1}^N \mathbb{E} \exp\left\{\sqrt{1 - \frac{1}{K}} \frac{\|\gamma\|_2}{N} |x_{ij}|\right\} \\ &= \exp\left\{\frac{\|\gamma\|_2^2}{2} \left(1 - \frac{1}{K}\right) \frac{8\zeta_1^2 \kappa^2}{N}\right\} \end{aligned}$$

and  $\mathbb{E}\{\mathbf{P}_K \mathbf{S}(\beta_j^*)\} = \mathbf{0}_K$ . Similarly, we have that

$$\Pr\left\{\|\mathbf{PS}(\beta_j^*)\|_2 \geq 2\sqrt{2}\zeta_1 \kappa \left\{\sqrt{\frac{K-1}{N}} + c_2 \sqrt{\left(1 - \frac{1}{K}\right) \frac{2\log(p)}{N}}\right\}\right\} \leq p^{-c_2^2} \tag{A7}$$

for a constant  $c_2 > 1$ .

Therefore, by (A6) and (A7),

$$\Pr\left\{\max_j \{\|\mathbf{PS}(\alpha^*)\|_2, \|\mathbf{PS}(\beta_j^*)\|_2\} \geq \Lambda_2\right\} \leq p^{-c_2^2}$$

with

$$\Lambda_2 = \max\{2\sqrt{2}\zeta_1 \kappa, 1\} \left\{\sqrt{\frac{K-1}{N}} + c_2 \sqrt{\left(1 - \frac{1}{K}\right) \frac{2\log(p)}{N}}\right\}.$$

Applying the union bound to (A5), it follows that

$$\Pr\left\{\tilde{P}\{\mathbf{PS}(\vartheta^*)\} \geq \tau\Lambda_1 + (1 - \tau)\Lambda_2\right\} \leq 2K(p+1)(pK)^{1-c_1^2} + p^{1-c_2^2},$$

and the desired result follows.  $\square$

#### Appendix A.4. Proof of Theorem 2

**Lemma A6.** Suppose that  $\lambda = c_0 \sqrt{\frac{\log(pK)}{N}}$ . Then  $\hat{\vartheta} - \vartheta^* \in \mathcal{U}$ , where

$$\mathcal{U} = \left\{\delta \in \mathbb{R}^{K(p+1)} \mid \frac{\tau}{1-\tau} \|\delta_{\mathcal{E}_+}\|_1 + \sum_{j \in \mathcal{G}_+} \|\delta_j\|_2 \geq C_0 \left(\frac{\tau}{1-\tau} \|\delta_{\mathcal{E}^c}\|_1 + \sum_{j \notin \mathcal{G}} \|\delta_j\|_2\right)\right\},$$

$C_0 = \frac{(c_0-1)}{(c_0+1)}$ ,  $\mathcal{E}^c$  denotes the complement of  $\mathcal{E}$ ,  $\mathcal{E}_+ = \mathcal{E} \cup \{l = 1 + (k-1)(p+1) \mid k = 1, \dots, K\}$ , and  $\mathcal{G}_+ = \mathcal{G} \cup \{0\}$ .

**Proof.** Since  $\hat{\vartheta} = \vartheta^* + \delta$  is the minimizer, we have that

$$\begin{aligned} L(\vartheta^*) + \lambda P(\beta^*) &\geq L(\hat{\vartheta}) + \lambda P(\hat{\beta}) \\ \lambda \{P(\beta^*) - P(\beta^* + \tilde{\delta})\} &\geq L(\vartheta^* + \delta) - L(\vartheta^*), \end{aligned} \tag{A8}$$

where  $\beta^*$  is the vector  $\vartheta^*$  without the  $a_k$  components, replacing  $\tilde{\delta}$  for  $\delta$ . Then



$$\begin{aligned}
 P(\boldsymbol{\beta}^*) - P(\boldsymbol{\beta}^* + \tilde{\boldsymbol{\delta}}) &= \tau(\|\boldsymbol{\beta}_{\mathcal{E}}^*\|_1 - \|\boldsymbol{\beta}_{\mathcal{E}}^* + \tilde{\boldsymbol{\delta}}_{\mathcal{E}}\|_1 - \|\tilde{\boldsymbol{\delta}}_{\mathcal{E}^c}\|_1) \\
 &\quad + (1 - \tau)\left(\sum_{j \in \mathcal{G}} \|\boldsymbol{\beta}_j^*\|_2 - \sum_{j \in \mathcal{G}} \|\boldsymbol{\beta}_j^* + \boldsymbol{\delta}_j\|_2 - \sum_{j \notin \mathcal{G}} \|\boldsymbol{\delta}_j^*\|_2\right) \\
 &\leq \tau(\|\tilde{\boldsymbol{\delta}}_{\mathcal{E}}\|_1 - \|\tilde{\boldsymbol{\delta}}_{\mathcal{E}^c}\|_1) + (1 - \tau)\left(\sum_{j \in \mathcal{G}} \|\boldsymbol{\delta}_j\|_2 - \sum_{j \notin \mathcal{G}} \|\boldsymbol{\delta}_j^*\|_2\right) \\
 &\leq \tau(\|\boldsymbol{\delta}_{\mathcal{E}_+}\|_1 - \|\boldsymbol{\delta}_{\mathcal{E}^c}\|_1) + (1 - \tau)\left(\sum_{j \in \mathcal{G}_+} \|\boldsymbol{\delta}_j\|_2 - \sum_{j \notin \mathcal{G}} \|\boldsymbol{\delta}_j\|_2\right).
 \end{aligned}$$

By the convexity of  $L$ ,

$$L(\boldsymbol{\theta}^* + \boldsymbol{\delta}) - L(\boldsymbol{\theta}^*) \geq \langle \mathbf{S}(\boldsymbol{\theta}^*), \boldsymbol{\delta} \rangle \geq -\bar{P}\{\mathbf{PS}(\boldsymbol{\theta}^*)\}P(\boldsymbol{\delta}) \geq -\frac{\lambda}{c_0}P(\boldsymbol{\delta}).$$

Note that

$$P(\boldsymbol{\delta}) = \tau(\|\boldsymbol{\delta}_{\mathcal{E}_+}\|_1 + \|\boldsymbol{\delta}_{\mathcal{E}^c}\|_1) + (1 - \tau)\left(\sum_{j \in \mathcal{G}_+} \|\boldsymbol{\delta}_j\|_2 + \sum_{j \notin \mathcal{G}} \|\boldsymbol{\delta}_j\|_2\right).$$

