

---

# Supplementary Material for Theoretical Insights Into Multiclass Classification: A High-dimensional Asymptotic View

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Contemporary machine learning applications often involve classification tasks with  
2 many classes. Despite their extensive use, a precise understanding of the statistical  
3 properties and behavior of classification algorithms is still missing, especially in  
4 modern regimes where the number of classes is rather large. In this paper, we take a  
5 step in this direction by providing the first asymptotically precise analysis of linear  
6 multiclass classification. Our theoretical analysis allows us to precisely character-  
7 ize how the test error varies over different training algorithms, data distributions,  
8 problem dimensions as well as number of classes, inter/intra class correlations  
9 and class priors. Specifically, our analysis reveals that the classification accuracy  
10 is highly distribution-dependent with different algorithms achieving optimal per-  
11 formance for different data distributions and/or training/features sizes. Unlike  
12 linear regression/binary classification, the test error in multiclass classification  
13 relies on intricate functions of the trained model (e.g., correlation between some of  
14 the trained weights) whose asymptotic behavior is difficult to characterize. This  
15 challenge is already present in simple classifiers, such as those minimizing a square  
16 loss. Our novel theoretical techniques allow us to overcome some of these chal-  
17 lenges. The insights gained may pave the way for a precise understanding of other  
18 classification algorithms beyond those studied in this paper.

## 19 Roadmap to the supplementary

20 This supplementary material contains a more complete version of the paper together with all the  
21 formal proofs. In particular, compared to the main Neurips submission this supplementary contains  
22 formal statements about the Weighted Least-squares approach for the Gaussian Mixture Model (see  
23 Section 3.3) and the Multinomial Logistic Model (see Section 4.3). We also carryout additional  
24 numerical experiments in Section A. Finally, we formally prove all of our results in the appendix of  
25 this supplementary.

## 26 1 Introduction

27 Multiclass classification is fundamental to a large number of real-world machine learning applications  
28 that demand the ability to automatically distinguish between thousands of different classes. Applica-  
29 tions include essentially any problem with categorical outputs spanning natural language processing  
30 [62], where a seq2seq decoder has to choose the correct word token, reinforcement learning [28, 46],  
31 where the agent has to choose the correct action, to recommendation systems, where the model should  
32 recommend the correct movie out of many other options. For instance, YouTube’s recommendation

33 system is modeled as an extreme multiclass problem with more than a million classes where each  
34 video corresponds to a viable class [12].

35 The growing list of applications motivate an in-depth exploration of multiclass classification algo-  
36 rithms. Despite their extensive use however, a precise understanding of the statistical properties and  
37 behavior of classification algorithms is still missing with many open questions: *What is the total*  
38 *and per class test accuracy? How does this quantity depend on various problem parameters such as*  
39 *data distributions, problem dimensions, etc.? What is the highest test accuracy achievable by any*  
40 *algorithm? What is the best algorithm for each scenario? Which algorithm achieves the highest*  
41 *accuracy on rare or minority classes? How does the answer to the above question change in modern*  
42 *regimes where the number of classes are large?*

43 Asymptotic analysis in modern high-dimensional regimes where the number of training data and  
44 feature sizes grow in tandem with each other provides a promising setting for precisely quantifying the  
45 accuracy of classification algorithms as a function of problem variables and resolving the questions  
46 above. However, despite the rich literature on precise high-dimensional estimation and more recently  
47 binary classification, multiclass classification is an under-explored venue possibly due to the difficulty  
48 of capturing the intricate dependencies between the classes even for relatively simple linear classifiers.

49 **Contributions.** We initiate a precise asymptotic study of linear multiclass classification in the  
50 modern high-dimensional regime, where the sizes of the training data and of the feature vectors  
51 grow large at a proportional rate. A key promise of such a precise analysis is that it allows us  
52 to accurately compare between different classification algorithms and data models. Compared to  
53 linear regression/binary classification, we identify the following crucial challenge: the test accuracy  
54 in multiclass classification relies on intricate cross-correlations between the trained weights of  
55 the classifier. This has two consequences that drive our analysis. First, in order to obtain sharp  
56 asymptotics on the test error of any classifier, it is a prerequisite to precisely quantify the asymptotics  
57 of these cross-correlations. Second, the test error does not depend on the correlations in closed-form  
58 expressions. Thus, to compare between different classifiers, we need efficient numerical and analytic  
59 means to evaluate the test error in terms of the correlation matrices. Interestingly, we show that  
60 these challenges are already present in simple classifiers, such as minimizing the square loss, and in  
61 stylized distributional settings, such as Gaussian features. Our contributions are as follows:

- 62 • We study two different data models: a Gaussian Mixtures Model (GMM) and a Multinomial  
63 Logit Model (MLM) with Gaussian features. For each one of them, we provide a precise  
64 characterization of total and class-wise test accuracy for three different training algorithms: a  
65 least-squares based classifier, a weighted least-squares based classifier, and a simple per class  
66 averaging estimator. For the least-squares based classifiers, we develop a new technique  
67 to overcome the technical challenge of characterizing the limiting behavior of the weights'  
68 cross-correlations. For the per class averaging classifier, we show that it is Bayes optimal  
69 for a GMM with equal priors.
- 70 • We discuss efficient means of evaluating the test accuracy as a function of the weights'  
71 cross-correlations. This, together with the derived asymptotic formulae for the latter, lead  
72 to the first precise high-dimensional characterization of how the total/class-wise accuracy  
73 varies for different training algorithms, data distributions, problem dimensions as well as  
74 number of classes, the inter/intra class correlations and class priors. For special problem  
75 geometries, we derive precise conditions on the data distribution and on the relative size of  
76 the training set over which each of the two studied algorithms dominates.
- 77 • We present and discuss extensive numerical simulations that corroborate our theoretical find-  
78 ings. For instance, with an eye towards making classification algorithms more fair/equitable,  
79 we use our precise characterization of the class-wise accuracy to demonstrate how different  
80 algorithms behave in the presence of rare/minority classes. We also empirically compare the  
81 algorithms studied in this paper to other popular losses such as cross-entropy minimization.  
82 This allows us to better understand the performance of various algorithms in modern regimes  
83 of large number of classes.

84 **Related Work.** There is a classical body of algorithmic work on multiclass classification, e.g.,  
85 [13, 36, 68, 5, 16] and several empirical studies of their comparative performance [56, 20, 1, 52]. A  
86 more recent extension of this line of work investigates the effect of the loss function in deep neural  
87 networks [25, 21, 33, 4]. Algorithms for extreme multiclass problems with huge number of classes

88 has also been studied in several [10, 70, 14, 55, 34] works. On the theory front, numerous works have  
 89 investigated consistency [71, 36, 64, 54, 53] and finite-sample behavior [32, 23, 1, 40, 11, 37, 45, 38]  
 90 of multiclass classification algorithms. Our work differs from this literature in that we are interested  
 91 in *precise* characterizations of the test accuracy rather than order-wise bounds. Here we focus on  
 92 linear classifiers, but we consider the modern high-dimensional regime in which both the sample size  
 93 and the features’ dimension are large.

94 Specifically, our theoretical approach to linear multiclass classification fits in the rapidly growing  
 95 literature on *sharp* high-dimensional asymptotics of convex optimization-based estimators [18, 59,  
 96 3, 2, 60, 51, 66, 31, 19, 17, 50, 48, 67, 8, 26, 6, 29]. Most of this line of work studies linear models  
 97 and regression problems. More recently there has been a surge of interest in sharp analysis of a  
 98 variety of methods tailored to binary classification models [65, 27, 7, 61, 44, 43, 30, 57, 63, 15, 49,  
 99 41, 47, 63, 42]. Nevertheless, none of these prior works have yet considered multiclass classification  
 100 settings. Our paper unveils the salient features of the multiclass setting and shows that corresponding  
 101 results from the binary setting do not directly apply here. We emphasize that this is the case even  
 102 for seemingly simple one-vs-all (OVA) classifiers, such as minimizing the square-loss, that involve  
 103 training a single binary classifier per class [56]. The key technical tool behind our sharp analysis  
 104 is the convex Gaussian min-max Theorem (CGMT) [66, 60, 50]). However, a “naive” application  
 105 of the CGMT on the original optimization of the classifier does not allow us to compute all the  
 106 necessary correlations between the classifier’s weights to precisely capture the total/class-wise errors.  
 107 Instead, our key idea is to formulate an artificial optimization problem, which captures the missing  
 108 correlations and at the same time conveniently allows us to leverage the CGMT.

109 **Notation.** We use  $[k]$  to denote  $\{1, \dots, k\}$ . We write  $e_\ell$  for the  $\ell$ -th standard basis vector in  
 110  $\mathbb{R}^k$ . We also write  $\mathbf{I}_k, \mathbf{0}_{k \times k}$  and  $\mathbf{1}_k$  for the  $k \times k$  identity and all-zeros matrices and the  $k \times 1$   
 111 all-ones vectors. For a vector  $\mathbf{c} \in \mathbb{R}^k$  we write  $\arg \max \mathbf{c}$  to denote the index of its largest entry,  
 112 i.e.,  $\arg \max \mathbf{c} = \arg \max_{j \in [k]} c_j$ . We use  $Q(x)$  for the tail of a standard Gaussian (Q-function).

113 Finally, we reserve variables  $G_0, G_1, \dots, G_k \stackrel{iid}{\sim} \mathcal{N}(0, 1)$  to denote i.i.d. standard Gaussians.

## 114 2 Problem formulation

115 In this paper we focus on multiclass classification problems with  $k$  classes. Specifically, we assume  
 116 the training data consists of  $n$  feature/label pairs  $\{(\mathbf{x}_i, Y_i)\}_{i=1}^n$  with  $\mathbf{x}_i \in \mathbb{R}^d$  representing the features  
 117 and  $Y_i \in \{1, 2, \dots, k\}$  the associated labels representing one of  $k$  classes. It will be convenient  
 118 to also model the labels as one-hot encoded vectors  $\mathbf{y}_i \in \mathbb{R}^k$  representing one of  $k$  classes with  
 119 one-hot encoding i.e.  $\mathbf{y}_i = e_{Y_i}$ . Therefore, when convenient we shall use  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  to represent  
 120 the training data. Throughout this paper we shall use  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  to denote  
 121 the matrix of features with the corresponding labels aggregated into a matrix of the form  $\mathbf{Y} =$   
 122  $[\mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_n] \in \mathbb{R}^{k \times n}$ . We shall also use  $\mathbf{Y}_\ell \in \mathbb{R}^n$  to denote the  $\ell$ -th row of  $\mathbf{Y}$ . In our analysis  
 123 we focus on training linear classifiers. We shall use  $\mathbf{W} \in \mathbb{R}^{k \times d}$  and  $\mathbf{b} \in \mathbb{R}^k$  to denote the weights and  
 124 biases of this linear model. The overall input-output relationship of the classifier in this case is a  
 125 function that maps an input vector  $\mathbf{x} \in \mathbb{R}^d$  into an output of size  $k$  via  $\mathbf{x} \mapsto \mathbf{W}\mathbf{x} + \mathbf{b} \in \mathbb{R}^k$  where a  
 126 training algorithm is used to train the corresponding weights  $\mathbf{W} \in \mathbb{R}^{k \times d}$  and biases  $\mathbf{b} \in \mathbb{R}^k$ . Next we  
 127 detail the data models and training algorithms that are formally studied in this paper. We end this  
 128 section by discussing how the test error can be calculated for the different data models.

### 129 2.1 Data Models

130 In our theoretical analysis we assume the training data  $\{(\mathbf{x}_i, Y_i)\}_{i=1}^n$  (alternatively  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ ) are  
 131 generated i.i.d. according to  $(\mathbf{x}, Y)/(\mathbf{x}, \mathbf{y})$ . We consider two models for the distribution of  $(\mathbf{x}, \mathbf{y})$   
 132 which we detail next. In both models we shall use mean/regressor vectors  $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k \in \mathbb{R}^d$  and aggregate  
 133 them into columns of a matrix of the form  $\mathbf{M} := [\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k] \in \mathbb{R}^{d \times k}$ . In the first model, these  
 134 vectors represent the mean of the features conditioned on the class, i.e.  $\boldsymbol{\mu}_\ell = \mathbb{E}[\mathbf{x}|Y = \ell]$ , where as in  
 135 the second model these vectors can be viewed as regressor coefficients. We shall refer to  $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k / \mathbf{M}$   
 136 as “mean” vectors/matrix in both models. We use  $\Sigma_{\boldsymbol{\mu}, \boldsymbol{\mu}} = \mathbf{M}^T \mathbf{M}$  to denote the Gramian matrix of  
 137 means. Furthermore, we shall use  $\mu_\ell := \|\boldsymbol{\mu}_\ell\|_{\ell_2}$  to denote the norm of the mean vector  $\boldsymbol{\mu}_\ell$ .

138 **Gaussian Mixture Model (GMM).** In this model each example  $(\mathbf{x}, Y)$  belongs to class  $\ell \in [k]$   
 139 with probability  $\pi_\ell$ , i.e.,  $\mathbb{P}\{Y = \ell\} = \pi_\ell$ . We let  $\boldsymbol{\pi} = [\pi_1 \ \pi_2 \ \dots \ \pi_k]^T \in \mathbb{R}^k$  denote the vector of  
 140 priors which of course obeys  $\boldsymbol{\pi} \geq \mathbf{0}$  and  $\mathbf{1}^T \boldsymbol{\pi} = 1$ . Also, we model the class conditional density of an  
 141 example in class  $\ell$  with an isotropic Gaussian centered at a mean vector  $\boldsymbol{\mu}_\ell$ . In particular, we say that  
 142 a data point  $(\mathbf{x}, Y)$  (or its one-hot encoded representation  $(\mathbf{x}, \mathbf{y})$ ) follows the GMM model when

$$\mathbb{P}\{Y = \ell\} = \pi_\ell \quad \text{and} \quad \mathbf{x} = \boldsymbol{\mu}_\ell + \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d). \quad (2.1)$$

143 We note that for a training set summarized by the feature and label matrices  $\mathbf{X}$  and  $\mathbf{Y}$  with columns  
 144 generated i.i.d. according to the above distribution we have:  $\mathbf{X} = \mathbf{M}\mathbf{Y} + \mathbf{Z}$  where  $\mathbf{Z} \in \mathbb{R}^{d \times n}$  is a  
 145 Gaussian noise matrix with i.i.d.  $\mathcal{N}(0, \sigma^2)$  entries.

146 **Multinomial Logit Model (MLM).** In this model we assume that feature vectors  $\mathbf{x}$  are distributed  
 147 i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  and that the conditional density of the class labels is given by the soft-max function.  
 148 Concretely, we say that a data point  $(\mathbf{x}, Y)$  (or its one-hot encoded representation  $(\mathbf{x}, \mathbf{y})$ ) follows  
 149 the multinomial logit model when

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d) \quad \text{and} \quad \mathbb{P}\{Y = \ell \mid \mathbf{x}\} = \frac{e^{\langle \boldsymbol{\mu}_\ell, \mathbf{x} \rangle}}{\sum_{j \in [k]} e^{\langle \boldsymbol{\mu}_j, \mathbf{x} \rangle}}. \quad (2.2)$$

## 150 2.2 Classification algorithms

151 As mentioned earlier in this paper we focus on training linear classifiers of the form  $\mathbf{x} \mapsto \mathbf{W}\mathbf{x} + \mathbf{b}$   
 152 with  $\mathbf{W} \in \mathbb{R}^{k \times d}$  denoting the weights and  $\mathbf{b} \in \mathbb{R}^k$  the offset values.

153 **Least-squares (LS).** In this approach ‘

154 **Class averaging (Avg).** This approach uses the following weight and offset values

$$\widehat{\mathbf{W}} := \frac{1}{n} \mathbf{Y} \mathbf{X}^T \quad \text{and} \quad \widehat{\mathbf{b}} := \frac{1}{n} \mathbf{Y} \mathbf{1}. \quad (2.3)$$

155 Let  $n_\ell$  be the number of training data from class  $\ell$  then, equivalently,

$$\widehat{\mathbf{w}}_\ell = \frac{n_\ell}{n} \left( \frac{1}{n_\ell} \sum_{i: Y_i = \ell} \mathbf{x}_i \right) \quad \text{and} \quad \widehat{\mathbf{b}}_\ell = \frac{n_\ell}{n}.$$

156 Therefore, this classifier picks weights according to the empirical mean of features of each class  
 157 multiplied by the relative frequency of that class and the offset value as the fraction of data points  
 158 from that class. We note that this algorithm has the same classification performance as the outcome  
 159 of the ridge-regularized least-squares with infinite regularization.

160 **Weighted Least-squares (WLS).** This is a variation of the Least-squares approach where we fit a  
 161 weighted least squares loss of the form

$$(\widehat{\mathbf{W}}, \widehat{\mathbf{b}}) := \arg \min_{\mathbf{W}, \mathbf{b}} \frac{1}{2n} \left\| (\mathbf{W} \mathbf{X} + \mathbf{b} \mathbf{1}_n^T - \mathbf{Y}) \mathbf{D} \right\|_F^2. \quad (2.4)$$

162 Here,  $\mathbf{D} \in \mathbb{R}^{n \times n}$  is a diagonal matrix with the  $i$ th diagonal entry equal to  $D_{ii} = \omega_\ell$  when the  $i$ -th data  
 163 point is from class  $\ell$  (i.e.  $Y_i = \ell$ ) and  $\omega_\ell \geq 0$ ,  $\ell \in [k]$  denote the weights. Aggregating the weights  
 164 into a vector of the form  $\boldsymbol{\omega} = [\omega_1 \ \omega_2 \ \dots \ \omega_k]^T \in \mathbb{R}^k$  we can rewrite  $\mathbf{D}$  in the form

$$\mathbf{D} = \text{diag}(\mathbf{Y}^T \boldsymbol{\omega}).$$

165 In this approach the loss associated to data points to class  $\ell$  are weighted by a factor  $\omega_\ell^2$ . For instance,  
 166 if the class priors are known, a natural choice would be  $\omega_\ell = 1/\sqrt{\pi_\ell}$ . Such a weighted approach  
 167 allows the classification algorithm to focus on rare/minority classes which are not well represented in  
 168 the training data.

169 **Cross-entropy (CE).** In this approach the best weight/offset values are determined by fitting a cross  
 170 entropy loss

$$(\widehat{\mathbf{W}}, \widehat{\mathbf{b}}) := \arg \min_{\mathbf{W}, \mathbf{b}} \frac{1}{n} \sum_{i=1}^n \log \left( \frac{\sum_{\ell=1}^k e^{\langle \widehat{\mathbf{w}}_\ell, \mathbf{x}_i \rangle + \widehat{\mathbf{b}}_\ell}}{e^{\langle \widehat{\mathbf{w}}_{Y_i}, \mathbf{x}_i \rangle + \widehat{\mathbf{b}}_{Y_i}}} \right).$$

171 Theoretical analysis for CE is substantially more involved and we defer it to future work. Nevertheless,  
 172 we compare with this classifier in our numerical simulations.

173 **2.3 Class-wise and total test classification error in the high-dimensional regime**

174 Let  $\widehat{\mathbf{W}}, \widehat{\mathbf{b}}$  denote the parameters of a trained classifier. Now consider a fresh data sample  $(\mathbf{x}, Y)$   
 175 generated according to the same distribution as the training data. The class predicted by the classifier  
 176 in this case is given by

$$\widehat{Y} = \arg \max_{j \in [k]} \langle \widehat{\mathbf{w}}_j, \mathbf{x} \rangle + \widehat{\mathbf{b}}_j.$$

177 Therefore, the classification error condition on the the true label being  $c$  (which we shall refer to as  
 178 the class-wise test error) is equal to

$$\mathbb{P}_{e|c} := \mathbb{P} \{ \widehat{Y} \neq Y | Y = c \} = \mathbb{P} \{ \langle \widehat{\mathbf{w}}_c, \mathbf{x} \rangle + \widehat{\mathbf{b}}_c \leq \max_{j \neq c} \langle \widehat{\mathbf{w}}_j, \mathbf{x} \rangle + \widehat{\mathbf{b}}_j \}. \quad (2.5)$$

179 The total classification error is given by

$$\mathbb{P}_e := \mathbb{P} \{ \widehat{Y} \neq Y \} = \mathbb{P} \left\{ \arg \max_{j \in [k]} \{ \langle \widehat{\mathbf{w}}_j, \mathbf{x} \rangle + \widehat{\mathbf{b}}_j \} \neq Y \right\} = \mathbb{P} \left\{ \langle \widehat{\mathbf{w}}_Y, \mathbf{x} \rangle + \widehat{\mathbf{b}}_Y \leq \max_{j \neq Y} \langle \widehat{\mathbf{w}}_j, \mathbf{x} \rangle + \widehat{\mathbf{b}}_j \right\}. \quad (2.6)$$

For both the GMM and MLM, the classification error depends on the vector of intercepts  $\widehat{\mathbf{b}} \in \mathbb{R}^k$  and the following key ‘‘correlation’’ matrices:

$$\Sigma_{\mathbf{w}, \mathbf{w}} := \widehat{\mathbf{W}} \widehat{\mathbf{W}}^T \quad \text{and} \quad \Sigma_{\mathbf{w}, \boldsymbol{\mu}} := \widehat{\mathbf{W}} \mathbf{M}.$$

180

$$\text{GMM:} \quad \mathbb{P}_e = \mathbb{P} \left\{ \arg \max (\sigma \mathbf{g} + \widehat{\mathbf{b}} + \Sigma_{\mathbf{w}, \boldsymbol{\mu}} \mathbf{e}_Y) \neq Y \right\}, \quad \text{where } \mathbf{g} \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}, \mathbf{w}}). \quad (2.7)$$

$$\text{MLM:} \quad \mathbb{P}_e = \mathbb{P} \left\{ \arg \max (\mathbf{g} + \widehat{\mathbf{b}}) \neq Y(\mathbf{h}) \right\}, \quad \text{where } \begin{bmatrix} \mathbf{g} \\ \mathbf{h} \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \Sigma_{\mathbf{w}, \mathbf{w}} & \Sigma_{\mathbf{w}, \boldsymbol{\mu}} \\ \Sigma_{\mathbf{w}, \boldsymbol{\mu}}^T & \Sigma_{\boldsymbol{\mu}, \boldsymbol{\mu}} \end{bmatrix} \right), \quad (2.8)$$

$$\text{and } \mathbb{P} \{ Y(\mathbf{h}) = \ell \} = e^{\mathbf{h}_\ell} / \sum_{j \in [k]} e^{\mathbf{h}_j}, \quad \ell \in [k].$$

181 **Calculating the class-wise/total misclassification errors.** The identities (2.7) and (2.8) as well as  
 182 similar ones for the class-wise test error demonstrate that the total/class-wise errors only depend on  
 183 the correlation matrices  $\Sigma_{\mathbf{w}, \mathbf{w}}$  and  $\Sigma_{\mathbf{w}, \boldsymbol{\mu}}$ , the offset values  $\widehat{\mathbf{b}}$  and the the class conditional means. For  
 184 instance, as we show in the supplementary for GMM the class-wise errors are given by

$$\mathbb{P}_{e|c} = 1 - \mathbb{P} \{ \mathbf{S}_c^{1/2} \mathbf{z} \geq \mathbf{t}_c \}, \quad (2.9)$$

185 where  $\mathbf{z}$  is a Gaussian random vector distributed as  $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{k-1})$ ,  $\mathbf{S}_c \in \mathbb{R}^{(k-1) \times (k-1)}$  is a symmetric  
 186 matrix such that its  $i, j$  element is given by  $[\mathbf{S}_c]_{ij} := \langle \widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_j, \widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_i \rangle$  and  $\mathbf{t}_c \in \mathbb{R}^{k-1}$  a vector  
 187 with entries  $[\mathbf{t}_c]_i := \langle \widehat{\mathbf{w}}_i - \widehat{\mathbf{w}}_c, \boldsymbol{\mu}_c \rangle + (\widehat{\mathbf{b}}_i - \widehat{\mathbf{b}}_c)$ . Similarly, based on (2.9) the total classification  
 188 error in GMM is equal to  $\mathbb{P}_e = \sum_{\ell=1}^k \pi_\ell \mathbb{P}_{e|c} = 1 - \sum_{\ell=1}^k \pi_\ell \mathbb{P} \{ \mathbf{S}_c^{1/2} \mathbf{z} \geq \mathbf{t}_c \}$ . As also detailed in  
 189 the supplementary, the class-wise/total test errors for MLM similarly depends on quantities of the  
 190 form  $\mathbb{P} \{ \mathbf{A} \mathbf{z} \geq \mathbf{t} \}$  with  $\mathbf{z}$  a standard Gaussian random vector,  $\mathbf{A}$  and  $\mathbf{t}$  depending only on correlation  
 191 matrices, conditional means and classifier offset-values. There are a variety of algorithmic approaches  
 192 to calculate  $\mathbb{P} \{ \mathbf{A} \mathbf{z} \geq \mathbf{t} \}$  once  $\mathbf{A}/\mathbf{t}$  is known based on MCMC methods. Analytic bounds on this  
 193 quantity have also been studied in the literature [24, 58]; see more details in the supplementary.

194 **High-dimensional regime.** The main technical contribution of this paper involves sharp asymptotic  
 195 formulae for the class-wise and total classification error of averaging and (weighted) LS algorithms  
 196 for GMM and MLM. Due to space constraints, we defer the statements of our main results for WLS,  
 197 as well as, all our proofs to the supplementary. All our results hold in the following high-dimensional  
 198 regime with finite  $k$ .

199 **Assumption 1** *We focus on a double asymptotic regime where  $n, d \rightarrow \infty$  at a fixed ratio  $\gamma = d/n > 0$ .*

200 For the (weighted) least-squares classifier, we focus here in the overdetermined regime  $\gamma < 1$ .  
 201 However, our approach is also directly applicable to regularized (or min-norm) LS/WLS in the  
 202 overparameterized regime  $\gamma > 1$ . For a sequence of random variables  $\mathcal{X}_{n,d}$  that converges in  
 203 probability to some constant  $c$  in the limit above, we simply write  $X_{n,d} \xrightarrow{P} c$ .

### 204 3 Results for Gaussian Mixture Model

205 In this section we discuss the asymptotics of the intercepts/correlation matrices for the averaging and  
 206 the LS classifiers for the Gaussian Mixture Model. We end this section by also characterizing the  
 207 Bayes optimal estimator in this model when the priors are balanced  $\pi_\ell = 1/k$  for  $\ell \in [k]$ .

#### 208 3.1 Class averaging classifier

209 **Proposition 3.1** Consider data generated according to GMM in an asymptotic regime with any  $\gamma > 0$ .  
 210 For the averaging estimator discussed in Section 2.2, the following high-dimensional limits hold

$$\widehat{\mathbf{b}} \xrightarrow{P} \boldsymbol{\pi}, \quad \boldsymbol{\Sigma}_{\mathbf{w}, \boldsymbol{\mu}} \xrightarrow{P} \text{diag}(\boldsymbol{\pi}) \cdot \boldsymbol{\Sigma}_{\boldsymbol{\mu}, \boldsymbol{\mu}}, \quad (3.1a)$$

$$\boldsymbol{\Sigma}_{\mathbf{w}, \mathbf{w}} \xrightarrow{P} \gamma \sigma^2 \cdot \text{diag}(\boldsymbol{\pi}) + \text{diag}(\boldsymbol{\pi}) \cdot \boldsymbol{\Sigma}_{\boldsymbol{\mu}, \boldsymbol{\mu}} \cdot \text{diag}(\boldsymbol{\pi}). \quad (3.1b)$$

211 The above result allows us to precisely characterize the behavior of the averaging estimator in the  
 212 high-dimensional regime. Let us consider a few special cases.

213 **Two classes.** Consider the special case with two classes with class priors  $\pi_1 = 1 - \pi_2 =: \pi$ . In this case  
 214 we can compute the class-wise misclassification probabilities  $\mathbb{P}_{e|1}$  and  $\mathbb{P}_{e|2}$  explicitly. Specifically  
 215 using (3.1), we have  $\mathcal{S}_1 = \|\pi \boldsymbol{\mu}_1 - (1 - \pi) \boldsymbol{\mu}_2\|_{\ell_2}^2 + \gamma \sigma^2$  and  $t_1 = (1 - 2\pi) + (1 - \pi) \langle \boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \rangle - \pi \|\boldsymbol{\mu}_1\|_{\ell_2}^2$ .

216 Substituting the latter two in (2.9) we arrive at  $\mathbb{P}_{e|1} \xrightarrow{P} Q\left(\frac{\pi \|\boldsymbol{\mu}_1\|_{\ell_2}^2 - (1 - \pi) \langle \boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \rangle + 2\pi - 1}{\sqrt{\|\pi \boldsymbol{\mu}_1 - (1 - \pi) \boldsymbol{\mu}_2\|_{\ell_2}^2 + \gamma \sigma^2}}\right)$ . In the case

217 of equal priors  $\pi = \pi_1 = \pi_2 = 1/2$ , antipodal and equal energy of the means, i.e.  $\boldsymbol{\mu}_1 = -\boldsymbol{\mu}_2$  and  
 218  $\mu := \|\boldsymbol{\mu}_1\|_{\ell_2} = \|\boldsymbol{\mu}_2\|_{\ell_2}$ , we can use the above to conclude that  $\mathbb{P}_{e|1} = \mathbb{P}_{e|2} = \frac{1}{2} \mathbb{P}_e = \frac{1}{2} Q\left(\sqrt{\frac{\mu^2}{\mu^2 + \gamma \sigma^2}}\right)$ .

219 This formula recovers the result of [47] for this special case. Also, as mentioned in [47], the formula  
 220 matches the Bayes optimal error computed in [39] for Gaussian mean vectors. This shows that the  
 221 class averaging method is Bayes optimal in this very simple setting. In Section 3.4, we generalize  
 222 this result to multiple classes: we show that the average estimator is (asymptotically) Bayes optimal  
 223 for balanced classes and equal-energy Gaussian means for any  $k \geq 2$ .

224 **Orthogonal means, equal priors and equal energy.** Next we focus on a special case with orthogo-  
 225 nal means  $\langle \boldsymbol{\mu}_i, \boldsymbol{\mu}_j \rangle = 0$ ,  $i \neq j \in [k]$  of equal energy  $\mu^2 := \|\boldsymbol{\mu}_i\|_{\ell_2}^2$  and of equal priors  $\pi_i = \pi = 1/k$  for  
 226  $i \in [k]$ . In this case, the class-wise miss-classification error converges to  $\mathbb{P}_{e|c} \xrightarrow{P} 1 - \mathbb{P}\{\mathcal{S}_c^{1/2} \mathbf{z} > \mathbf{t}\}$ ,  
 227 where  $\mathcal{S}_c = \pi(\pi \mu^2 + \gamma \sigma^2)(\mathbf{I}_{k-1} + \mathbf{1}_{k-1} \mathbf{1}_{k-1}^T)$  and  $\mathbf{t} = -\pi \mu^2 \mathbf{1}_{k-1}$ . Defining  $u_{\text{Avg}} := \frac{\mu^2}{\sigma} \sqrt{\frac{1}{\mu^2 + k \gamma \sigma^2}}$   
 228 after some algebraic manipulations the total classification error of the averaging estimator in this case  
 229 is given by  $\mathbb{P}_{c|e} = \mathbb{P}_{e, \text{Avg}} \xrightarrow{P} \mathbb{P}\{G_0 + \max_{j \in [k-1]} G_j \geq u_{\text{Avg}}\}$ .

#### 230 3.2 Least-squares classifier

231 This section focuses on characterizing the intercepts and correlation matrices for the least-squares  
 232 classifier. To present our results, we assume that the Gramian matrix has eigenvalue decomposition

$$\boldsymbol{\Sigma}_{\boldsymbol{\mu}, \boldsymbol{\mu}} = \mathbf{M}^T \mathbf{M} = \mathbf{V} \boldsymbol{\Sigma}^2 \mathbf{V}^T, \quad \boldsymbol{\Sigma} > 0 \in \mathbb{R}^{r \times r}, \quad \mathbf{V} \in \mathbb{R}^{k \times r}, \quad r \leq k. \quad (3.2)$$

233 with  $\boldsymbol{\Sigma}$  a diagonal positive-definite matrix and  $\mathbf{V}$  an orthonormal matrix obeying  $\mathbf{V}^T \mathbf{V} = \mathbf{I}_r$ .

234 **Theorem 3.2** Consider data generated according to GMM in an asymptotic regime with  $\gamma < 1$ . In  
 235 addition to (3.2), define the following two positive (semi)-definite matrices:  $\mathbf{P} := \text{diag}(\boldsymbol{\pi}) - \boldsymbol{\pi} \boldsymbol{\pi}^T \geq$   
 236  $\mathbf{0}_{k \times k}$  and  $\boldsymbol{\Delta} := \sigma^2 \mathbf{I}_r + \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{P} \mathbf{V} \boldsymbol{\Sigma} > \mathbf{0}_{r \times r}$ . Then, for the least-squares linear classifier  $(\widehat{\mathbf{W}}, \widehat{\mathbf{b}})$  the  
 237 following limits are true asymptotically

$$\widehat{\mathbf{b}} \xrightarrow{P} \boldsymbol{\pi} - \mathbf{P} \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T \boldsymbol{\pi}, \quad \boldsymbol{\Sigma}_{\mathbf{w}, \boldsymbol{\mu}} \xrightarrow{P} \mathbf{P} \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T, \quad (3.3a)$$

$$\boldsymbol{\Sigma}_{\mathbf{w}, \mathbf{w}} \xrightarrow{P} \frac{\gamma}{(1 - \gamma) \sigma^2} \mathbf{P} + \mathbf{P} \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\gamma}{(1 - \gamma) \sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{P}. \quad (3.3b)$$

238 The above result allows us to precisely characterize the behavior of the least-squares classifier in  
 239 the high-dimensional regime. In the supplementary, we specialize (3.3a) to the case of orthogonal

240 means. Compared to the weight vectors  $\widehat{\mathbf{w}}_i, i \in [k]$  of the class averaging classifier that are also  
 241 (asymptotically) orthogonal when means are orthogonal, this is *not* the case for LS. We show next  
 242 that these spurious correlations only hurt the classification error when classes are balanced.

243 **Proposition 3.3** Consider the case of orthogonal, equal energy-means  $\Sigma_{\mu, \mu} = \mu \mathbf{I}_k$ , balanced priors  
 244  $\pi_i = 1/k, i \in [k]$  and  $\gamma < 1$ . Setting  $u_{\text{LS}} := \frac{\mu^2}{\sigma} \sqrt{\frac{1-\gamma}{\mu^2+k\gamma\sigma^2}}$ , it holds that  $\mathbb{P}_{e, \text{LS}} \xrightarrow{P} \mathbb{P}\{G_0 +$   
 245  $\max_{j \in [k-1]} G_j \geq u_{\text{LS}}\}$ . Specifically, since  $u_{\text{LS}} = u_{\text{Avg}} \sqrt{1-\gamma} < u_{\text{Avg}}$ , the averaging estimator  
 246 strictly outperforms LS for all  $0 < \gamma < 1$  and  $k \geq 2$  in this setting.

### 247 3.3 Weighted least-squares classifier

248 We now focus on characterizing the intercepts/correlation matrices for the WLS classifier.

249 **Theorem 3.4** Consider data generated according to GMM and  $\gamma < 1$ . Consider a weighted LS  
 250 classifier with weights  $\mathbf{D} = \text{diag}(\omega_1, \dots, \omega_k)$  and let  $\eta$  be the unique solution to  $\sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{\omega_\ell^2 + \eta} = \gamma$ . Also  
 251 define  $\mathbf{P} := \text{diag}(\tilde{\pi}) - \tilde{\pi} \tilde{\pi}^T \geq \mathbf{0}_{k \times k}$  and  $\Delta := \sigma^2 \mathbf{I}_r + \Sigma \mathbf{V}^T \mathbf{P} \mathbf{V} \Sigma > \mathbf{0}_{r \times r}$  with the entries of  $\tilde{\pi}$  given  
 252 by  $\tilde{\pi}_\ell = \frac{1}{\gamma} \frac{\pi_\ell \omega_\ell^2}{\omega_\ell^2 + \eta}$ . Then, for the WLS linear classifier  $(\widehat{\mathbf{W}}, \widehat{\mathbf{b}})$  the following asymptotic limits hold

$$\widehat{\mathbf{b}} \xrightarrow{P} \tilde{\pi} - \mathbf{P} \mathbf{V} \Sigma \Delta^{-1} \Sigma \mathbf{V}^T \tilde{\pi}, \quad \Sigma_{\mathbf{w}, \mu} \xrightarrow{P} \mathbf{P} \mathbf{V} \Sigma \Delta^{-1} \Sigma \mathbf{V}^T, \quad (3.4a)$$

$$\Sigma_{\mathbf{w}, \mathbf{w}} \xrightarrow{P} \frac{\zeta}{\sigma^2} \mathbf{P} + \mathbf{P} \mathbf{V} \Sigma \Delta^{-1} \left( \Delta^{-1} - \frac{\zeta}{\sigma^2} \mathbf{I}_r \right) \Sigma \mathbf{V}^T \mathbf{P} + \frac{\eta \zeta}{\sigma^2} \mathbf{Q}. \quad (3.4b)$$

253 Here,  $\zeta := \frac{\gamma}{\eta \left( \sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{(\omega_\ell^2 + \eta)^2} \right)}$  and  $\mathbf{Q} \in \mathbb{R}^{k \times k}$  is a known matrix depending on various problem parameters

254 (precise value discussed in the proof).

255 Surprisingly, the effect of the weights is essentially equivalent to adjusting the class priors from  $\pi$   
 256 to  $\tilde{\pi}$  defined in the theorem (modulo the extra additive term in the cross correlation matrix  $\Sigma_{\mathbf{w}, \mathbf{w}}$ ).  
 257 This shows that weighted LS has similar performance to an un-weighted LS applied to a model  
 258 with different class priors  $\tilde{\pi}$ . This characterization allows us to precisely understand how different  
 259 weighting schemes can alter test accuracy for rare/minority classes.

### 260 3.4 Bayes estimator for the balanced Gaussian Mixture Model

261 To check how far the above algorithms are from the lowest misclassification error achievable by  
 262 any algorithm in this section, we consider a Bayesian setting with Gaussian mean vectors and  
 263 we derive the Bayes-optimal risk for the case of equal priors. Recall that the Bayes estimator  
 264  $\hat{Y} = \arg \max_{\ell \in [k]} \mathbb{P}\{Y = \ell \mid \mathbf{X}, \mathbf{Y}, \mathbf{x}\}$  minimizes the risk  $\mathbb{P}_e = \mathbb{P}\{\hat{Y} \neq Y\} = \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \mathbf{x}, Y} [\mathbb{1}[\hat{Y} \neq Y]]$ .

265 **Proposition 3.5** Consider  $\mu_i \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \frac{\mu^2}{d} \mathbf{I}_d)$  and  $\pi_i = 1/k$  for all  $i \in [k]$ . Set  $u_{\text{Bayes}} := \frac{\mu^2}{\sigma} \frac{1}{\sqrt{\mu^2+k\gamma\sigma^2}}$ .

266 Then, the Bayes risk converges to  $\mathbb{P}\{G_0 + \max_{\ell \in [k-1]} G_\ell \geq u_{\text{Bayes}}\}$ .

267 Under the Gaussian prior, the means are asymptotically orthogonal and equal-energy. As shown  
 268 earlier, in this setting,  $\mathbb{P}_{e, \text{Avg}} \xrightarrow{P} \mathbb{P}\{G_0 + \max_{\ell \in [k-1]} G_\ell \geq u_{\text{Avg}}\}$ . But,  $u_{\text{Avg}} = u_{\text{Bayes}}$ . Thus, the  
 269 averaging method is (asymptotically) Bayes optimal for equal-norm, orthogonal means and balanced  
 270 classes. An analogous result was recently derived in [39, 47], but only for binary classification.

## 271 4 Results for Multinomial Logit Model

272 In this section we discuss the asymptotics of the intercepts/correlation matrices for MLM. We present  
 273 results for arbitrary mean-vectors as well as special cases where the means are mutually orthogonal.  
 274 Recall the eigenvalue decomposition of the Grammian  $\Sigma_{\mu, \mu} = \mathbf{V} \Sigma \mathbf{V}^T$  in (3.2). In order to state our  
 275 results, it is convenient to introduce the following probability vectors in  $\mathbb{R}^k$  and  $\mathbb{R}^{k^2}$ :

$$\boldsymbol{\pi} := \mathbb{E} \left[ \frac{e^{\mathbf{V} \Sigma \mathbf{g}}}{\mathbf{1}_k^T e^{\mathbf{V} \Sigma \mathbf{g}}} \right] \in \mathbb{R}^k \quad \text{and} \quad \boldsymbol{\Pi} := \mathbb{E} \left[ \frac{(e^{\mathbf{V} \Sigma \mathbf{g}})(e^{\mathbf{V} \Sigma \mathbf{g}})^T}{(\mathbf{1}_k^T e^{\mathbf{V} \Sigma \mathbf{g}})^2} \right] \in \mathbb{R}^{k \times k}, \quad \text{where } \mathbf{g} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r). \quad (4.1)$$

276 Note that  $\boldsymbol{\pi}$  and  $\boldsymbol{\Pi}$  are the first and second moments of the soft-max mapping of  $\mathbf{V}\boldsymbol{\Sigma}\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}})$ .  
 277 In fact, for the MLM in (2.2) it holds that  $\mathbb{P}\{Y = \ell\} = \mathbb{E}[\mathbb{P}\{Y = \ell | \mathbf{x}\}] = \mathbb{E}\left[\frac{e^{e_\ell^T \mathbf{V}\boldsymbol{\Sigma}\mathbf{g}}}{\sum_{j \in [k]} e^{e_j^T \mathbf{V}\boldsymbol{\Sigma}\mathbf{g}}}\right] = \boldsymbol{\pi}_\ell$ ,  $\ell \in [k]$   
 278 since  $\mathbf{M}^T \mathbf{x}$  is distributed as  $\mathbf{V}\boldsymbol{\Sigma}\mathbf{g}$ . Thus,  $\boldsymbol{\pi}$  is the vector of class priors (which explains the slight  
 279 abuse of notation here in relation to our notation for the class priors of the GMM).

#### 280 4.1 Class averaging classifier

281 **Proposition 4.1** Consider data generated according to MLM in an asymptotic regime with any  $\gamma > 0$ .  
 282 For the averaging classifier, the following high-dimensional limits hold

$$\widehat{\mathbf{b}} \xrightarrow{P} \boldsymbol{\pi}, \quad \boldsymbol{\Sigma}_{\mathbf{w},\boldsymbol{\mu}} \xrightarrow{P} (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \cdot \boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}}, \quad (4.2a)$$

$$\boldsymbol{\Sigma}_{\mathbf{w},\mathbf{w}} \xrightarrow{P} \gamma \cdot \text{diag}(\boldsymbol{\pi}) + (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \cdot \boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}} (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}). \quad (4.2b)$$

283 Using Gaussian decomposition in (2.8) and checking from (4.2) that  $\boldsymbol{\Sigma}_{\mathbf{w},\mathbf{w}} - \boldsymbol{\Sigma}_{\mathbf{w},\boldsymbol{\mu}} \boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}}^\dagger \boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}}^T \xrightarrow{P}$   
 284  $\gamma \cdot \text{diag}(\boldsymbol{\pi})$  the test error obtains the following explicit form:

$$\mathbb{P}_{e,\text{Avg}} \xrightarrow{P} \mathbb{P}\left\{\arg \max \left\{\sqrt{\gamma} \cdot \text{diag}(\sqrt{\boldsymbol{\pi}}) \cdot \tilde{\mathbf{g}} + (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \cdot \mathbf{V}\boldsymbol{\Sigma} \cdot \mathbf{g} + \boldsymbol{\pi}\right\} \neq Y(\mathbf{g})\right\}, \quad (4.3)$$

285 where  $\tilde{\mathbf{g}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$ ,  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_r)$  and  $\mathbb{P}\{Y(\mathbf{g}) = c\} = e^{e_c^T \mathbf{V}\boldsymbol{\Sigma}\mathbf{g}} / \sum_{j \in [k]} e^{e_j^T \mathbf{V}\boldsymbol{\Sigma}\mathbf{g}}$ ,  $c \in [k]$ .

#### 286 4.2 Least-squares classifier

287 This section focuses on characterizing the intercepts and correlation matrices for the least-squares  
 288 classifier. We also use the result to characterize conditions under which LS outperforms averaging.

289 **Theorem 4.2** Consider data generated according to MLM in an asymptotic regime with  $0 < \gamma < 1$ .  
 290 Recall the notation in (4.1). For the LS classifier, the following high-dimensional limits hold.

$$\widehat{\mathbf{b}} \xrightarrow{P} \boldsymbol{\pi}, \quad \boldsymbol{\Sigma}_{\mathbf{w},\boldsymbol{\mu}} \xrightarrow{P} (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \cdot \boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}}, \quad (4.4a)$$

$$\boldsymbol{\Sigma}_{\mathbf{w},\mathbf{w}} \xrightarrow{P} \frac{\gamma}{1-\gamma} \cdot (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\pi}\boldsymbol{\pi}^T) + \frac{1-2\gamma}{1-\gamma} \cdot (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \cdot \boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}} \cdot (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}). \quad (4.4b)$$

It is interesting to observe that (4.4a) is identical to (4.2a). However, the cross-correlations in  $\boldsymbol{\Sigma}_{\mathbf{w},\mathbf{w}}$   
 differ. We prove below that this leads to an improved performance of the LS classifier for large  
 sample sizes. First, Theorem 4.2 can be used to check that

$$\boldsymbol{\Sigma}_{\mathbf{w},\mathbf{w}} - \boldsymbol{\Sigma}_{\mathbf{w},\boldsymbol{\mu}} \boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}}^\dagger \boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}}^T \xrightarrow{P} \frac{\gamma}{1-\gamma} (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\pi}\boldsymbol{\pi}^T - (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}} (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi})).$$

291 Thus, the only change in the test error formula compared to (4.3) is the term  $\gamma \cdot \text{diag}(\boldsymbol{\pi})$  substituted  
 292 by the matrix above.

293 **Proposition 4.3** Assume orthogonal, equal-energy means  $\boldsymbol{\Sigma}_{\boldsymbol{\mu},\boldsymbol{\mu}} = \mu^2 \mathbf{I}_k$ ,  $k \geq 2$ . Let  $\gamma_\star = \frac{\mu^2 k}{(k-1)^2} (1 -$   
 294  $k \mathbb{E}\left[\frac{e^{2\mu G_1}}{(\sum_{\ell \in [k]} e^{\mu G_\ell})^2}\right]) \in (0, 1)$ . Then, with probability 1 as  $n \rightarrow \infty$ ,  $\mathbb{P}_{e,\text{LS}} < \mathbb{P}_{e,\text{Avg}} \iff \gamma < \gamma_\star$ .

#### 295 4.3 Weighted least-squares classifier

296 Theorem 4.4 predicts the asymptotic performance of *weighted* least-squares for data generated  
 297 according to MLM.

298 **Theorem 4.4** Consider data generated according to MLM and  $\gamma < 1$ . Consider a weighted LS  
 299 classifier with weights  $\mathbf{D} = \text{diag}(\omega_1, \dots, \omega_k)$  and let  $\eta$  be the unique solution to  $\sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{\omega_\ell^2 + \eta} = \gamma$ . Also  
 300 define vector  $\boldsymbol{\nu} \in \mathbb{R}^k$  with entries given by  $\nu_\ell = \frac{1}{\gamma} \frac{\omega_\ell^2}{\omega_\ell^2 + \eta}$  and positive semidefnite matrix

$$\boldsymbol{\Delta} = \mathbb{E}\left[(\boldsymbol{\nu}^T \mathbf{v}) \mathbf{g} \mathbf{g}^T\right] - \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \boldsymbol{\nu} \boldsymbol{\nu}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V} \boldsymbol{\Sigma} > \mathbf{0}_{r \times r}, \quad (4.5)$$



301 where  $\mathbf{v} \in \mathbb{R}^k$  is a random vector with entries  $V_\ell = \frac{e^{\mathbf{v}^T \boldsymbol{\nu} \Sigma \mathbf{g}}}{\sum_{\ell' \in [k]} e^{\mathbf{v}^T \boldsymbol{\nu} \Sigma \mathbf{g}}}$  for  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_r)$ . Then, for the  
 302 WLS linear classifier  $(\widehat{\mathbf{W}}, \widehat{\mathbf{b}})$  the following asymptotic limits hold

$$\widehat{\mathbf{b}} \xrightarrow{P} \text{diag}(\boldsymbol{\nu})\boldsymbol{\pi} - \text{diag}(\boldsymbol{\nu})(\mathbf{I}_k - \boldsymbol{\pi}\boldsymbol{\nu}^T)(\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi})\mathbf{V}\boldsymbol{\Sigma}\widetilde{\boldsymbol{\Delta}}^{-1}\boldsymbol{\Sigma}\mathbf{V}^T(\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi})\boldsymbol{\nu}. \quad (4.6a)$$

$$\boldsymbol{\Sigma}_{\mathbf{w}, \boldsymbol{\mu}} \xrightarrow{P} \text{diag}(\boldsymbol{\nu})(\mathbf{I}_k - \boldsymbol{\pi}\boldsymbol{\nu}^T)(\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi})\mathbf{V}\boldsymbol{\Sigma}\widetilde{\boldsymbol{\Delta}}^{-1}\boldsymbol{\Sigma}\mathbf{V}^T \quad (4.6b)$$

303 The corresponding formula for the asymptotic limit of the cross-correlation matrix  $\boldsymbol{\Sigma}_{\mathbf{w}, \mathbf{w}}$  is given in  
 304 (J.39) in Section J.