Combining the above results, we have that

$$\begin{aligned}
 \lambda\{P(\boldsymbol{\theta}^*) - P(\boldsymbol{\theta}^* + \boldsymbol{\delta})\} &\geq \{L(\boldsymbol{\theta}^* + \boldsymbol{\delta}) - L(\boldsymbol{\theta}^*)\} \\
 (c + 1)\tau\|\boldsymbol{\delta}_{\mathcal{E}_+}\|_1 + (1 - \tau)\sum_{j \in \mathcal{G}_+} \|\boldsymbol{\delta}_j\|_2 &\geq (c - 1)\tau\|\boldsymbol{\delta}_{\mathcal{E}^c}\|_1 + (1 - \tau)\sum_{j \notin \mathcal{G}} \|\boldsymbol{\delta}_j\|_2 \\
 \frac{\tau}{1 - \tau}\|\boldsymbol{\delta}_{\mathcal{E}_+}\|_1 + \sum_{j \in \mathcal{G}_+} \|\boldsymbol{\delta}_j\|_2 &\geq C_0\left(\frac{\tau}{1 - \tau}\|\boldsymbol{\delta}_{\mathcal{E}^c}\|_1 + \sum_{j \notin \mathcal{G}} \|\boldsymbol{\delta}_j\|_2\right).
 \end{aligned}$$

□

**Lemma A7.** Assume that conditions (A1)-(A3) are satisfied. Then

$$\sup_{\mathbf{v} \in \mathcal{V}} \frac{|\Delta L(\mathbf{u}, \mathbf{v}) - \mathbb{E}\{\Delta L(\mathbf{u}, \mathbf{v})\}|}{\|\mathbf{v}\|_2} > \Lambda_3$$

with probability at most  $2(Kp)^{2(s_e + K)(1 - c_3^2)}$ , where

$$\Lambda_3 = (1 + \sqrt{2}c_3)\zeta_2 \sqrt{\frac{2(s_e + K) \log(pK)}{N}}$$

and  $\Delta L(\mathbf{u}, \mathbf{v}) = L(\mathbf{u} + \mathbf{v}) - L(\mathbf{u})$  for any  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{K(p+1)}$  and for some constant  $c_3 > 1$ .

**Proof.** Given any  $\mathbf{u} \in \mathbb{R}^{K(p+1)}$  and  $\mathbf{v} \in \mathcal{V}$  with  $\mathcal{V} = \{\mathbf{v} \in \mathbb{R}^{K(p+1)} | 0 < \|\mathbf{v}\|_0 \leq s_e + K\}$ ,

$$\begin{aligned}
 \Delta L(\mathbf{u}, \mathbf{v}) &= \frac{1}{N} \sum_{i=1}^N \mathbf{E}_i^\top \left( \phi_q\{(\mathbf{U} + \mathbf{V})^\top \mathbf{Z}_i\} - \phi_q\{\mathbf{U}^\top \mathbf{Z}_i\} \right) \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K e_{ik} \left( \phi_q\{\mathbf{Z}_i^\top (\mathbf{u}_k + \mathbf{v}_k)\} - \phi_q\{\mathbf{Z}_i^\top (\mathbf{u}_k)\} \right) \\
 &= \frac{1}{N} \sum_{i=1}^N d_i(\mathbf{u}, \mathbf{v}),
 \end{aligned}$$

where  $\mathbf{u} = \text{vec}\{\mathbf{U}\}$ ,  $\mathbf{v} = \text{vec}\{\mathbf{V}\}$ .

The bounded gradient implies the Lipschitz continuity of  $\phi_q$  so that  $|\phi_q(u + v) - \phi_q(u)| \leq |v|$ . Since  $e_{ik} \in \{0, 1\}$ , we have that

$$\begin{aligned} |d_i(\mathbf{u}, \mathbf{v})| &\leq \sum_{k=1}^K \left| e_{ik} \{ \phi_q \{ \mathbf{Z}_i^\top (\mathbf{u}_k + \mathbf{v}_k) \} - \phi_q(\mathbf{Z}_i^\top \mathbf{u}_k) \} \right| \\ &\leq \sum_{k=1}^K |e_{ik} \mathbf{Z}_i^\top \mathbf{v}_k| \leq \mathbf{E}_i^\top \text{vec} \{ \mathbf{V}^\top \mathbf{Z}_i \} \\ &= \mathbf{v}^\top (\mathbf{Z}_i \otimes \mathbf{I}_K) \mathbf{E}_i. \end{aligned}$$

Note that

$$\sum_{i=1}^N (\mathbf{v}^\top (\mathbf{Z}_i \otimes \mathbf{I}_K) \mathbf{E}_i)^2 = \|\text{diag}\{\text{vec}\{\mathbf{E}^\top\}\} (\mathbf{Z} \otimes \mathbf{I}_K) \mathbf{v}\|_2^2.$$

By Hoeffding’s inequality, we have that

$$\begin{aligned} &\Pr \left\{ \left| \frac{1}{N} \sum_{i=1}^N d_i(\mathbf{u}, \mathbf{v}) - \mathbb{E} \left( \frac{1}{N} \sum_{i=1}^N d_i(\mathbf{u}, \mathbf{v}) \right) \right| > t \right\} \\ &\leq 2 \exp \left\{ - \frac{2N^2 t^2}{4 \|\text{diag}\{\text{vec}\{\mathbf{E}^\top\}\} (\mathbf{Z} \otimes \mathbf{I}_K) \mathbf{v}\|_2^2} \right\} \\ &\leq 2 \exp \left\{ - \frac{Nt^2}{2\zeta_2^2 \|\mathbf{v}\|_2^2} \right\}. \end{aligned}$$

Thus  $\Pr\{R(\mathbf{v}) > \Lambda_3\} \leq 2(Kp)^{-(s_e+K)c_3^2}$  with

$$R(\mathbf{v}) = \frac{|\Delta L(\mathbf{u}, \mathbf{v}) - \mathbb{E}\{\Delta L(\mathbf{u}, \mathbf{v})\}|}{\|\mathbf{v}\|_2} \text{ and } \Lambda_3 = c_3 \zeta_2 \sqrt{\frac{2(s_e + K) \log(pk)}{N}}.$$