305 To see how the theorem above includes Theorem 4.2 as a special case, consider the case  $\omega_i = 1$ ,  $i \in [k]$ .  
 306 We show how this recovers the solution for (un-weighted) LS. First, note that in this case solving for  
 307  $\eta$  simply gives  $\eta = \frac{1}{\gamma} - 1$ . Thus,  $\boldsymbol{\nu} = \mathbf{1}_k$  and  $\widetilde{\boldsymbol{\pi}} = \boldsymbol{\pi}$ . Also, observe in (4.1) that  $(\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi})\mathbf{1}_k = \mathbf{0}$   
 308 and  $\mathbf{1}^T \boldsymbol{\nu} = 1$ . Thus, (4.5) reduces to  $\boldsymbol{\Delta} = \mathbb{E}[\mathbf{g}\mathbf{g}^T] = \mathbf{I}_r$ . With these, it can be readily checked that  
 309 (4.6a) and (4.6b) simplify to the expressions in (4.4a).

## 310 5 Numerical Results

311 This section validates our theory via numerical experiments and provides further insights on multiclass  
 312 classification. See also the supplementary for more extensive experiments. We study the class-  
 313 wise/total test misclassification error in both GMM and MLM for different sample sizes, number of  
 314 classes and class priors. In line with Section 2.2 we consider four algorithms: (i) Averaging (Avg),  
 315 (ii) LS, (iii) Weighted LS (WLS) with the  $i$ th class weighted by  $\omega_\ell^2 = 1/\pi_\ell$ , (iv) Cross-Entropy (CE).

316 Figures 1 and 2 focus on GMM with  $k = 9$  classes,  $d = 300$  and  $\|\boldsymbol{\mu}_i\|_{\ell_2}^2 = 15$ . To model different  
 317 class prior probabilities, we use the distribution  $\pi_1 = \pi_2 = \pi_3 = 0.5\pi_4 = 0.5\pi_5 = 0.5\pi_6 = 0.25\pi_7 =$   
 318  $0.25\pi_8 = 0.25\pi_9 = 1/21$ . We consider three scenarios: (a) orthogonal means, equal prior ( $\pi_i = 1/9$ );  
 319 (b) orthogonal means, different prior; (c) correlated means with pairwise correlation coefficient  
 320 equal to 0.5 (i.e.,  $\langle \boldsymbol{\mu}_i, \boldsymbol{\mu}_j \rangle / (\|\boldsymbol{\mu}_i\|_{\ell_2} \|\boldsymbol{\mu}_j\|_{\ell_2}) = 0.5$  for  $i \neq j$ ) and different priors as discussed above.

321 Figure 1 shows the test miss-classification errors as a function of  $\gamma := d/n$ . In all scenarios our  
 322 theoretical predictions are a near perfect match to the empirical performance. In scenario (a), class-  
 323 wise averaging achieves the lowest error as predicted by Proposition 3.5. However, in scenario  
 324 (b) where the means have different norms the averaging method has higher misclassification error  
 325 compared with CE, LS and WLS for large sample sizes (small  $\gamma$ ). We note that both LS and WLS  
 326 achieve lower errors compared with CE as the sample size grows. Scenario (c) is similar to (b)  
 327 however due to class correlations, the errors are uniformly higher. Figure 2 shows the corresponding  
 328 class-wise miss-classification errors for the smallest  $\gamma$  in Figure 1 ( $\gamma = 0.117$ ). In scenario (a), errors  
 329 are equal which is expected given the equal class priors. In scenarios (b) and (c) however, due to  
 330 different priors, large classes 7,8,9 achieve best accuracy. The performance difference is most visible  
 331 for the averaging approach. LS mitigates this issue to some extent, while WLS creates the flattest  
 332 class-wise errors suggesting that it can reduce the miss-classification error on small/minority classes.

333 Figure 3 focuses on orthogonal classes with varying number of classes  $k$  where  $\|\boldsymbol{\mu}_i\|_{\ell_2}^2 = 15$  and  
 334  $d \in \{50, 100, 200\}$  with  $kd/n = k\gamma$  fixed at  $k\gamma = 20/11$ . It plots the ratio of the empirical error  
 335 probability and our theoretical prediction as  $k$  grows until  $k = d$ . Two observations are worth  
 336 mentioning here. (1) The accuracy of our predictions noticeably improves as the problem dimension  
 337  $d, n$  grow as expected given the asymptotic nature of our analysis. Interestingly, the convergence  
 338 appears to be noticeably faster (as a function of  $d$ ) for the LS rather than the Averaging classifier. (2)  
 339 Our theoretical results formally require that  $k$  is fixed while  $d$  (and  $n$ ) grow large. Yet, the presented  
 340 experimental results suggest that they might also hold for large  $k$  under the shown scaling. This is a  
 341 fascinating research question that we believe is worth investigating further.

342 Figure 4 provides experiments on the MLM model with  $k = 9$  orthogonal classes. Unlike GMM, CE  
 343 achieves the best performance in MLM. In Figure 4 (a), classes have same norms  $\|\boldsymbol{\mu}_i\|_{\ell_2} = 10$ , while  
 344 in Figure 4 (b) we have quadrupled the norms of classes 7,8,9 and doubled the norms of classes 4,5,6.  
 345 This disparity between the norms seems to help improve the CE accuracy, but hurt LS/averaging  
 346 accuracy for small  $\gamma$ . Finally, Figure 4 (c) shows the class-wise probability of error associated with  
 347 (b) for  $\gamma = 0.117$  and demonstrates that LS outperforms averaging.

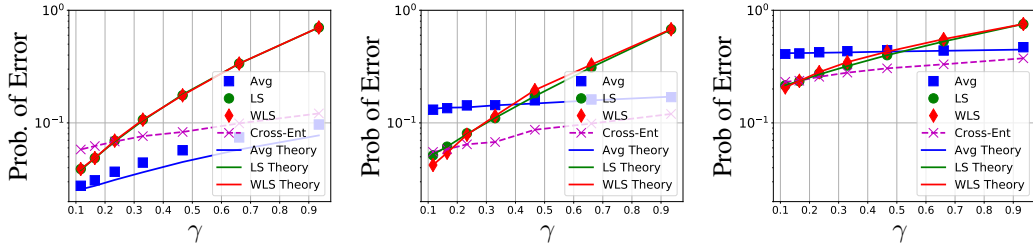


Figure 1: GMM with  $k = 9, d = 300$ . (a) orthogonal, equal prior, (b) orthogonal, different prior, (c) correlated, different prior.

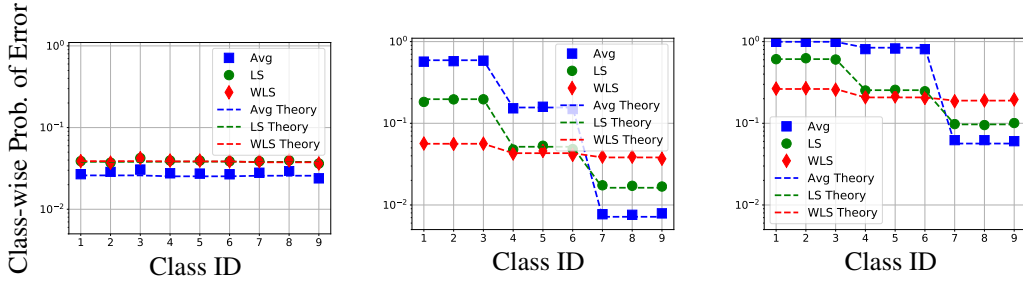


Figure 2: Class-wise probability of errors corresponding to Figure 1 with  $\gamma = 0.117$ .

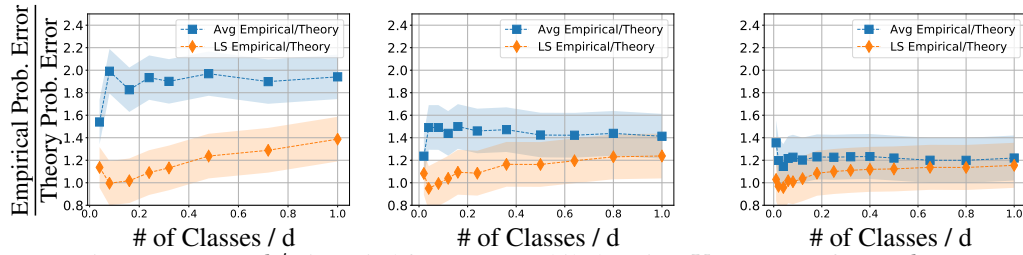


Figure 3: GMM,  $k/p$  is varied from 0 to 1 while keeping  $K\gamma$  constant for (a)  $d = 50$ , (b)  $d = 100$ , (c)  $d = 200$ .

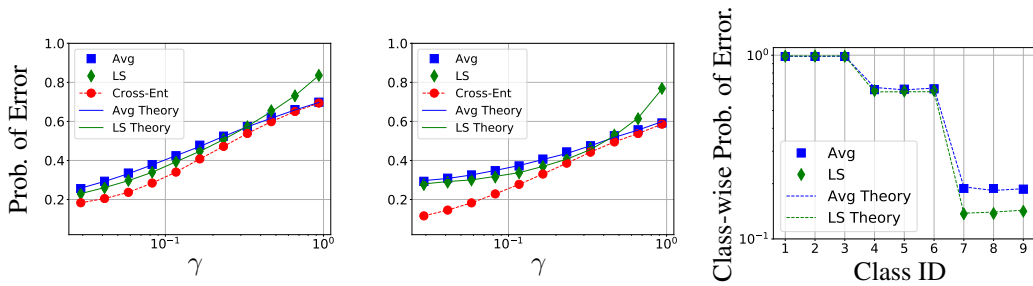


Figure 4: MLM with orthogonal means for (a) equal norms and (b) different norms. (c) Class-wise probability errors for (b).

## 348 6 Future Directions

349 This work aims at initiating a precise asymptotic study of multiclass classifiers that provides a  
 350 promising setting for resolving a rich set of open questions regarding the (comparative) performance  
 351 of classification algorithms as a function of the involved problem variables. As mentioned, even  
 352 understanding the statistical performance of one-vs-all multiclass classifiers does not follow directly  
 353 from the existing literature on binary classifiers. Extending the results of this paper to the one-vs-  
 354 all logistic and SVM classifiers would allow for a principled comparison among these different  
 355 choices. A possibly more challenging, albeit mathematically intriguing and practically relevant  
 356 task, is characterizing the asymptotics of more complicated (non-separable) losses, such as the  
 357 cross-entropy loss. For this, even characterizing the asymptotic behavior of the correlations  $\Sigma_{w,\mu}$

358 requires new ideas. The previously mentioned study of “extreme multiclass classification” in which  
359 the number of classes  $k$  is very large is another very fascinating direction.

## 360 **Broader Impact**

361 In this paper we develop a precise and asymptotically exact understanding of the statistical behavior  
362 of a variety of classification algorithms. In particular we precisely, characterize how the total and  
363 class-wise accuracy varies under different training algorithms, data distributions, problem dimensions,  
364 inter/intra class correlations and class priors. Despite being theoretical/foundational in nature it  
365 has potential for broader practical impact. In particular, our precise characterization of class-wise  
366 accuracy allows us to understand how different training algorithms impact accuracy of machine  
367 learning algorithms on rare/minority classes. Such a precise understanding may help guide the  
368 development of more fair/equitable algorithms. On the flip side, such insights may potentially also be  
369 used nefariously enabling the marginalization of rare/minority classes by developing algorithms that  
370 reduce their class-wise accuracy.

## 371 **References**

- 372 [1] ALLWEIN, E. L., SCHAPIRE, R. E., AND SINGER, Y. Reducing multiclass to binary: A  
373 unifying approach for margin classifiers. *Journal of machine learning research* 1, Dec (2000),  
374 113–141.
- 375 [2] AMELUNXEN, D., LOTZ, M., MCCOY, M. B., AND TROPP, J. A. Living on the edge: A  
376 geometric theory of phase transitions in convex optimization. *arXiv preprint arXiv:1303.6672*  
377 (2013).
- 378 [3] BAYATI, M., AND MONTANARI, A. The lasso risk for gaussian matrices. *Information Theory,*  
379 *IEEE Transactions on* 58, 4 (2012), 1997–2017.
- 380 [4] BOSMAN, A. S., ENGELBRECHT, A., AND HELBIG, M. Visualising basins of attraction for  
381 the cross-entropy and the squared error neural network loss functions. *Neurocomputing* (2020).
- 382 [5] BREDENSTEINER, E. J., AND BENNETT, K. P. Multicategory classification by support vector  
383 machines. In *Computational Optimization*. Springer, 1999, pp. 53–79.
- 384 [6] BU, Z., KLUSOWSKI, J., RUSH, C., AND SU, W. Algorithmic analysis and statistical estimation  
385 of slope via approximate message passing. In *Advances in Neural Information Processing*  
386 *Systems* (2019), pp. 9361–9371.
- 387 [7] CANDÈS, E. J., AND SUR, P. The phase transition for the existence of the maximum likelihood  
388 estimate in high-dimensional logistic regression. *arXiv preprint arXiv:1804.09753* (2018).
- 389 [8] CELENTANO, M., AND MONTANARI, A. Fundamental barriers to high-dimensional regression  
390 with convex penalties. *arXiv preprint arXiv:1903.10603* (2019).
- 391 [9] CHANDRASEKARAN, V., RECHT, B., PARRILO, P. A., AND WILLISKY, A. S. The convex  
392 geometry of linear inverse problems. *Foundations of Computational Mathematics* 12, 6 (2012),  
393 805–849.
- 394 [10] CHOROMANSKA, A., AGARWAL, A., AND LANGFORD, J. Extreme multi class classification.  
395 In *NIPS Workshop: eXtreme Classification, submitted* (2013).
- 396 [11] CORTES, C., KUZNETSOV, V., MOHRI, M., AND YANG, S. Structured prediction theory based  
397 on factor graph complexity. In *Advances in Neural Information Processing Systems* (2016),  
398 pp. 2514–2522.
- 399 [12] COVINGTON, P., ADAMS, J., AND SARGIN, E. Deep neural networks for youtube recom-  
400 mendations. In *Proceedings of the 10th ACM conference on recommender systems* (2016),  
401 pp. 191–198.
- 402 [13] CRAMMER, K., AND SINGER, Y. On the algorithmic implementation of multiclass kernel-based  
403 vector machines. *Journal of machine learning research* 2, Dec (2001), 265–292.
- 404 [14] DEMIRKAYA, A., CHEN, J., AND OYMAK, S. Exploring the role of loss functions in multiclass  
405 classification. In *2020 54th Annual Conference on Information Sciences and Systems (CISS)*  
406 (2020), IEEE, pp. 1–5.

- 407 [15] DENG, Z., KAMMOUN, A., AND THRAMOULIDIS, C. A model of double descent for  
408 high-dimensional binary linear classification. *arXiv preprint arXiv:1911.05822* (2019).
- 409 [16] DIETTERICH, T. G., AND BAKIRI, G. Solving multiclass learning problems via error-correcting  
410 output codes. *Journal of artificial intelligence research* 2 (1994), 263–286.
- 411 [17] DONOHO, D., AND MONTANARI, A. High dimensional robust m-estimation: Asymptotic  
412 variance via approximate message passing. *Probability Theory and Related Fields* 166, 3-4  
413 (2016), 935–969.
- 414 [18] DONOHO, D. L., MALEKI, A., AND MONTANARI, A. The noise-sensitivity phase transition  
415 in compressed sensing. *Information Theory, IEEE Transactions on* 57, 10 (2011), 6920–6941.
- 416 [19] EL KAROUI, N. On the impact of predictor geometry on the performance on high-dimensional  
417 ridge-regularized generalized robust regression estimators. *Probability Theory and Related*  
418 *Fields* 170, 1-2 (2018), 95–175.
- 419 [20] FÜRNKRANZ, J. Round robin classification. *Journal of Machine Learning Research* 2, Mar  
420 (2002), 721–747.
- 421 [21] GAJOWNICZEK, K., CHMIELEWSKI, L. J., ORŁOWSKI, A., AND ZĄBKOWSKI, T. Generalized  
422 entropy cost function in neural networks. In *International Conference on Artificial Neural*  
423 *Networks* (2017), Springer, pp. 128–136.
- 424 [22] GORDON, Y. *On Milman’s inequality and random subspaces which escape through a mesh in*  
425  $\mathbb{R}^n$ . Springer, 1988.
- 426 [23] GUERMEUR, Y. Combining discriminant models with new multi-class svms. *Pattern Analysis*  
427 *& Applications* 5, 2 (2002), 168–179.
- 428 [24] HASHORVA, E., AND HÜSLER, J. On multivariate gaussian tails. *Annals of the Institute of*  
429 *Statistical Mathematics* 55, 3 (2003), 507–522.
- 430 [25] HOU, L., YU, C.-P., AND SAMARAS, D. Squared earth mover’s distance-based loss for  
431 training deep neural networks. *arXiv preprint arXiv:1611.05916* (2016).
- 432 [26] HU, H., AND LU, Y. M. Asymptotics and optimal designs of slope for sparse linear regression.  
433 *arXiv preprint arXiv:1903.11582* (2019).
- 434 [27] HUANG, H. Asymptotic behavior of support vector machine for spiked population model. *The*  
435 *Journal of Machine Learning Research* 18, 1 (2017), 1472–1492.
- 436 [28] JANG, E., GU, S., AND POOLE, B. Categorical reparameterization with gumbel-softmax.  
437 *arXiv preprint arXiv:1611.01144* (2016).
- 438 [29] JAVANMARD, A., SOLTANOLKOTABI, M., AND HASSANI, H. Precise tradeoffs in adversarial  
439 training for linear regression. *arXiv preprint arXiv:2002.10477* (2020).
- 440 [30] KAMMOUN, A., AND ALOUINI, M.-S. On the precise error analysis of support vector machines.  
441 *Submitted to IEEE Transactions on information theory* (2019).
- 442 [31] KAROUI, N. E. Asymptotic behavior of unregularized and ridge-regularized high-dimensional  
443 robust regression estimators: rigorous results. *arXiv preprint arXiv:1311.2445* (2013).
- 444 [32] KOLTCHINSKII, V., PANCHENKO, D., ET AL. Empirical margin distributions and bounding  
445 the generalization error of combined classifiers. *The Annals of Statistics* 30, 1 (2002), 1–50.
- 446 [33] KUMAR, H., AND SASTRY, P. Robust loss functions for learning multi-class classifiers. In  
447 *2018 IEEE International Conference on Systems, Man, and Cybernetics (SMC)* (2018), IEEE,  
448 pp. 687–692.
- 449 [34] KUZNETSOV, V., MOHRI, M., AND SYED, U. Rademacher complexity margin bounds for  
450 learning with a large number of classes. In *ICML Workshop on Extreme Classification: Learning*  
451 *with a Very Large Number of Labels* (2015).
- 452 [35] LEDOUX, M., AND TALAGRAND, M. *Probability in Banach Spaces: isoperimetry and*  
453 *processes*, vol. 23. Springer, 1991.
- 454 [36] LEE, Y., LIN, Y., AND WAHBA, G. Multicategory support vector machines: Theory and  
455 application to the classification of microarray data and satellite radiance data. *Journal of the*  
456 *American Statistical Association* 99, 465 (2004), 67–81.

- 457 [37] LEI, Y., DOGAN, U., BINDER, A., AND KLOFT, M. Multi-class svms: From tighter data-  
458 dependent generalization bounds to novel algorithms. In *Advances in Neural Information*  
459 *Processing Systems* (2015), pp. 2035–2043.
- 460 [38] LEI, Y., DOGAN, Ü., ZHOU, D.-X., AND KLOFT, M. Data-dependent generalization bounds  
461 for multi-class classification. *IEEE Transactions on Information Theory* 65, 5 (2019), 2995–  
462 3021.
- 463 [39] LELARGE, M., AND MIOLANE, L. Asymptotic bayes risk for gaussian mixture in a semi-  
464 supervised setting. *arXiv preprint arXiv:1907.03792* (2019).
- 465 [40] LI, J., LIU, Y., YIN, R., ZHANG, H., DING, L., AND WANG, W. Multi-class learning:  
466 from theory to algorithm. In *Advances in Neural Information Processing Systems* (2018),  
467 pp. 1586–1595.
- 468 [41] LIANG, T., AND SUR, P. A precise high-dimensional asymptotic theory for boosting and  
469 min- $l_1$ -norm interpolated classifiers. *arXiv preprint arXiv:2002.01586* (2020).
- 470 [42] LOLAS, P. Regularization in high-dimensional regression and classification via random matrix  
471 theory. *arXiv preprint arXiv:2003.13723* (2020).
- 472 [43] MAI, X., LIAO, Z., AND COUILLET, R. A large scale analysis of logistic regression: asymptotic  
473 performance and new insights. In *ICASSP* (2019).
- 474 [44] MAI, X., LIAO, Z., AND COUILLET, R. A large scale analysis of logistic regression: Asymp-  
475 totic performance and new insights. In *ICASSP 2019-2019 IEEE International Conference on*  
476 *Acoustics, Speech and Signal Processing (ICASSP)* (2019), IEEE, pp. 3357–3361.
- 477 [45] MAURER, A. A vector-contraction inequality for rademacher complexities. In *International*  
478 *Conference on Algorithmic Learning Theory* (2016), Springer, pp. 3–17.
- 479 [46] MEI, J., XIAO, C., SZEPESVARI, C., AND SCHUURMANS, D. On the global convergence  
480 rates of softmax policy gradient methods. *arXiv preprint arXiv:2005.06392* (2020).
- 481 [47] MIGNACCO, F., KRZAKALA, F., LU, Y. M., AND ZDEBOROVÁ, L. The role of regularization  
482 in classification of high-dimensional noisy gaussian mixture. *arXiv preprint arXiv:2002.11544*  
483 (2020).
- 484 [48] MIOLANE, L., AND MONTANARI, A. The distribution of the lasso: Uniform control over  
485 sparse balls and adaptive parameter tuning. *arXiv preprint arXiv:1811.01212* (2018).
- 486 [49] MONTANARI, A., RUAN, F., SOHN, Y., AND YAN, J. The generalization error of max-margin  
487 linear classifiers: High-dimensional asymptotics in the overparametrized regime. *arXiv preprint*  
488 *arXiv:1911.01544* (2019).
- 489 [50] OYMAK, S., RECHT, B., AND SOLTANOLKOTABI, M. Sharp time–data tradeoffs for linear  
490 inverse problems. *IEEE Transactions on Information Theory* 64, 6 (2017), 4129–4158.
- 491 [51] OYMAK, S., THRAMOULIDIS, C., AND HASSIBI, B. The squared-error of generalized lasso:  
492 A precise analysis. *arXiv preprint arXiv:1311.0830* (2013).
- 493 [52] PAL, M., AND MATHER, P. Support vector machines for classification in remote sensing.  
494 *International journal of remote sensing* 26, 5 (2005), 1007–1011.
- 495 [53] PIRES, B. Á., AND SZEPESVÁRI, C. Multiclass classification calibration functions. *arXiv*  
496 *preprint arXiv:1609.06385* (2016).
- 497 [54] PIRES, B. A., SZEPESVARI, C., AND GHAVAMZADEH, M. Cost-sensitive multiclass classifi-  
498 cation risk bounds. In *International Conference on Machine Learning* (2013), pp. 1391–1399.
- 499 [55] RAWAT, A. S., CHEN, J., YU, F. X. X., SURESH, A. T., AND KUMAR, S. Sampled softmax  
500 with random fourier features. In *Advances in Neural Information Processing Systems* 32,  
501 H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds. Curran  
502 Associates, Inc., 2019, pp. 13834–13844.
- 503 [56] RIFKIN, R., AND KLAUTAU, A. In defense of one-vs-all classification. *Journal of machine*  
504 *learning research* 5, Jan (2004), 101–141.
- 505 [57] SALEHI, F., ABBASI, E., AND HASSIBI, B. The impact of regularization on high-dimensional  
506 logistic regression. *arXiv preprint arXiv:1906.03761* (2019).
- 507 [58] SATHE, Y., AND LINGRAS, S. A note on the inequalities for tail probability of the multivariate  
508 normal distribution. *Communications in Statistics-Theory and Methods* 9, 7 (1980), 711–715.