Next, we consider covering  $\mathcal{V}$  with  $\epsilon$ -balls such that for any  $\mathbf{v}_1$  and  $\mathbf{v}_2$  in the same ball,  $\left| \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|_2} - \frac{\mathbf{v}_2}{\|\mathbf{v}_2\|_2} \right| \leq \epsilon$ , where  $\epsilon$  is a small positive number. The number of  $\epsilon$ -balls required to cover a  $m$ -dimensional unit ball is bounded by  $(\frac{2}{\epsilon} + 1)^m$ . Then for those  $\frac{\mathbf{v}}{\|\mathbf{v}\|_2}$ , we require a covering number of at most  $(3(Kp)/\epsilon)^{s_e+K}$ . Let  $\mathcal{N}$  denote such an  $\epsilon$ -net. We have that

$$\Pr \left\{ \sup_{\mathbf{v} \in \mathcal{N}} R(\mathbf{v}) > \Lambda_3 \right\} \leq \left( \frac{3Kp}{\epsilon} \right)^{s_e+K} 2(Kp)^{-(s_e+K)c_3^2} = 2 \left\{ \frac{3}{\epsilon} (Kp)^{1-c_3^2} \right\}^{s_e+K}.$$

Furthermore, for any  $\mathbf{v}_1, \mathbf{v}_2 \in \mathcal{V}$ ,

$$\begin{aligned} |R(\mathbf{v}_1) - R(\mathbf{v}_2)| &\leq \frac{2}{N} \left\| \text{diag}\{\text{vec}\{\mathbf{E}^\top\}\} (\mathbf{Z} \otimes \mathbf{I}_K) \left( \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|_2} - \frac{\mathbf{v}_2}{\|\mathbf{v}_2\|_2} \right) \right\|_1 \\ &\leq \frac{2}{\sqrt{N}} \left\| \text{diag}\{\text{vec}\{\mathbf{E}^\top\}\} (\mathbf{Z} \otimes \mathbf{I}_K) \left( \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|_2} - \frac{\mathbf{v}_2}{\|\mathbf{v}_2\|_2} \right) \right\|_2 \\ &\leq 2\zeta_2 \epsilon. \end{aligned}$$

Therefore  $\sup_{\mathbf{v} \in \mathcal{V}} R(\mathbf{v}) \leq \sup_{\mathbf{v} \in \mathcal{N}} R(\mathbf{v}) + 2\zeta_2 \epsilon$ . Taking  $\epsilon = \sqrt{\frac{(s_e + K) \log(pk)}{2N}}$ , we have that

$$\begin{aligned} \Pr \left\{ \sup_{\mathbf{v} \in \mathcal{V}} R(\mathbf{v}) > \Lambda_3 \right\} &\leq \Pr \left\{ \sup_{\mathbf{v} \in \mathcal{N}} R(\mathbf{v}) > (c_3 - 1)\zeta_1 \sqrt{\frac{2(s_e + K) \log(pK)}{N}} \right\} \\ &\leq 2 \left\{ \sqrt{\frac{2N}{(s_e + K) \log(pK)}} 3(Kp)^{1-(c_3-1)^2} \right\}^{s_e+K} \\ &\leq 2 \left\{ (Kp)^{2-(c_3-1)^2} \right\}^{s_e+K}. \end{aligned}$$

Setting  $c_3 = 1 + \sqrt{2}c_4$  and  $c_4 > 1$ , we obtain the desired result that

$$\Pr \left\{ \sup_{\mathbf{v} \in \mathcal{V}} R(\mathbf{v}) > (1 + \sqrt{2}c_4)\zeta_2 \sqrt{\frac{2(s_e + K) \log(pK)}{N}} \right\} \leq 2(Kp)^{2(s_e+K)(1-c_4^2)}.$$

□

**Proof of Theorem 2.** Consider a disjoint partition on the coordinate set  $\delta = \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*$ , that is,  $\delta = \sum_{m=1}^M \mathbf{v}_m$  with  $\mathbf{v}_m \in \mathcal{V}$ . Note that, each subvector  $\mathbf{v}_m$  has at most  $s_e + K$  non-zero coordinates. Denote  $\mathbf{v}_0 = \mathbf{0}$  and  $\mathbf{u}_m = \boldsymbol{\theta}^* + \sum_{l=0}^{m-1} \mathbf{v}_l$  so that  $\mathbf{u}_1 = \boldsymbol{\theta}^*$  and  $\mathbf{u}_M + \mathbf{v}_M = \boldsymbol{\theta}^* + \delta$ . We have the decomposition

$$\begin{aligned} \Delta L(\boldsymbol{\theta}^*, \delta) &= L\left(\boldsymbol{\theta}^* + \sum_{m=1}^M \mathbf{v}_m\right) - L(\boldsymbol{\theta}^*) = \sum_{m=1}^M L\left(\boldsymbol{\theta}^* + \sum_{l=0}^m \mathbf{v}_l\right) - L\left(\boldsymbol{\theta}^* + \sum_{l=0}^{m-1} \mathbf{v}_l\right) \\ &= \sum_{m=1}^M L(\mathbf{u}_m + \mathbf{v}_m) - L(\mathbf{u}_m) = \sum_{m=1}^M \Delta L(\mathbf{u}_m, \mathbf{v}_m). \end{aligned}$$

By Lemma A7,

$$\sum_{m=1}^M \Delta L(\mathbf{u}_m, \mathbf{v}_m) \geq \sum_{m=1}^M \mathbb{E}\{\Delta L(\mathbf{u}_m, \mathbf{v}_m)\} - \Lambda_3 \|\mathbf{v}_m\|_2 = \mathbb{E}\{\Delta L(\boldsymbol{\theta}^*, \delta)\} - \Lambda_3 \|\delta\|_2$$

with high probability. By Lemma A5,  $\mathcal{L}$  is twice differentiable so that

$$\begin{aligned} \mathbb{E}\{\Delta L(\boldsymbol{\theta}^*, \delta)\} &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}\left(\mathbf{E}_i^\top \phi_q\{F_i(\boldsymbol{\theta}^* + \delta)\}\right) - \mathbb{E}\left(\mathbf{E}_i^\top \phi_q\{F_i(\boldsymbol{\theta}^*)\}\right) \\ &= \mathcal{L}(\boldsymbol{\theta}^* + \delta) - \mathcal{L}(\boldsymbol{\theta}^*) \\ &= \mathcal{S}(\boldsymbol{\theta}^*)^\top \delta + \frac{1}{2} \delta^\top \mathcal{H}(\boldsymbol{\theta}^*) \delta + o(\|\delta\|_2^2) \\ &\geq 0 + \frac{\zeta_3^2}{2} \|\delta\|_2^2 + o(\|\delta\|_2^2). \end{aligned}$$

Consequently,  $\Delta L(\boldsymbol{\theta}^*, \delta)$  is bounded below by  $\frac{\zeta_3^2}{2} \|\delta\|_2^2 - \Lambda_3 \|\delta\|_2$  with high probability.