- 509 [59] STOJNIC, M. Various thresholds for  $\ell_1$ -optimization in compressed sensing. *arXiv preprint*  
510 *arXiv:0907.3666* (2009).
- 511 [60] STOJNIC, M. A framework to characterize performance of lasso algorithms. *arXiv preprint*  
512 *arXiv:1303.7291* (2013).
- 513 [61] SUR, P., AND CANDÈS, E. J. A modern maximum-likelihood theory for high-dimensional  
514 logistic regression. *Proceedings of the National Academy of Sciences* 116, 29 (2019), 14516–  
515 14525.
- 516 [62] SUTSKEVER, I., VINYALS, O., AND LE, Q. V. Sequence to sequence learning with neural  
517 networks. In *Advances in neural information processing systems* (2014), pp. 3104–3112.
- 518 [63] TAHERI, H., PEDARSANI, R., AND THRAMPOULIDIS, C. Sharp asymptotics and optimal  
519 performance for inference in binary models. *arXiv preprint arXiv:2002.07284* (2020).
- 520 [64] TEWARI, A., AND BARTLETT, P. L. On the consistency of multiclass classification methods.  
521 *Journal of Machine Learning Research* 8, May (2007), 1007–1025.
- 522 [65] THRAMPOULIDIS, C., ABBASI, E., AND HASSIBI, B. Lasso with non-linear measurements is  
523 equivalent to one with linear measurements. In *Advances in Neural Information Processing*  
524 *Systems* (2015), pp. 3420–3428.
- 525 [66] THRAMPOULIDIS, C., OYMAK, S., AND HASSIBI, B. Regularized linear regression: A precise  
526 analysis of the estimation error. In *Proceedings of The 28th Conference on Learning Theory*  
527 (2015), pp. 1683–1709.
- 528 [67] WANG, S., WENG, H., AND MALEKI, A. Does slope outperform bridge regression? *arXiv*  
529 *preprint arXiv:1909.09345* (2019).
- 530 [68] WESTON, J., AND WATKINS, C. Multi-class support vector machines. Tech. rep., Citeseer,  
531 1998.
- 532 [69] WISLER, A., BERISHA, V., WEI, D., RAMAMURTHY, K., AND SPANIAS, A. Empirically-  
533 estimable multi-class classification bounds. In *2016 IEEE International Conference on Acoustics,*  
534 *Speech and Signal Processing (ICASSP)* (2016), IEEE, pp. 2594–2598.
- 535 [70] YEN, I. E.-H., HUANG, X., RAVIKUMAR, P., ZHONG, K., AND DHILLON, I. Pd-sparse:  
536 A primal and dual sparse approach to extreme multiclass and multilabel classification. In  
537 *International Conference on Machine Learning* (2016), pp. 3069–3077.
- 538 [71] ZHANG, T. Statistical analysis of some multi-category large margin classification methods.  
539 *Journal of Machine Learning Research* 5, Oct (2004), 1225–1251.

540 **Contents**

541	<b>A Additional Numerical Results</b>	<b>17</b>
542	<b>B Preliminaries</b>	<b>18</b>
543	B.1 Slepian’s inequality . . . . .	18
544	B.2 Useful results for MLM . . . . .	18
545	B.3 Block matrix inversions . . . . .	18
546	<b>C Calculating and bounding the missclassification error</b>	<b>19</b>
547	C.1 MLM . . . . .	19
548	C.1.1 Class-wise miss-classification error . . . . .	19
549	C.1.2 A useful expression for the total miss-classification error . . . . .	20
550	C.2 Class-wise and total miss-classification error for GMM . . . . .	20
551	C.3 Evaluating and bounding tail probabilities of multivariate Gaussians . . . . .	20
552	C.3.1 A special case: Rank-one update of Identity . . . . .	21
553	C.3.2 Slepian’s bound . . . . .	21
554	C.3.3 Union bound . . . . .	21
555	<b>D The Class-averaging estimator</b>	<b>22</b>
556	D.1 Proofs for GMM . . . . .	22
557	D.1.1 GMM: Proof of Proposition 3.1 . . . . .	22
558	D.2 Proofs for MLM . . . . .	22
559	D.2.1 Proof of Proposition 4.1 . . . . .	22
560	D.2.2 Orthogonal means . . . . .	24
561	<b>E On the Bayes risk of GMM: Proof of Proposition 3.5</b>	<b>24</b>
562	<b>F Proof outline for (weighted) least-squares: key ideas and challenges</b>	<b>26</b>
563	F.1 Background on CGMT . . . . .	28
564	<b>G Least-squares for GMM</b>	<b>28</b>
565	G.1 Proof of Theorem 3.2 . . . . .	28
566	G.1.1 Computing $\Sigma_{w,\mu}$ . . . . .	28
567	G.1.2 Computing $\Sigma_{w,w}$ . . . . .	31
568	G.2 Orthogonal means . . . . .	33
569	<b>H Least-squares for MLM</b>	<b>33</b>
570	H.1 Proof of Theorem 4.2 . . . . .	33
571	H.1.1 Computing $\Sigma_{w,\mu}$ . . . . .	33
572	H.1.2 Computing $\Sigma_{w,w}$ . . . . .	35
573	H.2 Orthogonal means and equal-energy . . . . .	36

574	H.3 Proof of Proposition 4.3 . . . . .	36
575	<b>I Weighted LS for GMM (Proof of Theorem 3.4)</b>	<b>40</b>
576	I.1 Computing $\Sigma_{w,\mu}$ . . . . .	40
577	I.2 Computing $\Sigma_{w,w}$ . . . . .	45
578	<b>J Weighted LS for MLM (Proof of Theorem 4.4)</b>	<b>49</b>
579	J.1 Deterministic Analysis of the AO . . . . .	51
580	J.3 Computing cross-correlations $\Sigma_{w,w}$ . . . . .	53



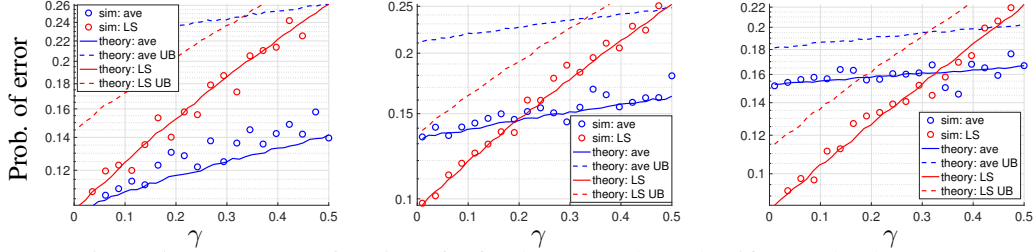


Figure 5: Test error as a function of  $\gamma$  for the Avg and LS classifiers under the GMM. Precise theoretical predictions compared to a theoretical upper bound (UB) obtained by Union bound that does not require knowledge of the off-diagonal entries of  $\Sigma_{w,w}$ . (a)  $\pi_1 = \pi_2 = \pi_3 = 1/3$ ; (b)  $2\pi_1 = \pi_2 = 2\pi_3$ ; (c)  $4\pi_1 = \pi_2 = 4\pi_3$

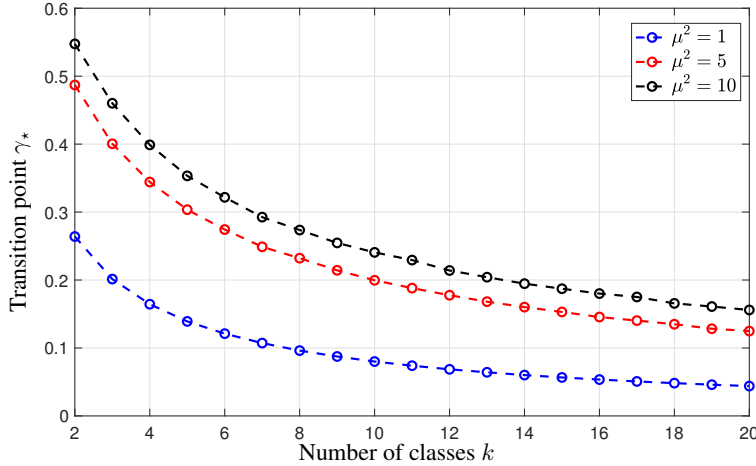


Figure 6: The threshold  $\gamma_*$  of Proposition 4.3 as a function of the number of classes  $k$  and the means' energy  $\mu$ . LS provably outperforms class-averaging for  $\gamma < \gamma_*$ .

## 581 A Additional Numerical Results

582 In this section, we provide further numerical experiments. The first fundamental question we ask is

*To what extent pairwise class correlations are necessary for performance prediction?*

583 In Figure 5, we attempt to answer this question in a GMM setup with  $k = 3$  classes. The solid lines  
 584 are exact performance predictions based on our theory for averaging and least-squares estimators.  
 585 The dashed lines are the theoretical upper bounds which does not require the knowledge of cross-  
 586 correlations between the classes (i.e. off-diagonal entries of  $\Sigma_{w,w}$  are unknown). These bounds  
 587 are calculated by considering the *worst case correlation matrix*  $\Sigma_{w,w}$  given  $\Sigma_{w,\mu}$  and the self-  
 588 correlations of the classes i.e. the diagonal entries of  $\Sigma_{w,w}$ . Both of these information can be obtained  
 589 by studying the properties of isolated least-squares on individual classes without understanding their  
 590 pairwise relations.

591 Figure 5 has two important takeaways. The first takeaway is the fact that pairwise correlations are  
 592 indeed critical for exact asymptotic analysis and naive approaches cannot reproduce the results of our  
 593 novel multiclass learning framework. There is a visible gap between upper and exact bounds and this  
 594 gap is particularly more visible for the averaging estimator. The second takeaway is that our upper  
 595 bounds are surprisingly strong. In most regimes, our upper bound is within a factor of 1.5 to 2 of the  
 596 actual bound. These observations suggest a simpler method to calculate reasonably strong theoretical  
 597 bounds.

598 The second fundamental question is

*When does least-squares provably outperform averaging?*

599 Our Proposition 4.3 provides a fundamental transition point in sample complexity above which  
600 least-squares is provably better than averaging under MLM. In Figure 6, we visualize  $\gamma_*$  as a function  
601 of different number of classes as well as different levels of mean energy. Least-squares outperform  
602 averaging in the region below the lines displayed in Figure 6. Our key message is that least-squares  
603 work better when the sample complexity is higher and the problem is less noisy. As the number  
604 of classes  $k$  increase, the problem becomes more difficult/noisy and we require a larger sample  
605 complexity to ensure classifier achieves a similar amount of accuracy as small  $k$ . Following this  
606 intuition, as  $k$  increases,  $\gamma_*$  shifts smaller due to larger sample requirement. Similarly energy  $\mu$   
607 directly controls the noise level of the problem i.e. larger  $\mu$  results in a larger signal-to-noise ratio.  
608 Thus, as we increase  $\mu$ ,  $\gamma_*$  increases as well because same test accuracy can be achieved with smaller  
609 sample size.

## 610 B Preliminaries

611 In this section we gather a few preliminary results that will be used later on in our proofs.

### 612 B.1 Slepian's inequality

**Lemma B.1 (Slepian's inequality [35])** Let  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{S})$  and  $\tilde{\mathbf{g}} \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$  such that for all  $i, j \in [k]$ :

$$\mathbf{S}_{ii} = \mathbf{R}_{ii}, \quad \text{and} \quad \mathbf{S}_{ij} \leq \mathbf{R}_{ij}.$$

Then, for any  $\mathbf{t} \in \mathbb{R}^k$  it holds that

$$\mathbb{P} \left\{ \bigcup_{j \in [k]} \{\mathbf{g}_j \geq \mathbf{t}_j\} \right\} \leq \mathbb{P} \left\{ \bigcup_{j \in [k]} \{\tilde{\mathbf{g}}_j \geq \mathbf{t}_j\} \right\}.$$

or equivalently,

$$\mathbb{P} \left\{ \bigcap_{j \in [k]} \{\mathbf{g}_j \leq \mathbf{t}_j\} \right\} \geq \mathbb{P} \left\{ \bigcap_{j \in [k]} \{\tilde{\mathbf{g}}_j \leq \mathbf{t}_j\} \right\}.$$

### 613 B.2 Useful results for MLM

614 **Lemma B.2 (Gaussian integration by parts)** Let  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)$  and random vector  $\mathbf{v} =$   
615  $[V_1, V_2, \dots, V_k]^T$  with entries:

$$\mathbf{v} = \frac{e^{\mathbf{V}\Sigma\mathbf{g}}}{\mathbf{1}_k^T e^{\mathbf{V}\Sigma\mathbf{g}}}, \quad V_\ell = \frac{e^{e_\ell^T \mathbf{V}\Sigma\mathbf{g}}}{\sum_{j \in [k]} e^{e_j^T \mathbf{V}\Sigma\mathbf{g}}}, \quad \ell \in [k]. \quad (\text{B.1})$$

616 Further recall the notation of  $\boldsymbol{\pi}$  and  $\boldsymbol{\Pi}$  in (4.1). The following statements are true:

617 (i)  $\mathbb{E}[\mathbf{v}] = \boldsymbol{\pi}$ .

(ii) For all  $i \in [r], \ell \in [k]$ :

$$\mathbb{E}[\mathbf{g}_i V_\ell] = (e_i^T \Sigma \mathbf{V}^T e_\ell) \boldsymbol{\pi}_\ell - e_i^T \Sigma \mathbf{V}^T \sum_{j \in [k]} e_j \boldsymbol{\Pi}_{ij},$$

and in matrix form:

$$\mathbb{E}[\mathbf{v}\mathbf{g}^T] = (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V}\Sigma$$

618 (ii).  $\mathbb{E}[\mathbf{v}\mathbf{g}^T] = (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V}\Sigma$ .

### 619 B.3 Block matrix inversions

620 **Lemma B.3 (Block matrix inversion)** Let  $\mathbf{T} = \begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{b}^T & \delta \end{bmatrix}$  be an invertible block matrix. Then

$$\mathbf{T}^{-1} \begin{bmatrix} \mathbf{f} \\ \epsilon \end{bmatrix} = \begin{bmatrix} \Delta^{-1} (\mathbf{f} - \frac{\epsilon}{\delta} \mathbf{b}) \\ \frac{\epsilon}{\delta} - \frac{1}{\delta} \mathbf{b}^T \Delta^{-1} (\mathbf{f} - \frac{\epsilon}{\delta} \mathbf{b}) \end{bmatrix} \quad (\text{B.2})$$

621 where  $\Delta = \mathbf{A} - \frac{1}{\delta} \mathbf{b}\mathbf{b}^T > \mathbf{0}$  is the Schur complement.

622 **C Calculating and bounding the missclassification error**

623 **C.1 MLM**

624 **C.1.1 Class-wise miss-classification error**

625 The class-wise misclassification error can be calculated as follows

$$\begin{aligned}
\mathbb{P}\{\widehat{Y} \neq Y|Y = c\} &= \mathbb{E}_{\mathbf{x}} \left[ \mathbb{P}\{\widehat{Y} \neq Y|Y = c, \mathbf{x}\} \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ \frac{\mathbb{P}\{\widehat{Y} \neq Y, Y = c|\mathbf{x}\}}{\mathbb{P}\{Y = c|\mathbf{x}\}} \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ \frac{\mathbb{P}\{\widehat{Y} \neq c|\mathbf{x}\}}{\mathbb{P}\{Y = c|\mathbf{x}\}} \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ \frac{\mathbb{P}\{\arg \max(\widehat{\mathbf{W}}\mathbf{x} + \widehat{\mathbf{b}}) \neq c|\mathbf{x}\}}{\mathbb{P}\{Y = c|\mathbf{x}\}} \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ \frac{1 - \prod_{\ell \neq c} \mathbb{1}_{\{\widehat{\mathbf{w}}_{\ell}^T \mathbf{x} + \widehat{b}_{\ell} \leq \widehat{\mathbf{w}}_c^T \mathbf{x} + \widehat{b}_c\}}}{\mathbb{P}\{Y = c|\mathbf{x}\}} \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ \left( 1 + \sum_{j \neq c} e^{(\boldsymbol{\mu}_j - \boldsymbol{\mu}_c)^T \mathbf{x}} \right) \left( 1 - \prod_{\ell \neq c} \mathbb{1}_{\{\widehat{\mathbf{w}}_{\ell}^T \mathbf{x} + \widehat{b}_{\ell} \leq \widehat{\mathbf{w}}_c^T \mathbf{x} + \widehat{b}_c\}} \right) \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ 1 - \prod_{\ell \neq c} \mathbb{1}_{\{\widehat{\mathbf{w}}_{\ell}^T \mathbf{x} + \widehat{b}_{\ell} \leq \widehat{\mathbf{w}}_c^T \mathbf{x} + \widehat{b}_c\}} \right] + \sum_{j \neq c} \mathbb{E}_{\mathbf{x}} \left[ e^{(\boldsymbol{\mu}_j - \boldsymbol{\mu}_c)^T \mathbf{x}} \left( 1 - \prod_{\ell \neq c} \mathbb{1}_{\{\widehat{\mathbf{w}}_{\ell}^T \mathbf{x} + \widehat{b}_{\ell} \leq \widehat{\mathbf{w}}_c^T \mathbf{x} + \widehat{b}_c\}} \right) \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ 1 - \prod_{\ell \neq c} \mathbb{1}_{\{\widehat{\mathbf{w}}_{\ell}^T \mathbf{x} + \widehat{b}_{\ell} \leq \widehat{\mathbf{w}}_c^T \mathbf{x} + \widehat{b}_c\}} \right] + \sum_{j \neq c} e^{\frac{\|\boldsymbol{\mu}_j - \boldsymbol{\mu}_c\|_{\ell_2}^2}{2}} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_c, \mathbf{I})} \left[ 1 - \prod_{\ell \neq c} \mathbb{1}_{\{\widehat{\mathbf{w}}_{\ell}^T \mathbf{x} + \widehat{b}_{\ell} \leq \widehat{\mathbf{w}}_c^T \mathbf{x} + \widehat{b}_c\}} \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ 1 - \prod_{\ell \neq c} \mathbb{1}_{\{\widehat{\mathbf{w}}_{\ell}^T \mathbf{x} + \widehat{b}_{\ell} \leq \widehat{\mathbf{w}}_c^T \mathbf{x} + \widehat{b}_c\}} \right] \\
&\quad + \sum_{j \neq c} e^{\frac{\|\boldsymbol{\mu}_j - \boldsymbol{\mu}_c\|_{\ell_2}^2}{2}} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ 1 - \prod_{\ell \neq c} \mathbb{1}_{\{\widehat{\mathbf{w}}_{\ell}^T (\boldsymbol{\mu}_j - \boldsymbol{\mu}_c + \mathbf{x}) + \widehat{b}_{\ell} \leq \widehat{\mathbf{w}}_c^T (\boldsymbol{\mu}_j - \boldsymbol{\mu}_c + \mathbf{x}) + \widehat{b}_c\}} \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ 1 - \prod_{\ell \neq c} \mathbb{1}_{\{(\widehat{\mathbf{w}}_{\ell} - \widehat{\mathbf{w}}_c)^T \mathbf{x} \leq \widehat{b}_c - \widehat{b}_{\ell}\}} \right] \\
&\quad + \sum_{j \neq c} e^{\frac{\|\boldsymbol{\mu}_j - \boldsymbol{\mu}_c\|_{\ell_2}^2}{2}} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ 1 - \prod_{\ell \neq c} \mathbb{1}_{\{(\widehat{\mathbf{w}}_{\ell} - \widehat{\mathbf{w}}_c)^T \mathbf{x} \leq (\widehat{\mathbf{w}}_{\ell} - \widehat{\mathbf{w}}_c)^T (\boldsymbol{\mu}_c - \boldsymbol{\mu}_j) + \widehat{b}_c - \widehat{b}_{\ell}\}} \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ \prod_{\ell \neq c} \mathbb{1}_{\{(\widehat{\mathbf{w}}_{\ell} - \widehat{\mathbf{w}}_c)^T \mathbf{x} \geq \widehat{b}_c - \widehat{b}_{\ell}\}} \right] \\
&\quad + \sum_{j \neq c} e^{\frac{\|\boldsymbol{\mu}_j - \boldsymbol{\mu}_c\|_{\ell_2}^2}{2}} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \prod_{\ell \neq c} \mathbb{1}_{\{(\widehat{\mathbf{w}}_{\ell} - \widehat{\mathbf{w}}_c)^T \mathbf{x} \geq (\widehat{\mathbf{w}}_{\ell} - \widehat{\mathbf{w}}_c)^T (\boldsymbol{\mu}_c - \boldsymbol{\mu}_j) + \widehat{b}_c - \widehat{b}_{\ell}\}} \right]
\end{aligned}$$

626 To continue define  $\mathbf{S}_c \in \mathbb{R}^{(k-1) \times (k-1)}$  to be a symmetric matrix such that its  $i, j$  element is  $[\mathbf{S}_c]_{ij} :=$   
627  $\langle \widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_j, \widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_i \rangle$  and  $\mathbf{t}_c^{(j)} \in \mathbb{R}^{k-1}$  a vector with entries  $[\mathbf{t}_c^{(j)}]_i := (\widehat{\mathbf{w}}_i - \widehat{\mathbf{w}}_c)^T (\boldsymbol{\mu}_c - \boldsymbol{\mu}_j) + \widehat{b}_c - \widehat{b}_i$ .

628 Thus,

$$\mathbb{P}_{e|c} := \mathbb{P}\{\widehat{Y} \neq Y | Y = c\} = \sum_{j=1}^k e^{-\frac{\|\mu_j - \mu_c\|_{\ell_2}^2}{2}} \mathbb{P}\{\mathbf{S}_c^{\frac{1}{2}} \mathbf{z} \geq \mathbf{t}_c^{(j)}\}$$

629 Also note that

$$\pi_c = \mathbb{P}\{Y = c\} = \mathbb{E}_{\mathbf{x}} \left[ \mathbb{P}\{Y = c | \mathbf{x}\} \right] = \mathbb{E}_{\mathbf{x}} \left[ \frac{1}{(1 + \sum_{j \neq c} e^{(\mu_j - \mu_c)^T \mathbf{x}})} \right]$$

### 630 C.1.2 A useful expression for the total miss-classification error

631 We can rewrite (2.8) as follows for  $\tilde{\mathbf{g}} \sim \mathcal{N}(\mathbf{0}_k, \mathbf{I}_k)$  and  $\mathbf{h} \sim \mathcal{N}(\mathbf{0}_k, \Sigma_{\mu, \mu})$

$$\begin{aligned} \mathbb{P}_e &= \mathbb{P}\left\{ \arg \max_{j \in [k]} (\mathbf{g}_j + \widehat{\mathbf{b}}_j) \neq Y(\mathbf{h}) \right\}, \\ &= \mathbb{P}\left\{ \arg \max \left( (\Sigma_{w, w} - \Sigma_{w, \mu} \Sigma_{\mu, \mu}^\dagger \Sigma_{w, \mu})^{1/2} \tilde{\mathbf{g}} + \Sigma_{w, \mu} \Sigma_{\mu, \mu}^\dagger \mathbf{h} + \widehat{\mathbf{b}} \right) \neq Y(\mathbf{h}) \right\} \\ &= \mathbb{P}\left\{ \arg \max \left( (\Sigma_{w, w} - \Sigma_{w, \mu} \Sigma_{\mu, \mu}^\dagger \Sigma_{w, \mu})^{1/2} \tilde{\mathbf{g}} + \Sigma_{w, \mu} \Sigma_{\mu, \mu}^\dagger \mathbf{h} + \widehat{\mathbf{b}} \right) \neq Y(\mathbf{h}) \right\} \\ &= 1 - \sum_{c \in [k]} \mathbb{E} \left[ \frac{e^{\mathbf{h}_c}}{\sum_{\ell \in [k]} e^{\mathbf{h}_\ell}} \mathbb{P}\{\mathbf{S}_c^{1/2} \mathbf{z} \geq \mathbf{t}(\mathbf{h})\} \right], \end{aligned} \quad (\text{C.1})$$

632 where in the last line  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{k-1})$ :

$$\begin{aligned} [\mathbf{t}(\mathbf{h})]_j &= \widehat{\mathbf{b}}_j - \widehat{\mathbf{b}}_c + [\Sigma_{w, \mu} \Sigma_{\mu, \mu}^\dagger \mathbf{h}]_j - [\Sigma_{w, \mu} \Sigma_{\mu, \mu}^\dagger \mathbf{h}]_c, \quad j \in [k-1] \\ [\mathbf{S}_c]_{i, j} &= (\mathbf{e}_c - \mathbf{e}_j)^T (\Sigma_{w, w} - \Sigma_{w, \mu} \Sigma_{\mu, \mu}^\dagger \Sigma_{w, \mu}) (\mathbf{e}_c - \mathbf{e}_i), \quad i, j \in [k-1]. \end{aligned}$$

### 633 C.2 Class-wise and total miss-classification error for GMM

634 The class-wise miss-classification error is given by

$$\mathbb{P}_{e|c} = \mathbb{P}(\exists j \neq c : \langle \widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_j, \mathbf{z} \rangle \leq \langle \widehat{\mathbf{w}}_j - \widehat{\mathbf{w}}_c, \boldsymbol{\mu}_c \rangle + (\widehat{\mathbf{b}}_j - \widehat{\mathbf{b}}_c)) \quad (\text{C.2})$$

$$= 1 - \mathbb{P}(\forall j \neq c : \langle \widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_j, \mathbf{z} \rangle \geq \langle \widehat{\mathbf{w}}_j - \widehat{\mathbf{w}}_c, \boldsymbol{\mu}_c \rangle + (\widehat{\mathbf{b}}_j - \widehat{\mathbf{b}}_c)). \quad (\text{C.3})$$

635 where we used that  $\mathbf{x} = \boldsymbol{\mu}_c + \mathbf{z}$ . Let  $\mathbf{S}_c \in \mathbb{R}^{(k-1) \times (k-1)}$  be a symmetric matrix such that its  $i, j$  element  
636 is  $[\mathbf{S}_c]_{i, j} := \langle \widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_j, \widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_i \rangle$  and  $\mathbf{t}_c \in \mathbb{R}^{k-1}$  a vector with entries  $[\mathbf{t}_c]_i := \langle \widehat{\mathbf{w}}_i - \widehat{\mathbf{w}}_c, \boldsymbol{\mu}_c \rangle + (\widehat{\mathbf{b}}_i - \widehat{\mathbf{b}}_c)$ .  
637 Then, we can rewrite (C.3) as

$$\mathbb{P}_{e|c} := 1 - \mathbb{P}\{\mathbf{S}_c^{1/2} \mathbf{z} \geq \mathbf{t}_c\}, \quad (\text{C.4})$$

638 where the inequality in the rightmost expression applies entry-wise.