Note that

$$\begin{aligned} P(\boldsymbol{\beta}^*) - P(\boldsymbol{\beta}^* + \delta) &\leq \tau(\|\delta_{\mathcal{G}_+}\|_1 - \|\delta_{\mathcal{G}^c}\|_1) + (1 - \tau) \left( \sum_{j \in \mathcal{G}_+} \|\delta_j\|_2 - \sum_{j \notin \mathcal{G}_+} \|\delta_j\|_2 \right) \\ &\leq \left( \tau \|\delta_{\mathcal{G}_+}\|_1 + (1 - \tau) \sum_{j \in \mathcal{G}_+} \|\delta_j\|_2 \right). \end{aligned}$$

From (A8),

$$\begin{aligned} L(\boldsymbol{\theta}^*) + \lambda P(\boldsymbol{\beta}^*) &\geq L(\hat{\boldsymbol{\theta}}) + \lambda P(\hat{\boldsymbol{\beta}}) \\ \lambda \{P(\boldsymbol{\beta}^*) - P(\boldsymbol{\beta}^* + \tilde{\boldsymbol{\delta}})\} &\geq L(\boldsymbol{\theta}^* + \boldsymbol{\delta}) - L(\boldsymbol{\theta}^*) \\ \lambda \left( \tau \|\boldsymbol{\delta}_{\mathcal{E}_+}\|_1 + (1 - \tau) \sum_{j \in \mathcal{G}_+} \|\boldsymbol{\delta}_j\|_2 \right) &\geq \frac{\zeta_3^2}{2} \|\boldsymbol{\delta}\|_2^2 - \Lambda_3 \|\boldsymbol{\delta}\|_2. \end{aligned}$$

Clearly,  $\|\boldsymbol{\delta}_{\mathcal{E}_+}\|_1 \leq \sqrt{s_e + K} \|\boldsymbol{\delta}_{\mathcal{E}_+}\|_2 \leq \sqrt{s_e + K} \|\boldsymbol{\delta}\|_2$  and  $\sum_{j \in \mathcal{G}_+} \|\boldsymbol{\delta}_j\|_2 \leq \sqrt{s_g + 1} \|\boldsymbol{\delta}\|_2$ . We conclude that

$$\begin{aligned} \frac{\zeta_3^2}{2} \|\boldsymbol{\delta}\|_2^2 &\leq \lambda \left( \tau \|\boldsymbol{\delta}_{\mathcal{E}_+}\|_1 + (1 - \tau) \sum_{j \in \mathcal{G}_+} \|\boldsymbol{\delta}_j\|_2 \right) + \Lambda_3 \|\boldsymbol{\delta}\|_2 \\ \|\boldsymbol{\delta}\|_2^2 &\leq 2\zeta_3^{-2} \left\{ \lambda \left( \tau \sqrt{s_e + K} + (1 - \tau) \sqrt{s_g + 1} \right) + \Lambda_3 \right\} \|\boldsymbol{\delta}\|_2, \end{aligned}$$

after which the desired result follows from straightforward algebraic manipulation.  $\square$

Appendix A.5. Proof of Lemma 3

**Proof.** Since

$$\text{vec}(\mathbf{F}^\top)^\top = \text{vec}\{(\mathbf{1}_N \boldsymbol{\alpha}^\top + \mathbf{X}\mathbf{B})^\top\}^\top = \boldsymbol{\alpha}^\top (\mathbf{1}_N^\top \otimes \mathbf{I}_K) + \text{vec}(\mathbf{B}^\top)^\top (\mathbf{X}^\top \otimes \mathbf{I}_K),$$

we have that

$$\begin{cases} \frac{\partial \text{vec}(\mathbf{F}^\top)^\top}{\partial \boldsymbol{\alpha}} = \frac{\partial \boldsymbol{\alpha}^\top (\mathbf{1}_N^\top \otimes \mathbf{I}_K)}{\partial \boldsymbol{\alpha}} = \frac{\partial \boldsymbol{\alpha}^\top}{\partial \boldsymbol{\alpha}} (\mathbf{1}_N^\top \otimes \mathbf{I}_K) = \mathbf{I}_K (\mathbf{1}_N^\top \otimes \mathbf{I}_K) = \mathbf{1}_N^\top \otimes \mathbf{I}_K \\ \frac{\partial \text{vec}(\mathbf{F}^\top)^\top}{\partial \text{vec}(\mathbf{B}^\top)} = \frac{\partial \text{vec}(\mathbf{B}^\top)^\top (\mathbf{X}^\top \otimes \mathbf{I}_K)}{\partial \text{vec}(\mathbf{B}^\top)} = \mathbf{I}_{pK} (\mathbf{X}^\top \otimes \mathbf{I}_K) = \mathbf{X}^\top \otimes \mathbf{I}_K. \end{cases}$$

The derivative with respect to  $\boldsymbol{\alpha}$  is

$$\begin{aligned} \mathbf{N}\mathbf{S}(\boldsymbol{\alpha}) &= N \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\alpha}} = \frac{\partial}{\partial \boldsymbol{\alpha}} \text{vec}\{\mathbf{E}^\top\}^\top \text{vec}\{\phi_q(\mathbf{F}^\top)\} \\ &= \frac{\partial \text{vec}(\mathbf{F}^\top)^\top}{\partial \boldsymbol{\alpha}} \frac{\partial \phi_q\{\text{vec}(\mathbf{F}^\top)^\top\}}{\partial \text{vec}(\mathbf{F}^\top)} \text{vec}\{\mathbf{E}^\top\} \\ &= (\mathbf{1}_N^\top \otimes \mathbf{I}_K) \text{diag}\{\text{vec}\{\phi'_q(\mathbf{F}^\top)\}\} \text{vec}\{\mathbf{E}^\top\} \\ &= \text{vec}\{\mathbf{I}_K \{ \mathbf{E} \circ \phi'_q(\mathbf{F}) \}^\top \mathbf{1}_N\} \\ &= \{ \mathbf{E} \circ \phi'_q(\mathbf{F}) \}^\top \mathbf{1}_N. \end{aligned}$$