639 Next we focus on calculating the total miss-classification error. To this aim by using the law of total  
640 probability we have

$$\mathbb{P}_e = \sum_{c=1}^k \pi_c \mathbb{P}_{e|c} = 1 - \sum_{c=1}^k \pi_c \mathbb{P}\{\mathbf{S}_c^{1/2} \mathbf{z} \geq \mathbf{t}_c\}.$$

### 641 C.3 Evaluating and bounding tail probabilities of multivariate Gaussians

642 In Sections C.1.2 and C.2, we expressed the class-wise probability of missclassification error for both  
643 GMM and MLM in the following convenient form

$$1 - \mathbb{P}\{\mathbf{A} \mathbf{g} \leq \mathbf{t}\} = \mathbb{P}\left\{ \bigcup_{i \in [k-1]} \{\mathbf{a}_i^T \mathbf{g} \geq t_i\} \right\}. \quad (\text{C.5})$$

644 Here,  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{k-1})$  and  $\mathbf{A} = [\mathbf{a}_1 \ \cdots \ \mathbf{a}_{k-1}]^T \in \mathbb{R}^{(k-1) \times (k-1)}$ ,  $\mathbf{t} \in \mathbb{R}^{k-1}$  are appropriate coefficient  
645 matrices; see (C.4) and (C.1).

646 The formulation above is convenient both in our theoretical analysis, as well as, in simulations. In  
647 the rest of this section, we briefly discuss some relevant tools that allow to further simplify or bound  
648 expressions in the form of (C.5).

649 **C.3.1 A special case: Rank-one update of Identity**

650 First, we discuss the case where the coefficient matrix  $\mathbf{A}$  and vector  $\mathbf{t}$  in (C.5) take the special form  
 651  $\mathbf{A} \propto \mathbf{I} + \mathbf{1}\mathbf{1}^T$  and  $\mathbf{t} \propto \mathbf{1}$ . This special case often appears in some of the stylized symmetric problem  
 652 settings studied in this paper, such as classification problems with orthogonal and equally-balanced  
 653 means.

654 **Lemma C.1** *Let  $\mathbf{A} = \mathbf{I}_k + \mathbf{1}_k\mathbf{1}_k^T$  and  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{A})$ . Then, for any  $t \in \mathbb{R}$ ,*

$$1 - \mathbb{P}\{\mathbf{g} \leq t\mathbf{1}_k\} = \mathbb{P}\{G_0 + \max_{i \in [k]} G_i \geq t\}, \quad G_0, G_1, \dots, G_k \stackrel{iid}{\sim} \mathcal{N}(0, 1). \quad (\text{C.6})$$

**Proof** It is easily checked that we can decompose for each  $i \in [n]$ :

$$\mathbf{g}_i = G_0 + G_i,$$

655 where  $G_0, G_1, \dots, G_k$  are iid standard normals. Indeed, it can be readily checked from this that  
 656  $\mathbb{E}[\mathbf{g}_i^2] = 2$  and  $\mathbb{E}[\mathbf{g}_i\mathbf{g}_j] = \mathbb{E}[G_0^2] = 1$ ,  $i \neq j$ , which is consistent with  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{A})$ .

Thus, we can write

$$1 - \mathbb{P}\{\mathbf{g} \leq t\mathbf{1}_k\} = \mathbb{P}\{\max_{i \in [k]} \mathbf{g}_i \geq t\} = \mathbb{P}\{G_0 + \max_{i \in [k]} G_i \geq t\},$$

657 to complete the proof. ■

658 **C.3.2 Slepian's bound**

659 When the matrix  $\mathbf{A}$  does not have the special structure assumed by Lemma C.1, it is not possible  
 660 in general to provide simple expressions as the one in (C.6). Yet, it might be possible to obtain  
 661 upper bounds of the same simple form. Such simple bounds can be useful for theoretical interpreta-  
 662 tions of otherwise complicated formulae, or can provide efficient means for quick (but, non-tight)  
 663 implementations.

In this section, we discuss Slepian's inequality (see B.1) as a useful tool in this direction. Assume for  
 simplicity that  $a = \min_{i \neq j \in [k]} \mathbf{A}_{ij} \geq 0$  and note that

$$\mathbf{A} \geq (\text{diag}(\mathbf{A}) - a\mathbf{I}) + a\mathbf{1}\mathbf{1}^T,$$

664 where the inequality above holds element-wise and equality is true for the diagonal elements. Then,  
 665 one can apply Slepian's Lemma B.1 to upper bound the conditional probability of error in (C.5) with  
 666 the following simple bound:

$$\begin{aligned} 1 - \mathbb{P}\{\mathbf{A}\mathbf{g} \leq \mathbf{t}\} &\leq 1 - \mathbb{P}\left\{\left((\text{diag}(\mathbf{A}) - a\mathbf{I}) + a\mathbf{1}\mathbf{1}^T\right)^{1/2} \mathbf{g} \leq \mathbf{t}\right\} \\ &\leq \mathbb{P}\left\{\bigcup_{j \in [k]} \left\{G_0 + G_j \sqrt{[\mathbf{A}]_{jj}/a} \geq [\mathbf{t}]_j/a\right\}\right\}, \quad G_0, G_1, \dots, G_k \stackrel{iid}{\sim} \mathcal{N}(0, 1). \end{aligned}$$

667 In the second line above, we used the Gaussian decomposition of Lemma C.1.

668 **C.3.3 Union bound**

669 Of course, it is also possible to apply (a simpler) union bound to upper bound the tail probability in  
 670 (C.5). Here, we show explicitly the result of applying union bound to the class-wise error probabilities  
 671 of the GMM model. Specifically, consider (C.2). An application of the union bound leads to the  
 672 following:

$$\begin{aligned} \mathbb{P}_{e|c} &\leq \sum_{j \neq c} \mathbb{P}_{G \sim \mathcal{N}(0,1)} \left\{ \|\widehat{\mathbf{w}}_j - \widehat{\mathbf{w}}_c\|_{\ell_2} G \leq \langle \widehat{\mathbf{w}}_j - \widehat{\mathbf{w}}_c, \boldsymbol{\mu}_c \rangle + (\widehat{\mathbf{b}}_j - \widehat{\mathbf{b}}_c) \right\} \\ &= \sum_{j \neq c} Q \left( \left\langle \frac{\widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_j}{\|\widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_j\|_{\ell_2}}, \boldsymbol{\mu}_c \right\rangle + \frac{\widehat{\mathbf{b}}_c - \widehat{\mathbf{b}}_j}{\|\widehat{\mathbf{w}}_c - \widehat{\mathbf{w}}_j\|_{\ell_2}} \right) = \sum_{j \neq c} Q \left( \frac{[\mathbf{t}_c]_j}{\sqrt{[\mathbf{S}_c]_{j,j}}} \right) \\ &\leq (k-1) \cdot Q(d_{\min}), \end{aligned} \quad (\text{C.7})$$

673 where in the last line, we denote  $d_{\min} := \min_{j \neq c} \left\{ \frac{[t_c]_j}{\sqrt{[S_c]_{j,j}}} \right\}$ . It is worth noting that the upper  
674 bound given in (C.7) still requires knowledge of the cross-correlations  $\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_c \rangle$ ,  $\ell \neq c$ .

675 **Oracle bound.** For completeness, we briefly discuss an oracle lower bound for the class-wise  
676 probability of error in GMM. Specifically, assume that the means  $\boldsymbol{\mu}_i$ ,  $i \in [n]$  are known. Then  
677 the problem of classifying a new sample  $\mathbf{x}$  is a  $k$ -ary hypothesis testing problem with Gaussian  
678 conditionals. Denote  $\mathbb{P}_{\text{genie, Bayes}}$  the Bayes error of this hypothesis testing problem. Clearly  
679  $\mathbb{P}_{\text{genie, Bayes}}$  is a lower bound on the error of any classifier that is trained on data. The paper [69]  
680 further lower bounds  $\mathbb{P}_{\text{genie, Bayes}}$  in terms of the Bayesian probability of errors between every two  
681 classes as follows:

$$\begin{aligned} \mathbb{P}_e &\geq \mathbb{P}_{\text{genie, Bayes}} \geq \frac{2}{k} \sum_{i \neq j} \pi_i \mathbb{P}_{\text{genie, ij}} & (\text{C.8}) \\ &= \frac{2}{k} \sum_{i \neq j} \frac{\pi_i}{\pi_i + \pi_j} \left\{ \pi_i \cdot Q \left( \frac{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|_{\ell_2}}{2} + \frac{\log(\pi_i/\pi_j)}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|_{\ell_2}} \right) + \pi_j \cdot Q \left( \frac{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|_{\ell_2}}{2} - \frac{\log(\pi_i/\pi_j)}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|_{\ell_2}} \right) \right\}. \end{aligned}$$

682 where  $\mathbb{P}_{\text{genie, ij}}$  is the Bayesian error between classes  $i$  and  $j$  with priors  $\frac{\pi_i}{\pi_i + \pi_j}$  and  $\frac{\pi_j}{\pi_i + \pi_j}$ . For the  
683 last equality we have used the well-known formula for the Bayesian probability of binary Gaussian  
684 hypothesis testing. In the case of equal-priors the genie lower bound above simplifies to:

$$\mathbb{P}_e \geq \frac{2}{k} \sum_{i \neq j} \pi_i \cdot Q \left( \frac{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|_{\ell_2}}{2} \right). \quad (\text{C.9})$$

685 Note that, in contrast to (C.9), our analysis allows for precise evaluations of the missclassification  
686 error  $\mathbb{P}_e$ .

## 687 D The Class-averaging estimator

### 688 D.1 Proofs for GMM

#### 689 D.1.1 GMM: Proof of Proposition 3.1

690 The first statement (3.1a) follows directly from the fact that  $\frac{1}{n} \mathbf{1}^T \mathbf{Y}_i = \frac{n_i}{n} \xrightarrow{P} \pi_i$ . For the next two  
691 statements note that

$$\widehat{\mathbf{w}}_i = \frac{1}{n} \mathbf{M} \mathbf{Y} \mathbf{Y}_i + \frac{1}{n} \mathbf{Z} \mathbf{Y}_i = \frac{1}{n} \sum_{j=1}^k \boldsymbol{\mu}_j (\mathbf{Y}_j^T \mathbf{Y}_i) + \frac{1}{n} \mathbf{Z} \mathbf{Y}_i = \frac{\|\mathbf{Y}_i\|_{\ell_2}^2}{n} \boldsymbol{\mu}_i + \frac{1}{n} \mathbf{Z} \mathbf{Y}_i, \quad (\text{D.1})$$

692 where in the last line we used orthogonality of the rows  $\mathbf{Y}_j$  of the matrix  $\mathbf{Y}$ :

$$\langle \mathbf{Y}_i, \mathbf{Y}_j \rangle = 0, \forall i \neq j \in [k]. \quad (\text{D.2})$$

693 To conclude simply use the facts that for all  $i \in [k]$ :

- 694 (i)  $\frac{\|\mathbf{Y}_i\|_{\ell_2}^2}{n} = \frac{n_i}{n} \xrightarrow{P} \pi_i$ .
- 695 (ii)  $\mathbf{Z} \mathbf{Y}_i \sim \sigma \|\mathbf{Y}_i\|_{\ell_2} \mathbf{g}_i$  with  $\mathbf{g}_i \stackrel{iid}{\sim} \mathcal{N}(0, I_d)$  because of (D.2)
- 696 (iii)  $\frac{\|\mathbf{g}_i\|_{\ell_2}}{\sqrt{n}} \xrightarrow{P} \sqrt{\gamma}$  and  $\frac{1}{\sqrt{n}} \langle \mathbf{g}_i, \boldsymbol{\mu}_j \rangle \xrightarrow{P} 0$ .

### 697 D.2 Proofs for MLM

#### 698 D.2.1 Proof of Proposition 4.1

699 Let us define  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)$  and random vector  $\mathbf{v} = [V_1, V_2, \dots, V_k]^T$  with entries:

$$\mathbf{v} = \frac{e^{\mathbf{V} \boldsymbol{\Sigma} \mathbf{g}}}{\mathbf{1}_k^T e^{\mathbf{V} \boldsymbol{\Sigma} \mathbf{g}}}, \quad V_i = \frac{e^{e_i^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{g}}}{\sum_{j \in [k]} e^{e_j^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{g}}}, \quad i \in [k]. \quad (\text{D.3})$$

700 We will prove the following three statements:

$$\widehat{\mathbf{b}} \xrightarrow{P} \mathbb{E}[\mathbf{v}] \quad (\text{D.4a})$$

$$\Sigma_{\mathbf{w}, \mu} \xrightarrow{P} \mathbb{E}[\mathbf{v}\mathbf{g}^T] \Sigma \mathbf{V}^T \quad (\text{D.4b})$$

$$\Sigma_{\mathbf{w}, \mathbf{w}} \xrightarrow{P} \gamma \cdot \text{diag}(\mathbb{E}[\mathbf{v}]) + \mathbb{E}[\mathbf{v}\mathbf{g}^T] \cdot \mathbb{E}[\mathbf{g}\mathbf{v}^T] \quad (\text{D.4c})$$

701 To see how these lead to (4.2), note that  $\mathbb{E}[\mathbf{v}] = \boldsymbol{\pi}$  by the definition in (4.1) and an application of  
 702 Gaussian integration by parts shows that  $\mathbb{E}[\mathbf{v}\mathbf{g}^T] = (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V} \Sigma$ . Therefore, in what follows,  
 703 we prove (D.4)

For the intercepts  $\widehat{\mathbf{b}}_\ell$ ,  $\ell \in [k]$  it holds that

$$\widehat{\mathbf{b}}_\ell = \frac{1}{n} \mathbf{1}^T \mathbf{Y}_\ell \xrightarrow{P} \mathbb{P}(Y = \ell) = \mathbb{E} \left[ \frac{e^{h_\ell}}{\sum_{j \in [k]} e^{h_j}} \right],$$

704 where  $\mathbf{h} \sim \mathcal{N}(\mathbf{0}_k, \Sigma_{\mu\mu})$ . To deduce the first statement in (D.4a), note that  $\mathbf{h} \stackrel{(D)}{=} \mathbf{V} \Sigma \mathbf{g}$ .

Continuing with the vectors  $\widehat{\mathbf{w}}_\ell$ ,  $\ell \in [k]$ , recall that  $\mathbf{w}_\ell = \frac{1}{n} \mathbf{X} \mathbf{Y}_\ell = \frac{1}{n} \sum_{i_\ell \in [n]} \mathbf{x}_{i_\ell} [\mathbf{Y}_\ell]_{i_\ell}$ . Consider the singular decomposition

$$\mathbf{M} = \mathbf{U} \Sigma \mathbf{V}^T = [\mathbf{u}_1 \quad \mathbf{u}_2 \quad \dots \quad \mathbf{u}_r] \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r) \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \dots \\ \mathbf{v}_r^T \end{bmatrix},$$

705 with  $\mathbf{U} \in \mathbb{R}^{d \times r}$ ,  $\Sigma \in \mathbb{R}^{r \times r}$ , and  $\mathbf{V} \in \mathbb{R}^{k \times r}$  where  $r = \text{rank}(\mathbf{M}) \leq k$ . Decompose  $\mathbf{X} \in \mathbb{R}^{d \times n}$  as  
 706  $\mathbf{X} = \mathbf{U} \mathbf{U}^T \mathbf{X} + \mathbf{P}^\perp \mathbf{X}$  with  $\mathbf{P}^\perp = \mathbf{I}_d - \mathbf{U} \mathbf{U}^T$ . With this notation we compute

$$\begin{aligned} \langle \mathbf{w}_\ell, \boldsymbol{\mu}_c \rangle &= \frac{1}{n} \sum_{i \in [n]} \mathbf{x}_i^T \boldsymbol{\mu}_c [\mathbf{Y}_\ell]_i = \frac{1}{n} \sum_{j=1}^r \sum_{i \in [n]} (\mathbf{x}_i^T \mathbf{u}_j) \cdot (\boldsymbol{\mu}_c^T \mathbf{u}_j) [\mathbf{Y}_\ell]_i + \frac{1}{n} \sum_{j=1}^r \sum_{i \in [n]} (\boldsymbol{\mu}_c^T \mathbf{P}^\perp \mathbf{x}_i) [\mathbf{Y}_\ell]_i \\ &\xrightarrow{P} \sum_{j=1}^r (e_c^T \mathbf{V} \Sigma e_j) \mathbb{E} \left[ \mathbf{g}_j \frac{(e^{e_j^T \mathbf{V} \Sigma \mathbf{g}})}{\sum_{\ell' \in [k]} e^{e_{\ell'}^T \mathbf{V} \Sigma \mathbf{g}}} \right] \\ &= \sum_{j=1}^r \mathbb{E}[V_\ell \mathbf{g}_j] (e_j^T \Sigma \mathbf{V}^T e_c). \end{aligned} \quad (\text{D.5})$$

707 Here, we have recognized that for every  $i \in [n] : \mathbf{U}^T \mathbf{x}_i \sim \mathbf{g}$ , and also, conditioned on  $\mathbf{x}_i : [\mathbf{Y}_\ell]_i \sim$   
 708  $\text{Bern}(e^{\boldsymbol{\mu}_\ell^T \mathbf{x}_i} / \sum_{\ell'} e^{\boldsymbol{\mu}_{\ell'}^T \mathbf{x}_i})$  and  $\boldsymbol{\mu}_\ell^T \mathbf{x}_i = e_\ell^T \mathbf{V} \Sigma \mathbf{U}^T \mathbf{x}_i \sim e_\ell^T \mathbf{V} \Sigma \mathbf{g}$ ,  $\ell \in [k]$ . This shows the second  
 709 statement in (D.4b) when expressed in matrix form.

710 We proceed similarly with the proof of the last statement in (4.2b) as follows:

$$\begin{aligned} \langle \mathbf{w}_\ell, \mathbf{w}_c \rangle &= \frac{1}{n^2} \sum_{i_\ell \in [n], i_c \in [n]} \mathbf{x}_{i_\ell}^T \mathbf{x}_{i_c} [\mathbf{Y}_\ell]_{i_\ell} [\mathbf{Y}_c]_{i_c} \\ &= \frac{1}{n^2} \sum_{j=1}^r \sum_{i_\ell \in [n], i_c \in [n]} (\mathbf{x}_{i_\ell}^T \mathbf{u}_j) (\mathbf{x}_{i_c}^T \mathbf{u}_j) [\mathbf{Y}_\ell]_{i_\ell} [\mathbf{Y}_c]_{i_c} + \frac{1}{n^2} \sum_{i_\ell \in [n], i_c \in [n]} (\mathbf{P}^\perp \mathbf{x}_{i_\ell})^T (\mathbf{P}^\perp \mathbf{x}_{i_c}) [\mathbf{Y}_\ell]_{i_\ell} [\mathbf{Y}_c]_{i_c} \end{aligned}$$

711 For  $i_\ell = i_c = i \in [n]$  note that

$$\frac{1}{n^2} \sum_{i \in [n]} \sum_{j=1}^r (\mathbf{x}_i^T \mathbf{u}_j) (\mathbf{x}_i^T \mathbf{u}_j) [\mathbf{Y}_\ell]_i [\mathbf{Y}_c]_i \xrightarrow{P} 0,$$

712 while, for  $i_\ell \neq i_c$ ,

$$\begin{aligned} \frac{1}{n^2} \sum_{i_\ell \neq i_c \in [n]} \sum_{j=1}^r (\mathbf{x}_{i_\ell}^T \mathbf{u}_j) [\mathbf{Y}_\ell]_{i_\ell} (\mathbf{x}_{i_c}^T \mathbf{u}_j) [\mathbf{Y}_c]_{i_c} &\xrightarrow{P} \sum_{j=1}^r \mathbb{E} \left[ \mathbf{g}_j \cdot \frac{(e^{e_\ell^T \mathbf{V} \Sigma \mathbf{g}})}{\sum_{\ell' \in [k]} e^{e_{\ell'}^T \mathbf{V} \Sigma \mathbf{g}}} \right] \mathbb{E} \left[ \mathbf{g}_j \cdot \frac{(e^{e_c^T \mathbf{V} \Sigma \mathbf{g}})}{\sum_{\ell' \in [k]} e^{e_{\ell'}^T \mathbf{V} \Sigma \mathbf{g}}} \right] \\ &= \sum_{j=1}^r \mathbb{E}[V_\ell \mathbf{g}_j] \mathbb{E}[\mathbf{g}_j V_c] = e_\ell^T \mathbf{E}[\mathbf{v}\mathbf{g}^T] \cdot \mathbf{E}[\mathbf{g}\mathbf{v}^T] e_c^T \end{aligned}$$

713 Furthermore,

$$\frac{1}{n} \sum_{i \in [n]} \|\mathbf{P}^\perp \mathbf{x}_i\|_{\ell_2}^2 [\mathbf{Y}_\ell]_i [\mathbf{Y}_c]_i^2 \xrightarrow{P} \gamma \cdot \mathbf{1}_{\ell,c} \cdot \mathbb{E} \left[ \left( \frac{e^{\mathbf{e}_\ell^T \mathbf{V} \Sigma \mathbf{g}}}{\sum_{\ell' \in [k]} e^{\mathbf{e}_{\ell'}^T \mathbf{V} \Sigma \mathbf{g}}} \right) \right] = \gamma \cdot \mathbf{e}_\ell^T \text{diag}(\mathbb{E}[\mathbf{v}]) \mathbf{e}_c.$$

714 Combining the last two displays results in (4.2b), as desired.

## 715 D.2.2 Orthogonal means

716 Here, we specialize the general result of Proposition 4.1 to the special case of orthogonal and means:  
717  $(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j) = 0, \forall i \neq j$ . Recall the notation  $\mu_i = \|\boldsymbol{\mu}_i\|_{\ell_2}, i \in [k]$ . Then, in this case the parameters in  
718 (4.1) are simply given by the following

$$\boldsymbol{\pi}_i := \mathbb{E} \left[ \frac{e^{\mu_i G_i}}{\sum_{\ell \in [k]} e^{\mu_\ell G_\ell}} \right], i \in [k] \quad \text{and} \quad \boldsymbol{\Pi}_{ij} := \mathbb{E} \left[ \frac{e^{\mu_i G_i} e^{\mu_j G_j}}{(\sum_{\ell \in [k]} e^{\mu_\ell G_\ell})^2} \right], i, j \in [k]. \quad (\text{D.6})$$

719 In this case (4.3) can be equivalently expressed as

$$\mathbb{P}_{e, \text{Avg}} \xrightarrow{P} \mathbb{P} \left( \arg \max_{\ell \in [k]} \{ \gamma \cdot \text{diag}(\boldsymbol{\pi}) \cdot \tilde{\mathbf{g}} + (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \cdot \Sigma \mathbf{g} \} \neq Y(\mathbf{g}) \right) \quad (\text{D.7})$$

$$= \mathbb{P} \left( \bigcup_{j \neq Y} \{ \gamma \cdot \boldsymbol{\pi}_\ell \cdot \tilde{\mathbf{g}}_\ell \geq \gamma \cdot \boldsymbol{\pi}_Y \cdot \tilde{\mathbf{g}}_Y + (\mathbf{e}_Y - \mathbf{e}_\ell)^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \Sigma \mathbf{g} + (\boldsymbol{\pi}_Y - \boldsymbol{\pi}_\ell) \} \right), \quad (\text{D.8})$$

720 where  $\mathbf{g}, \tilde{\mathbf{g}} \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$ ,  $\mathbb{P}(Y(\mathbf{g}) = c) = \frac{e^{\mu_c \boldsymbol{\theta} c}}{\sum_{\ell \in [k]} e^{\mu_\ell \boldsymbol{\theta} \ell}}$  and  $\Sigma = \text{diag}(\mu_1, \dots, \mu_k)$ .