Thus,

$$\begin{aligned} \|\mathbf{S}(\boldsymbol{\alpha})\|_v^u \|^2 &= \|\mathbf{S}(\mathbf{u}) - \mathbf{S}(\mathbf{v})\|_2^2 = N^{-2} \|(\mathbf{1}_N^\top \otimes \mathbf{I}_K) \text{vec}\{(\mathbf{E} \circ \phi'_q\{\mathbf{F}(\boldsymbol{\alpha})\})\|_v^u \}^\top\|_2^2 \\ &\leq N^{-2} \|\mathbf{1}_N^\top \otimes \mathbf{I}_K\|_2^2 \|\text{vec}\{(\mathbf{E} \circ \phi'_q\{\mathbf{F}(\boldsymbol{\alpha})\})\|_v^u \}^\top\|_2^2 \\ &= N^{-1} \sum_{k=1}^K \sum_{i=1}^N e_{ik}^2 (\phi'_q\{f_{ik}(u_k)\} - \phi'_q\{f_{ik}(v_k)\})^2 \\ &\leq N^{-1} \sum_{k=1}^K \left( \sum_{i=1}^N e_{ik} \right) L_q^2(u_k - v_k)^2 \\ &\leq N^{-1} n_{\max} L_q^2 \|\mathbf{u} - \mathbf{v}\|_2^2, \end{aligned}$$

where  $L_q = \frac{(q+1)^2}{q}$  is the Lipschitz constant of  $\phi'_q$ . We have that  $L_\alpha = \sqrt{\frac{n_{\max}}{N}} L_q$ .

The derivative with respect to  $\text{vec}(\mathbf{B}^\top)$  is

$$\begin{aligned} N \frac{\partial L(\theta)}{\partial \text{vec}(\mathbf{B}^\top)} &= \frac{\partial}{\partial \text{vec}(\mathbf{B}^\top)} \text{vec}\{\mathbf{E}^\top\}^\top \text{vec}\{\phi_q(\mathbf{F}^\top)\} \\ &= \frac{\partial \text{vec}(\mathbf{F}^\top)^\top}{\partial \text{vec}(\mathbf{B}^\top)} \frac{\partial \phi_q\{\text{vec}(\mathbf{F}^\top)^\top\}}{\partial \text{vec}(\mathbf{F}^\top)} \text{vec}\{\mathbf{E}^\top\} \\ &= (\mathbf{X}^\top \otimes \mathbf{I}_K) \text{diag}(\text{vec}\{\phi'_q(\mathbf{F}^\top)\}) \text{vec}\{\mathbf{E}^\top\} \\ &= \text{vec}\left(\mathbf{I}_K \{\mathbf{E} \circ \phi'_q(\mathbf{F})\}^\top \mathbf{X}\right) \\ &= \text{vec}\left(\{\mathbf{E} \circ \phi'_q(\mathbf{F})\}^\top \mathbf{X}\right). \end{aligned}$$

Therefore, the derivative with respect to  $\mathbf{B}$  is  $\mathbf{S}(\mathbf{B}) = N^{-1} \mathbf{X}^\top \{\mathbf{E} \circ \phi'_q(\mathbf{F})\}$ . Note that

$$\begin{aligned} \text{vec}\left(\mathbf{X}^\top \{\mathbf{E} \circ \phi'_q(\mathbf{F})\}\right) &= (\mathbf{I}_K \otimes \mathbf{X}^\top) \text{diag}\{\text{vec}(\mathbf{E})\} \text{vec}\{\phi'_q(\mathbf{F})\} \\ &= \left\{ \bigoplus_{k=1}^K \mathbf{X}^\top \text{diag}(e_k) \right\} \text{vec}\{\phi'_q(\mathbf{F})\} \end{aligned}$$

and

$$\sum_{i=1}^N \{e_{ik} \mathbf{X}_i^\top (\mathbf{u}_k - \mathbf{v}_k)\}^2 = \|\text{diag}(e_k) \mathbf{X} (\mathbf{u}_k - \mathbf{v}_k)\|_2^2 \leq \|\text{diag}(e_k) \mathbf{X}\|_2^2 \|\mathbf{u}_k - \mathbf{v}_k\|_2^2;$$

thus

$$\begin{aligned} N^2 \|\text{vec}\{\mathbf{S}(\mathbf{U}) - \mathbf{S}(\mathbf{V})\}\|_2^2 &= \sum_{k=1}^K \left\| \mathbf{X}^\top \text{diag}(e_k) \phi_q\{f_k(\mathbf{b}_k)\} \Big|_{\mathbf{v}_k}^{\mathbf{u}_k} \right\|_2^2 \\ &\leq \sum_{k=1}^K \|\mathbf{X}^\top \text{diag}(e_k)\|_2^2 \|\text{diag}(e_k) \phi_q\{f_k(\mathbf{b}_k)\}\|_{\mathbf{v}_k}^{\mathbf{u}_k}{}^2 \\ &\leq \sum_{k=1}^K \|\text{diag}(e_k) \mathbf{X}\|_2^2 \sum_{i=1}^N e_{ik} (\phi_q\{f_{ik}(\mathbf{u}_k)\} - \phi_q\{f_{ik}(\mathbf{v}_k)\})^2 \\ &\leq L_q^2 \sum_{k=1}^K \|\text{diag}(e_k) \mathbf{X}\|_2^2 \sum_{i=1}^N \{e_{ik} \mathbf{X}_i^\top (\mathbf{u}_k - \mathbf{v}_k)\}^2 \\ &\leq L_q^2 \sum_{k=1}^K \|\text{diag}(e_k) \mathbf{X}\|_2^4 \|\mathbf{u}_k - \mathbf{v}_k\|_2^2 \\ &\leq \max_k \left\{ \|\text{diag}(e_k) \mathbf{X}\|_2^2 \right\}^2 \|\text{vec}(\mathbf{U} - \mathbf{V})\|_2^2. \end{aligned}$$

We conclude that  $L_B = L_q N^{-1} \max_k \|\text{diag}(e_k) \mathbf{X}\|_2^2$ .  $\square$

Appendix A.6. Proof of Theorem 3

**Lemma A8.** The indicator function

$$\delta_{\mathcal{R}}(x) = \begin{cases} 0, & \text{if } x \in \mathcal{R} \\ \infty, & \text{if } x \notin \mathcal{R}, \end{cases}$$

where  $\mathcal{R} = \{x \in \mathbb{R}^p \mid \mathbf{1}_p^\top x = 0\}$ , has subdifferential

$$\partial\delta_{\mathcal{R}}(x) = \begin{cases} \{\mathbf{g} \in \mathbb{R}^p \mid \mathbf{g} = s\mathbf{1}_p, s \in \mathbb{R}\}, & \text{if } x \in \mathcal{R} \\ \emptyset, & \text{if } x \notin \mathcal{R}. \end{cases}$$

**Proof.** Suppose that  $x \in \mathcal{R}$ . Then  $\mathbf{g} \in \partial\delta_{\mathcal{R}}(x)$  if and only if both

$$\delta_{\mathcal{R}}(\mathbf{y}) \geq \delta_{\mathcal{R}}(x) + \langle \mathbf{g}, \mathbf{y} - x \rangle \text{ for all } \mathbf{y} \in \mathcal{R} \text{ and} \\ \omega^\top(\mathbf{y} - x) \leq 0.$$

Let  $z = \mathbf{y} - x$ . Then  $z \in \mathcal{R}$  since  $\mathbf{1}_p^\top(\mathbf{y} - x) = 0$ . Thus,  $\mathbf{g}^\top z \leq 0$ . If  $\mathbf{g}^\top z = 0$ , then  $\mathbf{g} \in \{\mathbf{g} \in \mathbb{R}^p \mid \mathbf{g} = s\mathbf{1}_p, s \in \mathbb{R}\}$ . If there exists  $\mathbf{g} \in \partial\delta_{\mathcal{R}}(x)$  satisfying  $\mathbf{g}^\top z < 0$  for some  $z \in \mathcal{R}$ , then  $-z \in \mathcal{R}$ , so we must have that  $\mathbf{g}^\top z > 0$ . This is a contradiction.

Now, for any  $x \notin \mathcal{R}$ , we have that  $\mathbf{g} \in \partial\delta_{\mathcal{R}}(x)$  if and only if both

$$\delta_{\mathcal{R}}(\mathbf{y}) \geq \delta_{\mathcal{R}}(x) + \langle \mathbf{g}, \mathbf{y} - x \rangle \text{ for all } \mathbf{y} \in \mathcal{R} \text{ and} \\ \omega^\top(x - \mathbf{y}) \geq \infty.$$

For  $x \notin \mathcal{R}$  and  $\mathbf{y} \in \mathcal{R}$ , since  $z = x - \mathbf{y} \in \mathbb{R}^p$  and  $\mathbf{g}^\top z \geq \infty$ , it must be that  $\mathbf{g} \in \emptyset$ .  $\square$

**Proof of Theorem 3.** It is sufficient to minimize the objective function

$$G(\boldsymbol{\beta}) = \frac{1}{2}\|\boldsymbol{\beta} - \boldsymbol{\beta}_*\|_2^2 + \rho_1\|\boldsymbol{\beta}\|_1 + \rho_2\|\boldsymbol{\beta}\|_2 + \delta_{\mathcal{R}}(\boldsymbol{\beta}),$$

where  $\mathcal{R} = \{x \in \mathbb{R}^K \mid \mathbf{1}_K^\top x = 0\}$ . Then the subdifferential of  $G(\boldsymbol{\beta})$  is

$$\partial G(\boldsymbol{\beta}) = \boldsymbol{\beta} - \boldsymbol{\beta}_* + \rho_1\partial\|\boldsymbol{\beta}\|_1 + \rho_2\partial\|\boldsymbol{\beta}\|_2 + \partial\delta_{\mathcal{R}}(\boldsymbol{\beta}).$$

For an optimal solution  $\boldsymbol{\beta}^* \in \mathcal{R}$ , we have that  $\mathbf{0}_p \in \partial G(\boldsymbol{\beta}^*)$  if and only if there exist  $\mathbf{u} \in \partial\|\boldsymbol{\beta}\|_1$ ,  $\mathbf{v} \in \partial\|\boldsymbol{\beta}\|_2$  and  $s \in \mathbb{R}$  such that  $\boldsymbol{\beta}^* = \boldsymbol{\beta}_* - \rho_1\mathbf{u} - \rho_2\mathbf{v} - s\mathbf{1}_p$ . Since  $\mathbf{1}^\top \boldsymbol{\beta}^* = 0$ , we have that  $s = p^{-1}\mathbf{1}_p^\top(\boldsymbol{\beta}_* - \rho_1\mathbf{u} - \rho_2\mathbf{v})$ , so

$$\boldsymbol{\beta}^* = \mathbf{P}_K(\boldsymbol{\beta}_* - \rho_1\mathbf{u} - \rho_2\mathbf{v}).$$

If  $\boldsymbol{\beta}^* = \mathbf{0}_p$ , then  $|u_j| < 1$  for  $j = 1, \dots, p$ ,  $\|\mathbf{v}\|_2 \leq 1$  and

$$\|\mathbf{P}_K(\boldsymbol{\beta}_* - \rho_1\mathbf{u})\|_2 = \rho_2\|\mathbf{P}_K\mathbf{v}\|_2 \leq \rho_2\|\mathbf{P}_K\|_2\|\mathbf{v}\|_2 = \rho_2\|\mathbf{v}\|_2 \leq \rho_2;$$

If  $\boldsymbol{\beta}^* \neq \mathbf{0}_K$ , then  $\mathbf{u} \in \partial\|\mathbf{x}\|_1$ ,  $\mathbf{v} = \frac{\boldsymbol{\beta}^*}{\|\boldsymbol{\beta}^*\|_2}$  and

$$\boldsymbol{\beta}^* = \mathbf{P}_K\left(\boldsymbol{\beta}_* - \rho_1\mathbf{u} - \rho_2\frac{\boldsymbol{\beta}^*}{\|\boldsymbol{\beta}^*\|_2}\right) \\ \left(1 + \frac{\rho_2}{\|\boldsymbol{\beta}^*\|_2}\right)\boldsymbol{\beta}^* = \mathbf{P}_K(\boldsymbol{\beta}_* - \rho_1\mathbf{u}).$$

Note that  $\boldsymbol{\beta}^* = \mathbf{P}_K\boldsymbol{\beta}^* \in \mathcal{R}$ . Taking the norm of both sides, we see that

$$\left(1 + \frac{\rho_2}{\|\boldsymbol{\beta}^*\|_2}\right)\|\boldsymbol{\beta}^*\|_2 = \|\mathbf{P}_K(\boldsymbol{\beta}_* - \rho_1\mathbf{u})\|_2 \\ \|\boldsymbol{\beta}^*\|_2 = \|\mathbf{P}_K(\boldsymbol{\beta}_* - \rho_1\mathbf{u})\|_2 - \rho_2 > 0.$$