## 721 E On the Bayes risk of GMM: Proof of Proposition 3.5

Without loss of generality in this proof we assume  $\sigma = 1$ . The general result follows by simply replacing  $(\mu, \sigma)$  with  $(\frac{\mu}{\sigma}, 1)$  and using the proof for  $\sigma = 1$ . Recall that the feature vectors  $\mathbf{x}_1, \dots, \mathbf{x}_n$  of the training data set are given by:

$$\mathbf{x}_i = \mathbf{M} \mathbf{y}_i + \mathbf{z}_i, \quad i \in [n],$$

722 where the matrix of means  $\mathbf{M} \in \mathbb{R}^{d \times k}$  has iid Gaussian entries with variance  $\mu^2/d$ ,  $\mathbf{z}_i \stackrel{iid}{\sim} \mathcal{N}(0, \mathbf{I}_d)$   
723 and  $\mathbf{y}_i \stackrel{iid}{\sim} \text{Unif}(\mathbf{e}_1, \dots, \mathbf{e}_k)$  with  $\mathbf{e}_j$  denoting the  $j^{\text{th}}$  canonical vector in  $\mathbb{R}^k$ . By definition here the  
724 Bayes estimator is the maximum-likelihood (ML) estimator. By applying the law of total probability  
725 and by successive application of the Bayes rule we have the following chain of reformulations of the  
726 ML:

$$\begin{aligned} \hat{\mathbf{y}}_{n+1} &= \arg \max_{e_j, j \in [k]} P(\mathbf{y} = e_j | \mathbf{X}, \mathbf{Y}, \mathbf{x}_{n+1}) \\ &= \arg \max_{e_j, j \in [k]} \int P(\mathbf{y} = e_j | \mathbf{M}, \mathbf{X}, \mathbf{Y}, \mathbf{x}_{n+1}) P(\mathbf{M} | \mathbf{X}, \mathbf{Y}, \mathbf{x}_{n+1}) d\mathbf{M} \\ &= \arg \max_{e_j, j \in [k]} \int \frac{P(\mathbf{x}_{n+1} | \mathbf{y} = e_j, \mathbf{M}, \mathbf{X}, \mathbf{Y}) \cdot P(\mathbf{y} = e_j | \mathbf{M}, \mathbf{X}, \mathbf{Y})}{P(\mathbf{x}_{n+1} | \mathbf{M}, \mathbf{X}, \mathbf{Y})} P(\mathbf{M} | \mathbf{X}, \mathbf{Y}, \mathbf{x}_{n+1}) d\mathbf{M} \\ &= \arg \max_{e_j, j \in [k]} \int P(\mathbf{x}_{n+1} | \mathbf{y} = e_j, \mathbf{M}) \frac{P(\mathbf{M} | \mathbf{X}, \mathbf{Y}, \mathbf{x}_{n+1})}{P(\mathbf{x}_{n+1} | \mathbf{M}, \mathbf{X}, \mathbf{Y})} d\mathbf{M} \quad (\text{E.1}) \end{aligned}$$

$$\begin{aligned} &= \arg \max_{e_j, j \in [k]} \int P(\mathbf{x}_{n+1} | \mathbf{y} = e_j, \mathbf{M}) \frac{P(\mathbf{M} | \mathbf{X}, \mathbf{Y})}{P(\mathbf{x}_{n+1} | \mathbf{X}, \mathbf{Y})} d\mathbf{M} \\ &= \arg \max_{e_j, j \in [k]} \int P(\mathbf{x}_{n+1} | \mathbf{y} = e_j, \mathbf{M}) P(\mathbf{M} | \mathbf{X}, \mathbf{Y}) d\mathbf{M} \quad (\text{E.2}) \end{aligned}$$

$$= \arg \max_{e_j, j \in [k]} \int P(\mathbf{x}_{n+1} | \mathbf{y} = e_j, \mathbf{M}) P(\mathbf{X} | \mathbf{M}, \mathbf{Y}) P(\mathbf{M}) d\mathbf{M} \quad (\text{E.3})$$

727 To arrive in (E.1) we used that  $P(\mathbf{y} = e_j | \mathbf{M}, \mathbf{X}, \mathbf{Y}) = \pi, \forall j \in [k]$  and  
728  $P(\mathbf{x}_{n+1} | \mathbf{y} = e_j, \mathbf{M}, \mathbf{X}, \mathbf{Y}) = P(\mathbf{x}_{n+1} | \mathbf{y} = e_j, \mathbf{M})$ . Also, (E.2) follows by recognizing



729 that  $P(\mathbf{x}_{n+1} \mid \mathbf{X}, \mathbf{Y}) > 0$  is independent of the variable of integration  $\mathbf{M}$  and of the optimization  
 730 variable  $j$ . For the same reasons, in (E.3) we have ignored the normalizing term  $P(\mathbf{X} \mid \mathbf{Y})$ .

731 Recalling that  $\mathbf{z}_{n+1} \sim \mathcal{N}(0, I_d)$ , we have that  $P(\mathbf{x}_{n+1} \mid \mathbf{y} = \mathbf{e}_j, \mathbf{M}) \propto \exp\left(-\|\mathbf{x}_{n+1} - \boldsymbol{\mu}_j\|_{\ell_2}^2 / 2\right)$   
 732 where  $\propto$  hides constant positive terms. Moreover, the posterior probability of the mean matrix given  
 733 the training data is given by

$$\begin{aligned} P(\mathbf{X} \mid \mathbf{M}, \mathbf{Y}) \cdot P(\mathbf{M}) &\propto \exp\left(-\frac{\|\mathbf{X} - \mathbf{M}\mathbf{Y}\|_{\ell_2}^2}{2}\right) \cdot \exp\left(-\frac{\|\mathbf{M}\|_{\ell_2}^2}{2(\mu^2/d)}\right) \\ &\propto \prod_{c=1}^k \left\{ \exp\left(-\frac{\|\boldsymbol{\mu}_c\|_{\ell_2}^2}{2(\mu^2/d)}\right) \cdot \prod_{i \in \mathcal{C}_c} \exp\left(-\frac{\|\mathbf{x}_i - \boldsymbol{\mu}_c\|_{\ell_2}^2}{2}\right) \right\}, \end{aligned} \quad (\text{E.4})$$

734 where we denote by  $\mathcal{C}_c$  the collection of training samples that belong to class  $c \in [k]$ , i.e.  $\mathcal{C}_c = \{i \in$   
 735  $[n] \mid \mathbf{y}_i = \mathbf{e}_c\}$ .

736 With these the objective function of the ML rule in (E.3) becomes:

$$\hat{\mathbf{y}}_{n+1} = \arg \max_{j \in [k]} \mathcal{I}(\ell, \mathcal{C}_j \cap n+1) \cdot \prod_{\substack{c=1 \\ c \neq j}}^k \mathcal{I}(c, \mathcal{C}_c) \quad (\text{E.5})$$

where for  $\ell \in [k]$  and a subset  $\mathcal{A} \subset [n+1]$  we denote

$$\mathcal{I}(\ell, \mathcal{A}) := \int d\boldsymbol{\mu}_c \exp\left(-\frac{\|\boldsymbol{\mu}_c\|_{\ell_2}^2}{2(\mu^2/d)}\right) \cdot \exp\left(-\sum_{i \in \mathcal{A}} \frac{\|\mathbf{x}_i - \boldsymbol{\mu}_c\|_{\ell_2}^2}{2}\right).$$

737 By completing the squares and invoking a gaussian integral it can be shown that

$$\begin{aligned} \mathcal{I}(\ell, \mathcal{A}) &:= \sqrt{\frac{(d/\mu^2 + |\mathcal{A}|)}{(2\varpi)^d}} \exp\left(-\frac{\left(1 - \frac{1}{d/\mu^2 + |\mathcal{A}|}\right) \sum_{i \in \mathcal{A}} \|\mathbf{x}_i\|_{\ell_2}^2 + \frac{1}{2(d/\mu^2 + |\mathcal{A}|)} \sum_{i \in \mathcal{A}} \langle \mathbf{x}_i, \sum_{\substack{j \in \mathcal{A} \\ j \neq i}} \mathbf{x}_j \rangle}{2}\right) \\ &:= \sqrt{\frac{(d/\mu^2 + |\mathcal{A}|)}{(2\varpi)^d}} \exp\left(-\frac{1}{2\left(\frac{d}{\mu^2} + |\mathcal{A}|\right)} \left( \left(\frac{d}{\mu^2} + |\mathcal{A}| - 1\right) \sum_{i \in \mathcal{A}} \|\mathbf{x}_i\|_{\ell_2}^2 - \sum_{i \in \mathcal{A}} \langle \mathbf{x}_i, \sum_{\substack{j \in \mathcal{A} \\ j \neq i}} \mathbf{x}_j \rangle \right)\right) \\ &:= \sqrt{\frac{(d/\mu^2 + |\mathcal{A}|)}{(2\varpi)^d}} \exp\left(-\frac{1}{2\left(\frac{d/n}{\mu^2} + |\mathcal{A}|/n\right)} \left( \left(\frac{d/n}{\mu^2} + \frac{|\mathcal{A}|}{n} - \frac{1}{n}\right) \sum_{i \in \mathcal{A}} \|\mathbf{x}_i\|_{\ell_2}^2 - \frac{1}{n} \sum_{i \in \mathcal{A}} \langle \mathbf{x}_i, \sum_{\substack{j \in \mathcal{A} \\ j \neq i}} \mathbf{x}_j \rangle \right)\right). \end{aligned}$$

738 Using this in (E.5) we have that

$$\hat{\mathbf{y}}_{n+1} = \arg \max_{j \in [k]} \mathcal{I}(j) \cdot \exp\left(-\frac{1}{2\left(\frac{d/n}{\mu^2} + \frac{n_j+1}{n}\right)} \left( \left(\frac{d/n}{\mu^2} + \frac{n_j}{n}\right) \|\mathbf{x}_{n+1}\|_{\ell_2}^2 - \frac{2}{n} \langle \mathbf{x}_{n+1}, \sum_{\ell \in \mathcal{C}_j} \mathbf{x}_\ell \rangle \right)\right), \quad (\text{E.6})$$

739 where  $\xi(n_c) := \frac{d/n}{\mu^2} + \frac{n_c}{n}$ ,  $c \in [k]$  and

$$\mathcal{I}(j) := \prod_{\substack{c=1 \\ c \neq j}}^k e^{-\frac{1}{2\xi(n_c)} \left( (\xi(n_c) - \frac{1}{n}) \sum_{i \in \mathcal{C}} \|\mathbf{x}_i\|_{\ell_2}^2 - \frac{1}{n} \sum_{i \in \mathcal{C}} \langle \mathbf{x}_i, \sum_{\substack{\ell \in \mathcal{C} \\ \ell \neq i}} \mathbf{x}_\ell \rangle \right)} \cdot e^{-\frac{1}{2(\xi(n_j) + \frac{1}{n})} \left( \xi(n_j) \sum_{i \in \mathcal{C}_j} \|\mathbf{x}_i\|_{\ell_2}^2 - \frac{1}{n} \sum_{i \in \mathcal{C}_j} \langle \mathbf{x}_i, \sum_{\substack{\ell \in \mathcal{C}_j \\ \ell \neq i}} \mathbf{x}_\ell \rangle \right)}.$$

740 Since the term  $\mathcal{I}$  is independent of  $j$  in the maximization in (E.6), we conclude that

$$\begin{aligned} \hat{\mathbf{y}}_{n+1} &= \arg \max_{j \in [k]} \log(\mathcal{I}(j)) - \frac{1}{2\left(\frac{d/n}{\mu^2} + \frac{n_j+1}{n}\right)} \left\{ \left(\frac{d/n}{\mu^2} + \frac{n_j}{n}\right) \|\mathbf{x}_{n+1}\|_{\ell_2}^2 - \frac{2}{n} \langle \mathbf{x}_{n+1}, \sum_{\ell \in \mathcal{C}_j} \mathbf{x}_\ell \rangle \right\} \\ &= \arg \max_{j \in [k]} \log(\mathcal{I}(j)) + \frac{1}{2\left(\frac{d/n}{\mu^2} + \frac{n_j+1}{n}\right)} \left\{ \frac{2}{n} \langle \mathbf{x}_{n+1}, \sum_{\ell \in \mathcal{C}_j} \mathbf{x}_\ell \rangle \right\} \end{aligned} \quad (\text{E.7})$$

741 Next, we evaluate the objective in (E.7) in the asymptotic limit  $n, d \rightarrow \infty, n/d = \gamma$ . First, since  
 742  $n_c/n \xrightarrow{P} \pi$ , note that  $\mathcal{I}(j) - \mathcal{I}(\ell) \xrightarrow{P} 0$  for all  $i, j \in [k]$ . Moreover, note that

$$\begin{aligned} \frac{1}{n} \langle \mathbf{x}_{n+1}, \sum_{\ell \in \mathcal{C}_j} \mathbf{x}_\ell \rangle &= \frac{1}{n} \langle \mathbf{M} \mathbf{y}_{n+1} + \mathbf{z}_{n+1}, n_j \boldsymbol{\mu}_j + \sum_{\ell \in \mathcal{C}_j} \mathbf{z}_\ell \rangle \\ &= \frac{n_j}{n} \langle \mathbf{M} \mathbf{y}_{n+1}, \boldsymbol{\mu}_j \rangle + \frac{n_j}{n} \langle \mathbf{z}_{n+1}, \boldsymbol{\mu}_j \rangle + \frac{1}{n} \sum_{\ell \in \mathcal{C}_j} \langle \mathbf{M} \mathbf{y}_{n+1}, \mathbf{z}_\ell \rangle + \frac{1}{n} \sum_{\ell \in \mathcal{C}_j} \langle \mathbf{z}_{n+1}, \mathbf{z}_\ell \rangle \end{aligned} \quad (\text{E.8})$$

743 For each one of the four terms in (E.8), we have the following by the CLT:

$$\begin{aligned} \frac{n_j}{n} \langle \mathbf{M} \mathbf{y}_{n+1}, \boldsymbol{\mu}_\ell \rangle &\xrightarrow{P} \pi \mu^2 \langle \mathbf{y}_{n+1}, \mathbf{e}_j \rangle \\ \frac{n_j}{n} \langle \mathbf{z}_{n+1}, \boldsymbol{\mu}_j \rangle &\xrightarrow{(D)} \mathcal{N}(0, \pi^2 r^2) \\ \frac{1}{n} \sum_{\ell \in \mathcal{C}_j} \langle \mathbf{M} \mathbf{y}_{n+1}, \mathbf{z}_\ell \rangle &\xrightarrow{P} 0 \\ \frac{1}{n} \sum_{\ell \in \mathcal{C}_j} \langle \mathbf{z}_{n+1}, \mathbf{z}_\ell \rangle &\xrightarrow{(D)} \mathcal{N}(0, \pi \gamma), \end{aligned}$$

744 where in the last line we used the fact that  $\frac{1}{\sqrt{n_j}} \sum_{\ell \in \mathcal{C}_j} \frac{\langle \mathbf{z}_{n+1}, \mathbf{z}_\ell \rangle}{\sqrt{n}} \xrightarrow{(D)} \mathcal{N}(0, \gamma)$ .

745 Therefore, in the asymptotic limit, the Bayes estimator is the solution to:

$$\hat{\mathbf{y}}_{n+1} = \arg \max_{\mathbf{e}_j, j \in [k]} \pi \mu^2 \langle \mathbf{y}_{n+1}, \mathbf{e}_j \rangle + \sqrt{\pi(\pi \mu^2 + \gamma)} G_j, \quad G_1, \dots, G_k \stackrel{iid}{\sim} \mathcal{N}(0, 1). \quad (\text{E.9})$$

746 As such, the probability of error is

$$\begin{aligned} \mathbb{P}_e &= \mathbb{P}(\hat{\mathbf{y}}_{n+1} \neq \mathbf{y}_{n+1}) = \mathbb{P}\left(\pi \mu^2 + \sqrt{\pi(\pi \mu^2 + \gamma)} G_0 \leq \max_{\ell \in [k-1]} \sqrt{\pi(\pi \mu^2 + \gamma)} G_\ell\right) \\ &= \mathbb{P}\left(G_0 + \max_{\ell \in [k-1]} G_\ell \geq \mu^2 \sqrt{\frac{\pi}{\pi \mu^2 + \gamma}}\right), \end{aligned} \quad (\text{E.10})$$

## 747 F Proof outline for (weighted) least-squares: key ideas and challenges

748 In this section, we provide a proof sketch for the analysis of the multiclass least-squares (LS) classifier.

749 Specifically, we discuss our approach towards specifying the high-dimensional limits of the key  
 750 quantities needed to evaluate the classification error:  $\hat{\mathbf{b}}$ ,  $\boldsymbol{\Sigma}_{w, \mu}$ , and  $\boldsymbol{\Sigma}_{w, w}$ . We note that our proofs  
 751 for the Weighted Least-Squares (WLS) classifiers follow the same general strategy, but in some parts  
 752 require significantly more involved and intricate analysis and derivations.

753 Our proofs follow the following general steps. See proof of each theorem for complete details and  
 754 derivations.

755 **Step I: Decomposing the loss across classes.** Recall from Section 2.2 that the multiclass LS  
 756 classifier produces a linear classifier  $\mathbf{x} \mapsto \mathbf{W} \mathbf{x} + \mathbf{b}$  via a least-squares fit to the training data:

$$(\widehat{\mathbf{W}}, \widehat{\mathbf{b}}) := \frac{1}{2n} \|\mathbf{W} \mathbf{X} + \mathbf{b} \mathbf{1}_n^T - \mathbf{Y}\|_F^2. \quad (\text{F.1})$$

757 Notice that the objective function above is separable. That is,

$$\frac{1}{2n} \|\mathbf{W} \mathbf{X} + \mathbf{b} \mathbf{1}_n^T - \mathbf{Y}\|_F^2 = \frac{1}{2n} \sum_{\ell=1}^k \|\mathbf{X}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2}^2.$$

758 Hence, for each  $\ell \in [k]$ ,

$$(\widehat{\mathbf{w}}_\ell, \widehat{b}_\ell) = \arg \min_{\mathbf{w}_\ell, b_\ell} \frac{1}{2n} \|\mathbf{X}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2}^2. \quad (\text{F.2})$$

759 This decomposition is convenient for analysis as it is easier to compute the statistical properties of  
 760 the simple single-output LS in (F.2) compared to the multi-output objective in (F.1). Indeed, as we  
 761 show, this simplification will eventually allow us to compute the high-dimensional behavior of the  
 762 intercepts  $\widehat{\mathbf{b}}_\ell$ ,  $\ell \in [k]$ , the mean-correlations  $\langle \widehat{\mathbf{w}}_\ell, \boldsymbol{\mu}_c \rangle$ ,  $\ell, c \in [k]$ , the norms  $\|\widehat{\mathbf{w}}_\ell\|_{\ell_2}$ ,  $\ell \in [k]$  and the  
 763 LS training loss  $\|\mathbf{X}^T \widehat{\mathbf{w}}_\ell + \widehat{\mathbf{b}}_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2}$ .

764 **Step II: Reduction to an Auxiliary Optimization (AO) problem via CGMT.** To calculate the  
 765 high-dimensional behavior (F.2) we use the Convex Gaussian min-max Theorem (CGMT) [60, 66]  
 766 framework. We provide a brief introduction of the CGMT machinery in Section F.1. Roughly stated  
 767 this framework allows us to replace a *Primary Optimization* (PO) problem of the form (F.2) with an  
 768 *Auxiliary Optimization* (AO) problem that is simpler to analyze, but is predictive of the behavior of  
 769 the latter. For instance for the PO in (F.2) in the GMM after some algebraic manipulations the AO  
 770 problem takes the form

$$\frac{1}{2} \left( \min_{\mathbf{w}_\ell, \mathbf{b}_\ell} \frac{1}{\sqrt{n}} \|\sigma \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + \mathbf{b}_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2} + \frac{1}{\sqrt{n}} \sigma \mathbf{h}^T \mathbf{w}_\ell \right)_+^2 \quad (\text{F.3})$$

771 where  $(x)_+ := \max(0, x)$  and  $\mathbf{g} \in \mathbb{R}^n$  and  $\mathbf{h} \in \mathbb{R}^d$  are two independent Gaussian random vectors  
 772 distributed as  $\mathcal{N}(\mathbf{0}, \mathbf{I}_n)$  and  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ .

773 **Step III: Simplification of the AO and computing  $\Sigma_{w, \mu}$  and  $\mathbf{b}$ .** In this step we carry out a series  
 774 of intricate calculations to further simplify (F.3) and characterize its various asymptotic properties.  
 775 In particular, in this step we compute the high-dimensional behavior of the intercepts  $\widehat{\mathbf{b}}_\ell$ ,  $\ell \in [k]$ ,  
 776 the mean-correlations  $\langle \widehat{\mathbf{w}}_\ell, \boldsymbol{\mu}_c \rangle$ ,  $\ell, c \in [k]$ , the norms  $\|\widehat{\mathbf{w}}_\ell\|_{\ell_2}$ ,  $\ell \in [k]$  and the LS training loss  
 777  $\|\mathbf{X}^T \widehat{\mathbf{w}}_\ell + \widehat{\mathbf{b}}_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2}$ . For instance, for GMM these calculations allow us to conclude that

$$\widehat{\mathbf{b}} \xrightarrow{P} \boldsymbol{\pi} - \mathbf{P} \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T \boldsymbol{\pi} \quad \text{and} \quad \boldsymbol{\Sigma}_{w, \mu} \xrightarrow{P} \mathbf{P} \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T,$$

778 where  $\mathbf{P} := \text{diag}(\boldsymbol{\pi}) - \boldsymbol{\pi} \boldsymbol{\pi}^T \geq \mathbf{0}_{k \times k}$  and  $\boldsymbol{\Delta} := \sigma^2 \mathbf{I}_r + \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{P} \mathbf{V} \boldsymbol{\Sigma} > \mathbf{0}_{r \times r}$ .

779 **Step IV: Computing  $\Sigma_{w, w}$  and capturing cross-correlations.** The final and most involved part  
 780 of our analysis is characterizing the asymptotic behavior of  $\Sigma_{w, w}$ . To see why this is particularly  
 781 challenging note that the reduction from (F.1) to (F.2) ‘‘breaks’’ the dependence of all  $\widehat{\mathbf{w}}_1, \widehat{\mathbf{w}}_2, \dots, \widehat{\mathbf{w}}_k$   
 782 on the *same* feature matrix  $\mathbf{X}$ . Capturing this dependence is crucial in determining the ‘‘cross-  
 783 correlations’’  $\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_c \rangle$ ,  $\ell \neq c$ . As noted in Section 2.3 the matrix  $\Sigma_{w, w}$  is needed to calculate the  
 784 class-wise and total miss-classification errors. Unfortunately, the CGMT is *not* directly applicable  
 785 to the multi-output LS optimization in (F.1). Our idea to circumvent this challenge builds on the  
 786 following simple observation: the vector  $\widehat{\mathbf{w}}_{\ell, c} = \widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_c$  is itself the solution to another simple  
 787 single-output LS problem.