Substituting this result back into the  $\beta^* \neq \mathbf{0}_K$  case, we have that

$$\beta^* = \left\{ 1 - \frac{\rho_2}{\|\mathbf{P}_K(\beta_* - \rho_1 \mathbf{u})\|_2} \right\} \mathbf{P}_K(\beta_* - \rho_1 \mathbf{u}).$$

Combining the above two cases gives the desired result.  $\square$

#### Appendix A.7. Proof of Theorem 4

**Proof.** Denote the objective function by

$$G(b) = \frac{1}{2}(b-t)^2 + \varrho\{|b| + |b+s|\}.$$

When  $s = 0$ , we obtain a lasso problem with

$$b^* = \operatorname{argmin}_{b \in \mathbb{R}} \frac{1}{2}(b-t)^2 + 2\varrho|x| = \mathcal{S}(t, 2\varrho).$$

When  $s \neq 0$ , the subdifferential of  $G(b)$  is

$$\partial G(b) = b - t + \varrho\{\partial|x| + \partial|x+s|\}.$$

We see that  $0 \in \partial G(b^*)$  if and only if there exist  $u \in \partial|b|$  and  $v \in \partial|b+s|$  with

$$b^* = b - \varrho(u + v).$$

If  $b^* = 0$ , then  $|u| < 1$  and  $v = \operatorname{sign}(s)$ , hence

$$b^* = 0 \text{ if } |t - \varrho \operatorname{sign}(s)| \leq \varrho.$$

If  $s > 0$ , then  $\operatorname{sign}(s) = 1$  and  $0 \leq t \leq 2\varrho$ . If  $s < 0$ , then  $\operatorname{sign}(s) = -1$ , and  $-2\varrho \leq t \leq 0$ . Note that if  $t \neq 0$ , then  $\operatorname{sign}(s) = \operatorname{sign}(t)$  or  $\operatorname{sign}(s)\operatorname{sign}(t) = 1$ .

When  $b^* = -s$ , then  $u = -\operatorname{sign}(s)$  and  $|v| < 1$ , hence

$$b^* = -s \text{ if } |t + s + \varrho \operatorname{sign}(s)| \leq \varrho.$$

If  $s > 0$ , then  $\operatorname{sign}(s) = 1$  and  $-(s + 2\lambda) \leq t \leq -s < 0$ . If  $s < 0$ , then  $\operatorname{sign}(s) = -1$  and  $0 < -s \leq t \leq -(s - 2\lambda)$ . Note that  $\operatorname{sign}(s) = -\operatorname{sign}(t)$  is equivalent to  $\operatorname{sign}(s)\operatorname{sign}(t) = -1$ .

Let  $C(s, t) = \frac{1 - \operatorname{sign}(s)\operatorname{sign}(t)}{2}|s| \geq 0$ . We can summarize the two cases above as

$$b^* = -C(s, t) \text{ if } 0 \leq C(s, t) \leq |t| \leq C(s, t) + 2\varrho. \quad (\text{A9})$$

If  $b^* \neq 0, -s$ , then  $u = \operatorname{sign}(b^*)$  and  $v = \operatorname{sign}(b^* + s)$ , thus

$$\begin{aligned} b^* &= t - \varrho\{\operatorname{sign}(b^*) + \operatorname{sign}(b^* + s)\} \\ b^* + s &= t + s - \varrho\{\operatorname{sign}(b^*) + \operatorname{sign}(b^* + s)\}. \end{aligned}$$

If  $\operatorname{sign}(b^*) = -\operatorname{sign}(b^* + s) = 1$ , then  $b^*(b^* + s) < 0$  or  $0 < t < -s$ . Thus  $b^* = t > 0$  if  $0 < t < -s$ . If  $\operatorname{sign}(b^*) = -\operatorname{sign}(b^* + s) = -1$ , then  $b^*(b^* + s) < 0$  or  $-s < t < 0$ . Thus  $b^* = t < 0$  if  $-s < t < 0$ . Rewriting the two cases above, we have that

$$b^* = t \quad \text{if } 0 < |t| < C(s, t). \quad (\text{A10})$$

If  $\text{sign}(b^*) = \text{sign}(b^* + s) = 1$ , then

$$\begin{aligned} \min\{b^*, b^* + s\} &> 0 \\ t - 2\varrho + \frac{s - |s|}{2} &> 0 \\ \text{sign}(t)|t| &> \text{sign}(t)\left(\frac{|s|}{2} + 2\varrho\right) - \frac{s}{2} > 0. \end{aligned}$$

Note that  $t > 0$  and  $\text{sign}(x) = \text{sign}(t)$ . If  $\text{sign}(b^*) = \text{sign}(b^* + s) = -1$ , then

$$\begin{aligned} \max\{b^*, b^* + s\} &< 0 \\ t + 2\varrho + \frac{s + |s|}{2} &> 0 \\ \text{sign}(t)|t| &< \text{sign}(t)\left(\frac{|s|}{2} + 2\varrho\right) - \frac{s}{2} < 0. \end{aligned}$$

Note that  $t < 0$  and  $\text{sign}(x) = \text{sign}(t)$ . Rewriting the two cases above, we have that

$$b^* = t - 2\varrho \text{sign}(t) \text{ if } |t| > 2\varrho + C(s, t). \quad (\text{A11})$$

Summarizing (A9)–(A11),

$$b^* = \begin{cases} t, & |t| < C(s, t), \\ -C(s, t), & C(s, t) \leq |t| \leq C(s, t) + 2\varrho, \\ \text{sign}(t)(|t| - 2\varrho), & |t| > C(s, t) + 2\varrho, \end{cases}$$

with  $C(s, t) = \frac{1 - \text{sign}(s)\text{sign}(t)}{2} |s| \geq 0$ . On one hand, when  $s \neq 0$ ,

$$b^* = t - \mathcal{S}(t, C(s, t)) + \mathcal{S}\{\mathcal{S}(t, C(s, t)), 2\varrho\}.$$

On the other hand, when  $s = 0$ , it follows that  $b^* = \mathcal{S}(t, 2\varrho)$  since  $\mathcal{S}(z, 0) = z$ .  $\square$