788 **Lemma F.1** For  $\ell \neq c \in [k]$  let  $\widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_c$  be the  $\ell$  and  $c$ -th row of  $\widehat{\mathbf{W}}$  which is the solution to the  
 789 multi-output least-squares minimization (F.1). Denote  $\widehat{\mathbf{w}}_{\ell, c} := \widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_c$ . Then,  $\widehat{\mathbf{w}}_{\ell, c}$  is a minimizer in  
 790 the following single-output least-squares problem:

$$\widehat{\mathbf{w}}_{\ell, j} = \arg \min_{\mathbf{w}, \mathbf{b}} \frac{1}{2n} \|\mathbf{Y}_\ell + \mathbf{Y}_j - \mathbf{X}^T \mathbf{w} - \mathbf{b} \mathbf{1}_n\|_{\ell_2}^2. \quad (\text{F.4})$$

791 Thanks to Lemma I.1, we can use the CGMT to characterize the limiting behavior of  $\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_c\|_{\ell_2}$ .  
 792 These calculations are similar to (but, in certain cases, such as for weighted least-squares, more  
 793 involved than) those in Steps II and III above. Now note that an asymptotic characterization of  
 794  $\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_c\|_{\ell_2}$  immediately yields the asymptotic characterization of  $\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_c \rangle$  as

$$\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_c \rangle = \frac{\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_c\|_{\ell_2}^2 - \|\widehat{\mathbf{w}}_\ell\|_{\ell_2}^2 - \|\widehat{\mathbf{w}}_c\|_{\ell_2}^2}{2}, \quad (\text{F.5})$$

795 and  $\|\widehat{\mathbf{w}}_\ell\|_{\ell_2}, \|\widehat{\mathbf{w}}_c\|_{\ell_2}$  are already computed in Step IV. For instance, for GMM the analysis in this step  
 796 allow us to calculate the asymptotic behavior of  $\Sigma_{w, w}$  as follows

$$\boldsymbol{\Sigma}_{w, w} \xrightarrow{P} \frac{\gamma}{(1 - \gamma)\sigma^2} \mathbf{P} + \mathbf{P} \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\gamma}{(1 - \gamma)\sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{P}.$$

797 **F.1 Background on CGMT**

798 The CGMT is an extension of Gordon’s Gaussian min-max inequality (GMT) [22]. In the context of  
 799 high-dimensional inference problems, Gordon’s inequality was first successfully used in the study of  
 800 sharp phase-transitions in noiseless Compressed Sensing [59, 9, 2, 59]. More recently, [60] (see also  
 801 [2, Sec. 10.3]) discovered that Gordon’s inequality is essentially tight for certain convex problems. A  
 802 concrete and general formulation of this idea was given by [66] and was called the CGMT.

803 In order to summarize the essential ideas, consider the following two Gaussian processes:

$$X_{\mathbf{w}, \mathbf{u}} := \mathbf{u}^T \mathbf{G} \mathbf{w} + \psi(\mathbf{w}, \mathbf{u}), \quad (\text{F.6a})$$

$$Y_{\mathbf{w}, \mathbf{u}} := \|\mathbf{w}\|_{\ell_2} \mathbf{g}^T \mathbf{u} + \|\mathbf{u}\|_{\ell_2} \mathbf{h}^T \mathbf{w} + \psi(\mathbf{w}, \mathbf{u}), \quad (\text{F.6b})$$

804 where:  $\mathbf{G} \in \mathbb{R}^{n \times d}$ ,  $\mathbf{g} \in \mathbb{R}^n$ ,  $\mathbf{h} \in \mathbb{R}^d$ , they all have entries iid Gaussian; the sets  $\mathcal{S}_{\mathbf{w}} \subset \mathbb{R}^d$  and  $\mathcal{S}_{\mathbf{u}} \subset \mathbb{R}^n$   
 805 are compact; and,  $\psi : \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}$ . For these two processes, define the following (random) min-max  
 806 optimization programs, which are referred to as the *primary optimization* (PO) problem and the  
 807 *auxiliary optimization* AO:

$$\Phi(\mathbf{G}) = \min_{\mathbf{w} \in \mathcal{S}_{\mathbf{w}}} \max_{\mathbf{u} \in \mathcal{S}_{\mathbf{u}}} X_{\mathbf{w}, \mathbf{u}}, \quad (\text{F.7a})$$

$$\phi(\mathbf{g}, \mathbf{h}) = \min_{\mathbf{w} \in \mathcal{S}_{\mathbf{w}}} \max_{\mathbf{u} \in \mathcal{S}_{\mathbf{u}}} Y_{\mathbf{w}, \mathbf{u}}. \quad (\text{F.7b})$$

808 if the sets  $\mathcal{S}_{\mathbf{w}}$  and  $\mathcal{S}_{\mathbf{u}}$  are convex and *bounded*, and  $\psi$  is continuous *convex-concave* on  $\mathcal{S}_{\mathbf{w}} \times \mathcal{S}_{\mathbf{u}}$ ,  
 809 then, for any  $\nu \in \mathbb{R}$  and  $t > 0$ , it holds

$$\mathbb{P}(|\Phi(\mathbf{G}) - \nu| > t) \leq 2\mathbb{P}(|\phi(\mathbf{g}, \mathbf{h}) - \nu| > t). \quad (\text{F.8})$$

810 In words, concentration of the optimal cost of the AO problem around  $q^*$  implies concentration of  
 811 the optimal cost of the corresponding PO problem around the same value  $q^*$ . Asymptotically, if we  
 812 can show that  $\phi(\mathbf{g}, \mathbf{h}) \xrightarrow{P} q^*$ , then we can conclude that  $\Phi(\mathbf{G}) \xrightarrow{P} q^*$ . Moreover, starting from  
 813 (F.8) and under appropriate strict convexity conditions, the CGMT shows that concentration of the  
 814 optimal solution of the AO problem implies concentration of the optimal solution of the PO around  
 815 the same value. For example, if minimizers of (F.7b) satisfy  $\|\mathbf{w}_{\phi}(\mathbf{g}, \mathbf{h})\|_{\ell_2} \xrightarrow{P} \alpha^*$  for some  $\alpha^* > 0$ ,  
 816 then, the same holds true for the minimizers of (F.7a):  $\|\mathbf{w}_{\Phi}(\mathbf{G})\|_{\ell_2} \xrightarrow{P} \alpha^*$ . Thus, one can analyze  
 817 the AO to infer corresponding properties of the PO, the premise being of course that the former is  
 818 simpler to handle than the latter. In [50], the authors introduce a principled machinery that allows to  
 819 (a) express a quite general family of convex inference optimization problems in the form of the PO  
 820 and (b) simplify the AO from a (random) optimization over vector variables to an easier optimization  
 821 over only few scalar variables, termed the “scalarized AO”.

822 **G Least-squares for GMM**

823 **G.1 Proof of Theorem 3.2**

824 **G.1.1 Computing  $\Sigma_{\mathbf{w}, \mu}$**

$$\begin{aligned} \min_{\mathbf{W} \in \mathbb{R}^{k \times d}, \mathbf{b} \in \mathbb{R}^k} \frac{1}{2n} \|\mathbf{W} \mathbf{X} + \mathbf{b} \mathbf{1}_n^T - \mathbf{Y}\|_F^2 &= \sum_{\ell=1}^k \min_{\mathbf{w}_{\ell}, b_{\ell}} \frac{1}{2n} \|\mathbf{X}^T \mathbf{w}_{\ell} + b_{\ell} \mathbf{1}_n - \mathbf{Y}_{\ell}\|_{\ell_2}^2 \\ &= \sum_{\ell=1}^k \min_{\mathbf{w}_{\ell}, b_{\ell}} \frac{1}{2n} \|\mathbf{Y}^T \mathbf{M}^T \mathbf{w}_{\ell} + \mathbf{Z}^T \mathbf{w}_{\ell} + b_{\ell} \mathbf{1}_n - \mathbf{Y}_{\ell}\|_{\ell_2}^2 \end{aligned}$$

825 Define

$$\mathcal{L}_{PO}(\mathbf{w}_{\ell}, b_{\ell}) := \frac{1}{2n} \|\mathbf{Y}^T \mathbf{M}^T \mathbf{w}_{\ell} + \mathbf{Z}^T \mathbf{w}_{\ell} + b_{\ell} \mathbf{1}_n - \mathbf{Y}_{\ell}\|_{\ell_2}^2 \quad (\text{G.1})$$

826 **Identifying the AO.** To continue further note that by duality we have

$$\min_{\mathbf{w}_{\ell}, b_{\ell}} \mathcal{L}_{PO}(\mathbf{w}_{\ell}, b_{\ell}) = \min_{\mathbf{w}_{\ell}, b_{\ell}} \max_{\mathbf{s}} \frac{1}{n} \left( \mathbf{s}^T \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_{\ell} + \mathbf{s}^T \mathbf{Z}^T \mathbf{w}_{\ell} + b_{\ell} \mathbf{s}^T \mathbf{1}_n - \mathbf{s}^T \mathbf{Y}_{\ell} - \frac{\|\mathbf{s}\|_{\ell_2}^2}{2} \right)$$

827 Note that the above is jointly convex in  $(\mathbf{w}_\ell, b_\ell)$  and concave in  $\mathbf{s}$ . We consider the Auxiliary  
828 Optimization (AO) problem

$$\min_{\mathbf{w}_\ell, b_\ell} \max_{\mathbf{s}} \frac{1}{n} \left( \mathbf{s}^T \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + \sigma \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g}^T \mathbf{s} + \sigma \|\mathbf{s}\|_{\ell_2} \mathbf{h}^T \mathbf{w}_\ell + b_\ell \mathbf{s}^T \mathbf{1}_n - \mathbf{s}^T \mathbf{Y}_\ell - \frac{\|\mathbf{s}\|_{\ell_2}^2}{2} \right),$$

829 where  $\mathbf{g} \in \mathbb{R}^n$  and  $\mathbf{h} \in \mathbb{R}^d$  are independent Gaussian random vectors with i.i.d.  $\mathcal{N}(0, 1)$  entries.  
830 Maximizing over the direction of  $\mathbf{s}$  and its norm  $\beta = \|\mathbf{s}\|_{\ell_2}$  we arrive at

$$\begin{aligned} \min_{\mathbf{w}_\ell, b_\ell} \max_{\beta \geq 0} \frac{1}{n} & \left( \beta \|\sigma \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2} + \beta \sigma \mathbf{h}^T \mathbf{w}_\ell - \frac{\beta^2}{2} \right) \\ & = \min_{\mathbf{w}_\ell, b_\ell} \frac{1}{2n} \left( \|\sigma \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2} + \sigma \mathbf{h}^T \mathbf{w}_\ell \right)_+^2 \\ & = \frac{1}{2} \left( \min_{\mathbf{w}_\ell, b_\ell} \frac{1}{\sqrt{n}} \|\sigma \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2} + \frac{1}{\sqrt{n}} \sigma \mathbf{h}^T \mathbf{w}_\ell \right)_+^2 \end{aligned}$$

831 **Scalarization of the AO.** For convenience, define

$$\bar{\phi}_{AO, \ell} := \min_{\mathbf{w}_\ell, b_\ell} \frac{1}{\sqrt{n}} \|\sigma \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2} + \frac{\sigma}{\sqrt{n}} \mathbf{h}^T \mathbf{w}_\ell. \quad (\text{G.2})$$

832 To continue, consider the singular value decomposition

$$\mathbf{M} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = [\mathbf{u}_1 \quad \mathbf{u}_2 \quad \dots \quad \mathbf{u}_r] \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r) \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \dots \\ \mathbf{v}_r^T \end{bmatrix}, \quad (\text{G.3})$$

with  $\mathbf{U} \in \mathbb{R}^{d \times r}$ ,  $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$ , and  $\mathbf{V} \in \mathbb{R}^{k \times r}$  where  $r = \text{rank}(\mathbf{M}) \leq k$ . We further decompose  $\mathbf{w}_\ell$  in its projections on the orthogonal columns  $\mathbf{u}_1, \dots, \mathbf{u}_r$  of  $\mathbf{U}$ :

$$\mathbf{w}_\ell = \sum_{i=1}^r \alpha_i \mathbf{u}_i + \alpha_0 \mathbf{w}_\ell^\perp,$$

833 where  $\|\mathbf{w}_\ell^\perp\|_{\ell_2} = 1$  and  $\mathbf{U}^T \mathbf{w}_\ell^\perp = \mathbf{0}$ ,  $\alpha_0 \geq 0$  and we denote

$$\alpha_i := \mathbf{u}_i^T \mathbf{w}_\ell, i \in [r]. \quad (\text{G.4})$$

834 We also define  $\boldsymbol{\alpha} = [\alpha_1 \quad \alpha_2 \quad \dots \quad \alpha_k]^T$ . In this notation, we have

$$\begin{aligned} \bar{\phi}_{AO, \ell}(\mathbf{g}, \mathbf{h}) & := \min_{\alpha_0 \geq 0, \boldsymbol{\alpha} \in \mathbb{R}^r, b_\ell} \frac{1}{\sqrt{n}} \left\| \sigma \sqrt{\alpha_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2} \mathbf{g} + \mathbf{Y}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell \right\|_{\ell_2} \\ & \quad + \sum_{i=1}^r \alpha_i \sigma \frac{\mathbf{h}^T \mathbf{u}_i}{\sqrt{n}} + \frac{\alpha_0 \sigma}{\sqrt{n}} \min(\mathbf{h}^T \mathbf{w}_\ell^\perp) \\ & = \min_{\alpha_0 \geq 0, \boldsymbol{\alpha} \in \mathbb{R}^r, b_\ell} \frac{1}{\sqrt{n}} \left\| \sigma \sqrt{\alpha_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2} \mathbf{g} + \mathbf{Y}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell \right\|_{\ell_2} \\ & \quad + \sum_{i=1}^r \alpha_i \sigma \frac{\mathbf{h}^T \mathbf{u}_i}{\sqrt{n}} - \alpha_0 \sigma \frac{\|\mathbf{h}^\perp\|_{\ell_2}}{\sqrt{n}}, \end{aligned} \quad (\text{G.5})$$

835 where in the second line we denote  $\mathbf{h}^\perp$  the projection of  $\mathbf{h}$  onto the complement subspace of the span  
836 of  $\mathbf{u}_1, \dots, \mathbf{u}_r$  and we recalled that  $\|\mathbf{w}_\ell^\perp\|_{\ell_2} = 1$  and  $\langle \mathbf{w}_\ell^\perp, \mathbf{u}_i \rangle = 0, i \in [r]$ .

837 **Convergence of the AO.** First, note that

$$\begin{aligned} \frac{1}{n} \|\mathbf{Y}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2}^2 & = \frac{1}{n} \|\mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + b_\ell \mathbf{1}_n\|_{\ell_2}^2 \\ & = \frac{1}{n} \|\mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell)\|_{\ell_2}^2 + b_\ell^2 + \frac{2}{n} b_\ell \mathbf{1}_n^T \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) \\ & = \text{trace} \left( (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell)^T \text{diag} \left( \frac{n_1}{n}, \frac{n_2}{n}, \dots, \frac{n_k}{n} \right) (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) \right) \\ & \quad + b_\ell^2 + 2b_\ell \left[ \frac{n_1}{n} \quad \frac{n_2}{n} \quad \dots \quad \frac{n_k}{n} \right] (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) \end{aligned}$$

838 Thus

$$\begin{aligned} \frac{1}{n} \|\mathbf{Y}^T \mathbf{V} \Sigma \boldsymbol{\alpha} + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2}^2 &\xrightarrow{P} \text{trace} \left( (\mathbf{V} \Sigma \boldsymbol{\alpha} - \mathbf{e}_\ell)^T \text{diag}(\boldsymbol{\pi}) (\mathbf{V} \Sigma \boldsymbol{\alpha} - \mathbf{e}_\ell) \right) \\ &\quad + b_\ell^2 + 2b_\ell \boldsymbol{\pi}^T (\mathbf{V} \Sigma \boldsymbol{\alpha} - \mathbf{e}_\ell) \\ &= \boldsymbol{\alpha}^T (\Sigma \mathbf{V}^T \text{diag}(\boldsymbol{\pi}) \mathbf{V} \Sigma) \boldsymbol{\alpha} - 2\pi_\ell \boldsymbol{\alpha}^T \Sigma \mathbf{V}^T \mathbf{e}_\ell + 2b_\ell \boldsymbol{\alpha}^T \Sigma \mathbf{V}^T \boldsymbol{\pi} \\ &\quad + b_\ell^2 - 2b_\ell \pi_\ell + \pi_\ell \end{aligned}$$

839 At this point, observe that we have reduced the AO to an optimization problem over only  $r + 2$  scalar  
840 variables. Using the law of large numbers, the fact that  $\|\mathbf{h}^\perp\|_{\ell_2}$  concentrates around  $\sqrt{d-r}$  and  
841  $(d-r)/n \xrightarrow{P} \gamma$ , as well as the limit calculation above it is not hard to see that for fixed  $\alpha_0, b_\ell$  and  
842  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_r]^T \in \mathbb{R}^r$ , the objective function in (G.5) converges to the following:

$$\begin{aligned} \mathcal{D}_\ell(\alpha_0, \boldsymbol{\alpha}, b_\ell) \\ := \sqrt{\alpha_0^2 \sigma^2 + \boldsymbol{\alpha}^T (\sigma^2 \mathbf{I}_r + \Sigma \mathbf{V}^T \text{diag}(\boldsymbol{\pi}) \mathbf{V} \Sigma) \boldsymbol{\alpha} - 2\boldsymbol{\alpha}^T (\pi_\ell \Sigma \mathbf{V}^T \mathbf{e}_\ell - b_\ell \Sigma \mathbf{V}^T \boldsymbol{\pi}) + b_\ell^2 - 2b_\ell \pi_\ell + \pi_\ell} \\ - \alpha_0 \sigma \sqrt{\gamma}, \end{aligned} \quad (\text{G.6})$$

843 We will show in the next paragraph that the argument inside the square-root in (G.6) is a convex  
844 quadratic over  $(\alpha_0, \boldsymbol{\alpha}, b_\ell)$  (see (G.19)). Thus, the function  $\mathcal{D}_\ell(\alpha_0, \boldsymbol{\alpha}, b_\ell)$  is jointly convex. Using  
845 uniform convergence of convex functions over compact sets, we arrive at

$$\bar{\phi}_{AO,\ell}(\mathbf{g}, \mathbf{h}) \xrightarrow{P} \min_{\alpha_0 \geq 0, \boldsymbol{\alpha}, b_\ell} \mathcal{D}_\ell(\alpha_0, \boldsymbol{\alpha}, b_\ell). \quad (\text{G.7})$$

846 **Deterministic Analysis.** Here, we analyze the deterministic scalar minimization on the RHS of  
847 (G.7). Define

$$\mathbf{A} := \begin{bmatrix} \sigma^2 \mathbf{I}_r + \Sigma \mathbf{V}^T \text{diag}(\boldsymbol{\pi}) \mathbf{V} \Sigma & \Sigma \mathbf{V}^T \boldsymbol{\pi} \\ \boldsymbol{\pi}^T \mathbf{V} \Sigma & 1 \end{bmatrix} \quad \text{and} \quad \mathbf{c}_\ell = \begin{bmatrix} \Sigma \mathbf{V}^T \mathbf{e}_\ell \\ 1 \end{bmatrix} \quad (\text{G.8})$$

848 and observe that we can write

$$\mathcal{D}_\ell(\alpha_0, \boldsymbol{\alpha}, b_\ell) = \sqrt{\alpha_0^2 \sigma^2 + \pi_\ell + [\boldsymbol{\alpha}^T \quad b_\ell] \mathbf{A} \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - 2\pi_\ell \mathbf{c}_\ell^T \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - \alpha_0 \sigma \sqrt{\gamma}} \quad (\text{G.9})$$

849 First, note that the matrix  $\mathbf{A}$  is positive definite. This can be checked by computing the Schur  
850 complement of  $\mathbf{A}$ :

$$\boldsymbol{\Delta} := \sigma^2 \mathbf{I}_r + \Sigma \mathbf{V}^T \mathbf{P} \mathbf{V} \Sigma := \sigma^2 \mathbf{I}_r + \Sigma \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\pi} \boldsymbol{\pi}^T) \mathbf{V} \Sigma > \mathbf{0}_{r \times r}. \quad (\text{G.10})$$

851 Positive definiteness above holds because  $\mathbf{P} := (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\pi} \boldsymbol{\pi}^T) \geq \mathbf{0}_{k \times k}$ . Thus the term under the  
852 square-root in (G.9) is a strictly convex quadratic. Thus,  $\mathcal{D}_\ell$  is jointly convex in its arguments.

853 To simplify the RHS of (G.7) we proceed by minimizing  $\mathcal{D}_\ell(\alpha_0, \boldsymbol{\alpha}, b_\ell)$  over  $(\boldsymbol{\alpha}, b_\ell)$  which is equal  
854 to

$$\begin{aligned} \begin{bmatrix} \hat{\alpha} \\ \hat{b}_\ell \end{bmatrix} &= \pi_\ell \mathbf{A}^{-1} \mathbf{c}_\ell = \pi_\ell \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\boldsymbol{\pi}^T \mathbf{V} \Sigma & 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\Delta}^{-1} & \mathbf{0} \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{I} & -\Sigma \mathbf{V}^T \boldsymbol{\pi} \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} \Sigma \mathbf{V}^T \mathbf{e}_\ell \\ 1 \end{bmatrix} \\ &= \pi_\ell \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\boldsymbol{\pi}^T \mathbf{V} \Sigma & 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\Delta}^{-1} & \mathbf{0} \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} -\Sigma \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell) \\ 1 \end{bmatrix} \\ &= \pi_\ell \begin{bmatrix} -\boldsymbol{\Delta}^{-1} \Sigma \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell) \\ 1 + \boldsymbol{\pi}^T \mathbf{V} \Sigma \boldsymbol{\Delta}^{-1} \Sigma \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell) \end{bmatrix} \end{aligned} \quad (\text{G.11})$$

855 Thus, the minimum value attained is

$$-\pi_\ell^2 \begin{bmatrix} -(\boldsymbol{\pi} - \mathbf{e}_\ell)^T \mathbf{V} \Sigma & 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\Delta}^{-1} & \mathbf{0} \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} -\Sigma \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell) \\ 1 \end{bmatrix} = -\pi_\ell^2 \left( 1 + (\boldsymbol{\pi} - \mathbf{e}_\ell)^T \mathbf{V} \Sigma \boldsymbol{\Delta}^{-1} \Sigma \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell) \right).$$

856 Using the above (G.7) reduces to

$$\bar{\phi}_{AO,\ell}(\mathbf{g}, \mathbf{h}) \xrightarrow{P} \min_{\alpha_0 \geq 0} \sqrt{\alpha_0^2 \sigma^2 + \pi_\ell - \pi_\ell^2 \left( 1 + (\boldsymbol{\pi} - \mathbf{e}_\ell)^T \mathbf{V} \Sigma \boldsymbol{\Delta}^{-1} \Sigma \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell) \right)} - \alpha_0 \sigma \sqrt{\gamma}. \quad (\text{G.12})$$

857 Setting the derivative with respect to  $\alpha_0$  to zero we arrive at

$$\frac{\alpha_0 \sigma^2}{\sqrt{\alpha_0^2 \sigma^2 + \pi_\ell - \pi_\ell^2 \left(1 + (\boldsymbol{\pi} - \mathbf{e}_\ell)^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell)\right)}} = \sigma \sqrt{\gamma}.$$

858 Thus,

$$\widehat{\alpha}_0 = \frac{1}{\sigma} \sqrt{\frac{\gamma}{1-\gamma}} \sqrt{\pi_\ell (1 - \pi_\ell) - \pi_\ell^2 (\boldsymbol{\pi} - \mathbf{e}_\ell)^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell)}. \quad (\text{G.13})$$

859 Plugging the latter into (G.12) we arrive at

$$\bar{\phi}_{AO,\ell}(\mathbf{g}, \mathbf{h}) \xrightarrow{P} \sqrt{1-\gamma} \sqrt{\pi_\ell (1 - \pi_\ell) - \pi_\ell^2 (\boldsymbol{\pi} - \mathbf{e}_\ell)^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell)}$$

**Asymptotic predictions.** First, from (G.11) the bias term converges as follows:

$$\widehat{b}_\ell \xrightarrow{P} \pi_\ell \left(1 + \boldsymbol{\pi}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell)\right).$$

860 Thus,

$$\widehat{\mathbf{b}} \xrightarrow{P} \text{diag}(\boldsymbol{\pi}) \left(\mathbf{1}_k + (\boldsymbol{\pi} \mathbf{1}_k^T - \mathbf{I}_k) \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T \boldsymbol{\pi}\right)$$

861 Recall from (G.4) that  $\boldsymbol{\alpha} = \mathbf{U}^T \mathbf{w}_\ell$ . Thus, the correlations  $\langle \boldsymbol{\mu}_i, \mathbf{w}_\ell \rangle$ ,  $i \in [k]$  converge as follows:

$$\mathbf{M}^T \mathbf{w}_\ell = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{w}_\ell \xrightarrow{P} \mathbf{V} \boldsymbol{\Sigma} \widehat{\boldsymbol{\alpha}} = -\pi_\ell \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell). \quad (\text{G.14})$$

862 Here, convergence applies element-wise to the entries of the involved random vectors. Moreover, from  
863 the analysis above we can predict the limit of the norm  $\|\mathbf{w}_\ell\|_{\ell_2}$ . For this, note that  $\|\mathbf{w}_\ell\|_{\ell_2}^2 = \widehat{\alpha}_0^2 + \widehat{\boldsymbol{\alpha}}^T \widehat{\boldsymbol{\alpha}}$ .  
864 Thus,

$$\|\mathbf{w}_\ell\|_{\ell_2}^2 \xrightarrow{P} \frac{\gamma}{(1-\gamma)\sigma^2} \pi_\ell (1 - \pi_\ell) + \pi_\ell^2 (\boldsymbol{\pi} - \mathbf{e}_\ell)^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\gamma}{(1-\gamma)\sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_\ell) \quad (\text{G.15})$$

### 865 G.1.2 Computing $\Sigma_{w,w}$

866 In the previous section we used the CGMT to [predict the bias  $\widehat{b}_\ell$ , the correlations  $\langle \boldsymbol{\mu}_i, \widehat{\mathbf{w}}_\ell \rangle$ ,  $i \in [k]$   
867 and the norm  $\|\widehat{\mathbf{w}}_\ell\|_{\ell_2}$  for all  $\ell \in [k]$  members of the multi-output classifier. Here, we show how to  
868 compute the limits of the cross-correlations  $\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j \rangle$ ,  $\ell \neq j \in [k]$ .