## References

1. Haralick, R.M.; Shanmugam, K.; Dinstein, I.H. Textural features for image classification. *IEEE Trans. Syst. Man Cybern.* **1973**, SMC-3, 610–621. [[CrossRef](#)]
2. Wang, X.; Zhang, H.H.; Wu, Y. Multiclass probability estimation with support vector machines. *J. Comput. Graph. Stat.* **2019**, *28*, 586–595. [[CrossRef](#)]
3. Hansen, J.H.; Hasan, T. Speaker recognition by machines and humans: A tutorial review. *IEEE Signal Process. Mag.* **2015**, *32*, 74–99. [[CrossRef](#)]
4. Duda, R.O.; Hart, P.E.; Stork, D.G. *Pattern Classification*; John Wiley & Sons: New York, NY, USA, 2012.
5. Hastie, T.; Tibshirani, R.; Friedman, J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*; Springer Science & Business Media: New York, NY, USA, 2009.
6. Cortes, C.; Vapnik, V. Support-vector networks. *Mach. Learn.* **1995**, *20*, 273–297. [[CrossRef](#)]
7. Cristianini, N.; Shawe-Taylor, J. *An Introduction to Support Vector Machines and Other Kernel-Based Learning Methods*; Cambridge University Press: Cambridge, UK, 2000.
8. Marron, J.S.; Todd, M.J.; Ahn, J. Distance-weighted discrimination. *J. Am. Stat. Assoc.* **2007**, *102*, 1267–1271. [[CrossRef](#)]
9. Qiao, X.; Zhang, H.H.; Liu, Y.; Todd, M.J.; Marron, J.S. Weighted distance weighted discrimination and its asymptotic properties. *J. Am. Stat. Assoc.* **2010**, *105*, 401–414. [[CrossRef](#)]
10. Marron, J. Distance-weighted discrimination. *Wiley Interdiscip. Rev. Comput. Stat.* **2015**, *7*, 109–114. [[CrossRef](#)]



11. Zhang, L.; Lin, X. Some considerations of classification for high dimension low-sample size data. *Stat. Methods Med. Res.* **2013**, *22*, 537–550. [[CrossRef](#)]
12. Wang, B.; Zou, H. Another look at distance-weighted discrimination. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **2018**, *80*, 177–198. [[CrossRef](#)]
13. Liu, Y.; Zhang, H.H.; Wu, Y. Hard or soft classification? Large-margin unified machines. *J. Am. Stat. Assoc.* **2011**, *106*, 166–177. [[CrossRef](#)]
14. Huang, H.; Liu, Y.; Du, Y.; Perou, C.M.; Hayes, D.N.; Todd, M.J.; Marron, J.S. Multiclass distance-weighted discrimination. *J. Comput. Graph. Stat.* **2013**, *22*, 953–969. [[CrossRef](#)]
15. Wang, B.; Zou, H. A multcategory kernel distance weighted discrimination method for multiclass classification. *Technometrics* **2019**, *61*, 396–408. [[CrossRef](#)]
16. Wang, B.; Zou, H. Sparse distance weighted discrimination. *J. Comput. Graph. Stat.* **2016**, *25*, 826–838. [[CrossRef](#)]
17. Wang, L.; Shen, X. On L1-norm multiclass support vector machines: Methodology and theory. *J. Am. Stat. Assoc.* **2007**, *102*, 583–594. [[CrossRef](#)]
18. Zhang, X.; Wu, Y.; Wang, L.; Li, R. Variable selection for support vector machines in moderately high dimensions. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **2016**, *78*, 53–76. [[CrossRef](#)]
19. Peng, B.; Wang, L.; Wu, Y. An error bound for L1-norm support vector machine coefficients in ultra-high dimension. *J. Mach. Learn. Res.* **2016**, *17*, 8279–8304.
20. Simon, N.; Friedman, J.; Hastie, T.; Tibshirani, R. A sparse-group lasso. *J. Comput. Graph. Stat.* **2013**, *22*, 231–245. [[CrossRef](#)]
21. Friedman, J.; Hastie, T.; Tibshirani, R. A note on the group lasso and a sparse group lasso. *arXiv* **2010**, arXiv:1001.0736.
22. Cai, T.T.; Zhang, A.; Zhou, Y. Sparse group lasso: Optimal sample complexity, convergence rate, and statistical inference. *arXiv* **2019**, arXiv:1909.09851.
23. Yu, D.; Zhang, L.; Mizera, I.; Jiang, B.; Kong, L. Sparse wavelet estimation in quantile regression with multiple functional predictors. *Comput. Stat. Data Anal.* **2019**, *136*, 12–29. [[CrossRef](#)]
24. He, Q.; Kong, L.; Wang, Y.; Wang, S.; Chan, T.A.; Holland, E. Regularized quantile regression under heterogeneous sparsity with application to quantitative genetic traits. *Comput. Stat. Data Anal.* **2016**, *95*, 222–239. [[CrossRef](#)]
25. Huang, H. Large dimensional analysis of general margin based classification methods. *arXiv* **2019**, arXiv:1901.08057.
26. Huang, H.; Yang, Q. Large scale analysis of generalization error in learning using margin based classification methods. *arXiv* **2020**, arXiv:2007.10112.
27. Lam, X.Y.; Marron, J.; Sun, D.; Toh, K.C. Fast algorithms for large-scale generalized distance weighted discrimination. *J. Comput. Graph. Stat.* **2018**, *27*, 368–379. [[CrossRef](#)]
28. Sun, D.; Toh, K.C.; Yang, L. A convergent 3-block semiproximal alternating direction method of multipliers for conic programming with 4-type constraints. *SIAM J. Optim.* **2015**, *25*, 882–915. [[CrossRef](#)]
29. Parikh, N.; Boyd, S. Proximal algorithms. *Found. Trends Optim.* **2014**, *1*, 127–239. [[CrossRef](#)]
30. Asahchop, E.L.; Branton, W.G.; Krishnan, A.; Chen, P.A.; Yang, D.; Kong, L.; Zochodne, D.W.; Brew, B.J.; Gill, M.J.; Power, C. HIV-associated sensory polyneuropathy and neuronal injury are associated with miRNA-455-3p induction. *JCI Insight* **2018**, *3*, e122450. [[CrossRef](#)]
31. Hsu, D.; Kakade, S.; Zhang, T. A tail inequality for quadratic forms of subgaussian random vectors. *Electron. Commun. Probab.* **2012**, *17*, 52. [[CrossRef](#)]

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



© 2020 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).