**Lemma G.1** For  $\ell \neq j \in [k]$  let  $\widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j$  be solutions to the least-squares minimization

$$(\widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j, \widehat{b}_\ell, \widehat{b}_j) = \arg \min_{\mathbf{w}_\ell, \mathbf{w}_j, b_\ell, b_j} \left\{ \frac{1}{2n} \|\mathbf{Y}_\ell - \mathbf{X}^T \mathbf{w}_\ell - b_\ell \mathbf{1}_n\|_{\ell_2}^2 + \frac{1}{2n} \|\mathbf{Y}_j - \mathbf{X}^T \mathbf{w}_j - b_j \mathbf{1}_n\|_{\ell_2}^2 \right\}.$$

869 Denote  $\widehat{\mathbf{w}}_{\ell,j} := \widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j$  and  $\widehat{b}_{\ell,j} := \widehat{b}_\ell + \widehat{b}_j$ . Then,  $(\widehat{\mathbf{w}}_{\ell,j}, \widehat{b}_{\ell,j})$  is a minimizer in the following  
870 least-squares problem:

$$(\widehat{\mathbf{w}}_{\ell,j}, \widehat{b}_{\ell,j}) = \arg \min_{\mathbf{w}, b} \frac{1}{2n} \|\mathbf{Y}_\ell + \mathbf{Y}_j - \mathbf{X}^T \mathbf{w} - b \mathbf{1}_n\|_{\ell_2}^2 \quad (\text{G.16})$$

**Proof** Clearly the minimization in (G.16) is convex. Thus, it suffices to prove that  $\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j$  satisfies the KKT conditions. First, by optimality of  $\widehat{\mathbf{w}}_\ell$ , we have that

$$\mathbf{X} (\mathbf{Y}_\ell - \mathbf{X}^T \widehat{\mathbf{w}}_\ell - \widehat{b}_\ell \mathbf{1}_n) = 0$$

Similarly, for  $\widehat{\mathbf{w}}_j$ :

$$\mathbf{X} (\mathbf{Y}_j - \mathbf{X}^T \widehat{\mathbf{w}}_j - \widehat{b}_j \mathbf{1}_n) = 0.$$

Adding the equations on the above displays we find that

$$\mathbf{X} (\mathbf{Y}_\ell + \mathbf{Y}_j - \mathbf{X}^T (\widehat{\mathbf{w}}_j + \widehat{\mathbf{w}}_\ell) - (\widehat{b}_j + \widehat{b}_\ell) \mathbf{1}_n) = 0.$$

871 Recognize that this coincides with the optimality condition for (G.16). Thus, the proof is complete. ■

872 Thanks to Lemma I.1, we can use the CGMT to characterize the limiting behavior of  $\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j\|_{\ell_2}$ .  
 873 Observe that this immediately gives the limit of  $\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j \rangle$  since

$$\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j \rangle = \frac{\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j\|_{\ell_2}^2 - \|\widehat{\mathbf{w}}_\ell\|_{\ell_2}^2 - \|\widehat{\mathbf{w}}_j\|_{\ell_2}^2}{2}. \quad (\text{G.17})$$

874 The analysis of (G.16) is very similar to that of (I.1); thus, most details are omitted. Similar to (G.5)  
 875 we can relate (G.16) with the following AO problem:

$$\begin{aligned} \bar{\phi}_{AO,\ell,j}(\mathbf{g}, \mathbf{h}) := & \min_{\beta_0 \geq 0, \boldsymbol{\beta} \in \mathbb{R}^r, b_{\ell,j}} \frac{1}{\sqrt{n}} \left\| \sigma \sqrt{\beta_0^2 + \|\boldsymbol{\beta}\|_{\ell_2}^2} \mathbf{g} + \mathbf{Y}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\beta} + b_{\ell,j} \mathbf{1}_n - \mathbf{Y}_\ell - \mathbf{Y}_j \right\|_{\ell_2} \\ & + \sigma \sum_{i=1}^r \beta_i \frac{\mathbf{h}^T \mathbf{u}_i}{\sqrt{n}} - \sigma \beta_0 \frac{\|\mathbf{h}^\perp\|_{\ell_2}}{\sqrt{n}}, \end{aligned} \quad (\text{G.18})$$

where we have decomposed

$$\mathbf{w}_{\ell,j} = \sum_{i=1}^r \beta_i \mathbf{u}_i + \beta_0 \mathbf{w}_{\ell,j}^\perp,$$

876 with  $\|\mathbf{w}_{\ell,j}^\perp\|_{\ell_2} = 1$  and  $\mathbf{U}^T \mathbf{w}_{\ell,j}^\perp = \mathbf{0}_r$ .

877 Using a calculation similar to the one leading to (G.9) we can show that (G.18) converges point-wise  
 878 in  $\beta_0, \boldsymbol{\beta} = [\beta_1, \dots, \beta_r], b$  to the following:

$$\mathcal{D}_\ell(\beta_0, \boldsymbol{\beta}, b_{\ell,j}) = \sqrt{\beta_0^2 \sigma^2 + \pi_\ell + \pi_j + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A} \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2 \mathbf{d}_{\ell,j}^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - \beta_0 \sigma \sqrt{\gamma}}, \quad (\text{G.19})$$

where  $\mathbf{A}$  is as in (G.8) and we have further defined

$$\mathbf{d}_{\ell,j} := \begin{bmatrix} \pi_\ell \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{e}_\ell + \pi_j \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{e}_j \\ \pi_\ell + \pi_j \end{bmatrix}.$$

879 Thus, similar to (G.11) we can compute the minimizer of the deterministic

$$\begin{bmatrix} \widehat{\boldsymbol{\beta}} \\ \widehat{b}_{\ell,j} \end{bmatrix} = \begin{bmatrix} -\boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\pi_\ell (\boldsymbol{\pi} - \mathbf{e}_\ell) + \pi_j (\boldsymbol{\pi} - \mathbf{e}_j)) \\ \pi_\ell + \pi_j + \boldsymbol{\pi}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\pi_\ell (\boldsymbol{\pi} - \mathbf{e}_\ell) + \pi_j (\boldsymbol{\pi} - \mathbf{e}_j)) \end{bmatrix} \quad (\text{G.20})$$

880 and

$$\widehat{\beta}_0 = \frac{1}{\sigma} \sqrt{\frac{\gamma}{1-\gamma}} \sqrt{\pi_\ell + \pi_j - (\pi_\ell + \pi_j)^2 - (\pi_\ell (\boldsymbol{\pi} - \mathbf{e}_\ell) + \pi_j (\boldsymbol{\pi} - \mathbf{e}_j))^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\pi_\ell (\boldsymbol{\pi} - \mathbf{e}_\ell) + \pi_j (\boldsymbol{\pi} - \mathbf{e}_j))}, \quad (\text{G.21})$$

881 where recall that  $\boldsymbol{\Delta}$  is as in (G.10).

882 From the CGMT, we have that  $\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j\|_{\ell_2}^2 \xrightarrow{P} \widehat{\beta}_0^2 + \|\boldsymbol{\beta}\|_{\ell_2}^2$ . Combining this with the calculations  
 883 above, we conclude that

$$\begin{aligned} & \|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j\|_{\ell_2}^2 \xrightarrow{P} \\ & \frac{\gamma}{(1-\gamma)\sigma^2} (\pi_\ell + \pi_j) (1 - \pi_\ell - \pi_j) \\ & + (\pi_\ell (\boldsymbol{\pi} - \mathbf{e}_\ell) + \pi_j (\boldsymbol{\pi} - \mathbf{e}_j))^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\gamma}{(1-\gamma)\sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T (\pi_\ell (\boldsymbol{\pi} - \mathbf{e}_\ell) + \pi_j (\boldsymbol{\pi} - \mathbf{e}_j)) \end{aligned} \quad (\text{G.22})$$

884 Finally, using (G.22) and (G.15) in (G.17) it follows that

$$\langle \mathbf{w}_\ell, \mathbf{w}_j \rangle \xrightarrow{P} \pi_\ell \pi_j \left( -\frac{\gamma}{(1-\gamma)\sigma^2} + (\boldsymbol{\pi} - \mathbf{e}_\ell)^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\gamma}{(1-\gamma)\sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} - \mathbf{e}_j) \right). \quad (\text{G.23})$$



885 **G.2 Orthogonal means**

886 Here, we specialize the asymptotic predictions of Theorem 3.2 to the case of orthogonal means  
 887  $\langle \boldsymbol{\mu}_i, \boldsymbol{\mu}_j \rangle = 0, i \neq j$ .

**Corollary G.2 (Orthogonal means)** Consider the case of orthogonal means, i.e.  $\langle \boldsymbol{\mu}_i, \boldsymbol{\mu}_j \rangle = 0, \forall i \neq j$  and  $\gamma < 1$  with Euclidean norms given by  $\mu_i = \|\boldsymbol{\mu}_i\|_{\ell_2}$ . Define the following parameters for  $i \in [k]$ :

$$\rho_i := \pi_i \sigma^2 / (\sigma^2 + \pi_i \mu_i^2) \quad \text{and} \quad \beta_i = \rho_i \sigma^2 / \left( \sigma^2 - \sum_{i=1}^k \pi_i \rho_i \mu_i^2 \right).$$

888 Then, the following asymptotic limits hold for the least-squares classifier, for all  $i, j \in [k]$ :

$$\widehat{\mathbf{b}}_i \xrightarrow{P} \beta_i, \quad \langle \widehat{\mathbf{w}}_i, \boldsymbol{\mu}_j \rangle \xrightarrow{P} \frac{1}{\sigma} (\mathbb{1}_{ij} - \beta_i) \rho_j \mu_j, \quad (\text{G.24a})$$

$$\langle \widehat{\mathbf{w}}_i, \widehat{\mathbf{w}}_j \rangle \xrightarrow{P} \frac{1}{\sigma^4} \beta_i \beta_j \sum_{\ell=1}^k \rho_\ell^2 \mu_\ell^2 - \frac{1}{\sigma^4} \beta_i \rho_j^2 \mu_j^2 - \frac{1}{\sigma^4} \beta_j \rho_i^2 \mu_i^2 - \frac{\gamma \beta_i \rho_j}{(1-\gamma)\sigma^2} + \frac{\mathbb{1}_{ij}}{\sigma^2} \left( \frac{\gamma}{(1-\gamma)} \rho_i + \frac{1}{\sigma^2} \rho_i^2 \mu_i^2 \right) \quad (\text{G.24b})$$

889 Furthermore, if the means have equal norms  $\mu := \mu_i$  and the classes are balanced:  $\pi_i = 1/k, i \in [k]$ ,  
 890 then, setting  $u_{\text{LS}} := \frac{\mu^2}{\sigma} \sqrt{\frac{1-\gamma}{\mu^2 + k\gamma\sigma^2}}$ , it holds that

$$\mathbb{P}_e = \mathbb{P} \left\{ G_0 + \max_{j \in [k-1]} G_j \geq u_{\text{LS}} \right\}, \quad G_0, G_1, \dots, G_{k-1} \stackrel{iid}{\sim} \mathcal{N}(0, 1). \quad (\text{G.25})$$

891 **Proof** This is a direct corollary of Theorem 3.2. Indeed, (G.24) can be derived from (3.3a) after  
 892 substituting  $\mathbf{V} = \mathbf{I}_k, \boldsymbol{\Sigma} = \text{diag}(\mu_1, \mu_2, \dots, \mu_k)$  and some algebra steps that we omit for brevity.

Instead, we outline below how to conclude (G.26) from (G.24). Assume that  $\mu_i = \mu, \forall i \in [k]$  and  $\pi_i = \pi = 1/k, \forall i \in [k]$ . Recall from (2.9) that  $\mathbb{P}(\text{error} | \mathbf{y} = \mathbf{e}_c) = 1 - \mathbb{P}(\mathbf{S}_c^{1/2} \mathbf{z} > \mathbf{t})$ , and using (G.24) it can be checked that

$$\mathbf{S}_c = \frac{\pi}{1 + \pi \mu^2} \left( \frac{\pi \mu^2}{1 + \pi \mu^2} + \frac{\gamma}{1 - \gamma} \right) (\mathbf{I}_k + \mathbf{1}_k \mathbf{1}_k^T) \quad \text{and} \quad \mathbf{t} = -\frac{\pi \mu^2}{1 + \pi \mu^2} \mathbf{1}.$$

893 Thus, setting

$$u_{\text{LS}} := \mu^2 \sqrt{\frac{\pi}{\pi \mu^2 + \left( \frac{\gamma}{1-\gamma} \right) (1 + \pi \mu^2)}} = \mu^2 \sqrt{\frac{1-\gamma}{\mu^2 + \gamma/\pi}}, \quad (\text{G.26})$$

894 and applying Lemma (C.1), the probability of error is given by the advertised expression.  $\blacksquare$

895 **H Least-squares for MLM**

896 **H.1 Proof of Theorem 4.2**

897 **H.1.1 Computing  $\Sigma_{w, \mu}$**

898 Assume that  $\mathbf{X}, \mathbf{Y}$  are generated from the MLM.

899 Fix any  $\ell \in [k]$ . The classifier parameters  $\widehat{\mathbf{w}}_\ell, \widehat{b}_\ell$  minimize the following objective function  
 900  $\mathcal{L}_{PO}(\mathbf{w}_\ell, b_\ell) := \frac{1}{2n} \|\mathbf{X}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2}^2$ .

901 **Identifying the AO.** To continue further note that by duality we have

$$\min_{\mathbf{w}_\ell, b_\ell} \mathcal{L}_{PO}(\mathbf{w}_\ell, b_\ell) = \min_{\mathbf{w}_\ell, b_\ell} \max_{\mathbf{s}} \frac{1}{n} \left( \mathbf{s}^T \mathbf{X}^T \mathbf{w}_\ell + b_\ell \mathbf{s}^T \mathbf{1}_n - \mathbf{s}^T \mathbf{Y}_\ell - \frac{\|\mathbf{s}\|_{\ell_2}^2}{2} \right) \quad (\text{H.1})$$

902 and the optimization is jointly convex in  $(\mathbf{w}_\ell, b_\ell)$  and concave in  $\mathbf{s}$ . Here, note that  $\mathbf{Y}_\ell$  depends on the  
 903 Gaussian matrix  $\mathbf{X}$ . Thus, before applying the CGMT, we need to break this dependence as follows.  
 904 Consider the singular value decomposition

$$M = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = [\mathbf{u}_1 \quad \mathbf{u}_2 \quad \dots \quad \mathbf{u}_r] \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r) \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \dots \\ \mathbf{v}_r^T \end{bmatrix}, \quad (\text{H.2})$$

with  $\mathbf{U} \in \mathbb{R}^{d \times r}$ ,  $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$ , and  $\mathbf{V} \in \mathbb{R}^{k \times r}$  where  $r = \text{rank}(M) \leq k$ . For every  $i \in [n]$ , we decompose  $\mathbf{x}_i$  in its projection on the subspace spanned orthogonal columns  $\mathbf{u}_1, \dots, \mathbf{u}_r$  as follows:

$$\mathbf{x}_i = \mathbf{U}\mathbf{U}^T \mathbf{X}_i + \mathbf{P}^\perp \mathbf{X}_i = \mathbf{U}\tilde{\mathbf{g}}_i + \mathbf{P}^\perp \mathbf{x}_i,$$

905 where  $\mathbf{P}^\perp = \mathbf{I}_r - \mathbf{U}\mathbf{U}^T$ , and we denote

$$\tilde{\mathbf{G}} := [\tilde{\mathbf{g}}_1 \quad \tilde{\mathbf{g}}_2 \quad \dots \quad \tilde{\mathbf{g}}_n], \quad \tilde{\mathbf{g}}_i := \mathbf{U}^T \mathbf{x}_i \in \mathbb{R}^r, i \in [n]. \quad (\text{H.3})$$

906 Recalling that  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$  note that

$$\tilde{\mathbf{g}}_i \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r) \quad \text{and} \quad \tilde{\mathbf{g}}_i \perp \mathbf{P}^\perp \mathbf{x}_i. \quad (\text{H.4})$$

907 Further recall that for all  $i \in [n]$ , conditioned on  $\mathbf{x}_i$ :

$$[\mathbf{Y}_\ell]_i \sim \text{Bern}\left(\frac{e^{\boldsymbol{\mu}_\ell^T \mathbf{x}_i}}{\sum_{\ell' \in [k]} e^{\boldsymbol{\mu}_{\ell'}^T \mathbf{x}_i}}\right) \sim \text{Bern}(V_\ell), \quad (\text{H.5})$$

908 where we used (H.3) and the SVD decomposition of  $M$ . In this notation, we can rewrite the PO as  
 909 follows:

$$\min_{\mathbf{w}_\ell, b_\ell} \max_{\mathbf{s}} \frac{1}{n} \left( \mathbf{s}^T \mathbf{X}^T \mathbf{P}^\perp \mathbf{w}_\ell + \mathbf{s}^T \tilde{\mathbf{G}}^T \mathbf{U}^T \mathbf{w}_\ell + b_\ell \mathbf{s}^T \mathbf{1}_n - \mathbf{s}^T \mathbf{Y}_\ell - \frac{\|\mathbf{s}\|_{\ell_2}^2}{2} \right)$$

910 From (H.4) and (H.5) notice that  $\mathbf{Y}_\ell$  depends only on  $\tilde{\mathbf{G}}$  and  $\tilde{\mathbf{G}}$  is independent of  $\mathbf{X}^T \mathbf{P}^\perp$ . Therefore,  
 911 the corresponding Auxiliary Optimization (AO) problem becomes

$$\min_{\mathbf{w}_\ell, b_\ell} \max_{\mathbf{s}} \frac{1}{n} \left( \|\mathbf{P}^\perp \mathbf{w}_\ell\|_{\ell_2} \mathbf{g}^T \mathbf{s} + \|\mathbf{s}\|_{\ell_2} \mathbf{h}^T \mathbf{P}^\perp \mathbf{w}_\ell + \mathbf{s}^T \tilde{\mathbf{G}}^T \mathbf{U}^T \mathbf{w}_\ell + b_\ell \mathbf{s}^T \mathbf{1}_n - \mathbf{s}^T \mathbf{Y}_\ell - \frac{\|\mathbf{s}\|_{\ell_2}^2}{2} \right), \quad (\text{H.6})$$

912 where  $\mathbf{g} \in \mathbb{R}^n$  and  $\mathbf{h} \in \mathbb{R}^d$  are iid Gaussian vectors independent of everything else.

913 **Scalarization of the AO.** Maximizing over the direction of  $\mathbf{s}$  and denoting its norm  $\beta = \|\mathbf{s}\|_{\ell_2} \geq 0$   
 914 we arrive at

$$\begin{aligned} \min_{\mathbf{w}_\ell, b_\ell} \max_{\beta \geq 0} \frac{1}{n} & \left( \beta \left( \|\mathbf{P}^\perp \mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + \mathbf{G}^T \mathbf{U}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell \right)_{\ell_2} + \beta \mathbf{h}^T \mathbf{P}^\perp \mathbf{w}_\ell - \frac{\beta^2}{2} \right) \\ & = \min_{\mathbf{w}_\ell, b_\ell} \frac{1}{2n} \left( \|\mathbf{P}^\perp \mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + \tilde{\mathbf{G}}^T \mathbf{U}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell \right)_{\ell_2}^2 + \mathbf{h}^T \mathbf{P}^\perp \mathbf{w}_\ell \\ & = \frac{1}{2} \left( \min_{\mathbf{w}_\ell, b_\ell} \frac{1}{\sqrt{n}} \left\| \|\mathbf{P}^\perp \mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + \tilde{\mathbf{G}}^T \mathbf{U}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell \right\|_{\ell_2} + \frac{1}{\sqrt{n}} \mathbf{h}^T \mathbf{P}^\perp \mathbf{w}_\ell \right)_+^2 \end{aligned} \quad (\text{H.7})$$

In the remaining, we focus in the inner minimization above. Let us denote

$$\mathbf{a} := \mathbf{U}^T \mathbf{w}_\ell \quad \text{and} \quad \alpha_0 = \|\mathbf{P}^\perp \mathbf{w}_\ell\|_{\ell_2}.$$

915 Notice that  $\mathbf{a} \perp \mathbf{P}^\perp \mathbf{w}_\ell$  and thus the orthogonal decomposition  $\mathbf{w}_\ell = \mathbf{U}\mathbf{a} + \mathbf{P}^\perp \mathbf{w}_\ell$ . With this  
 916 observation, we can optimize over the direction of  $\mathbf{P}^T \mathbf{w}_\ell$  in (H.7) by aligning it with  $\mathbf{P}^T \mathbf{h}$ . With  
 917 this, the minimization in (H.7) reduces to the following

$$\min_{\mathbf{a}, \alpha_0 \geq 0, b_\ell} \frac{1}{\sqrt{n}} \left\| \alpha_0 \mathbf{g} + \tilde{\mathbf{G}}^T \mathbf{a} + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell \right\|_{\ell_2} - \alpha_0 \frac{1}{\sqrt{n}} \|\mathbf{P}^\perp \mathbf{h}\|_{\ell_2}. \quad (\text{H.8})$$

918 **Convergence of the AO.** First, we argue on point-wise convergence of the objective function in  
 919 (H.8). Fix  $\mathbf{a}$ ,  $\alpha_0$  and  $\mathbf{b}_\ell$ . From the WLLN,  $\frac{1}{\sqrt{n}} \|\mathbf{P}^\perp \mathbf{h}\|_{\ell_2} \xrightarrow{P} \sqrt{\gamma}$  and

$$\frac{1}{n} \|\alpha_0 \mathbf{g} + \tilde{\mathbf{G}}^T \mathbf{a} + \mathbf{b}_\ell \mathbf{1}_n - \mathbf{Y}_\ell\|_{\ell_2}^2 = \frac{1}{n} \sum_{i=1}^n (\alpha_0 \mathbf{g}_i + \mathbf{a}^T \tilde{\mathbf{g}}_i + \mathbf{b}_\ell - [\mathbf{Y}_\ell]_i)^2 \xrightarrow{P} \mathbb{E} \left[ (\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2 \right], \quad (\text{H.9})$$

920 where the expectation is over  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)$  (with some abuse of notation) and

$$Y_\ell \sim \text{Bern}(V_\ell) \quad \text{and} \quad V_\ell = \frac{e^{\mathbf{e}_\ell^T \mathbf{V} \Sigma \mathbf{g}}}{\sum_{\ell'=1}^k e^{\mathbf{e}_{\ell'}^T \mathbf{V} \Sigma \mathbf{g}}}. \quad (\text{H.10})$$

921 Therefore, point-wise on  $\mathbf{a}$ ,  $\alpha_0$  and  $\mathbf{b}_\ell$ , the objective of the AO converges to

$$\mathcal{D}_\ell(\alpha_0, \alpha, \mathbf{b}_\ell) := \sqrt{\mathbb{E} \left[ (\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2 \right]} - \alpha_0 \sqrt{\gamma}. \quad (\text{H.11})$$

922 Next, with an argument based on convexity and compactness similar to that in ‘‘Convergence analysis  
 923 of the AO’’ in Section G it can be argued that the convergence above is uniform. Thus,

$$(\text{H.8}) \xrightarrow{P} \min_{\alpha_0 \geq 0, \alpha, \mathbf{b}_\ell} \mathcal{D}_\ell(\alpha_0, \alpha, \mathbf{b}_\ell). \quad (\text{H.12})$$

**Deterministic analysis of the AO.** Here, we solve the deterministic minimization problem in  
 (H.12). Optimization over  $\mathbf{b}_\ell$  is straightforward. By setting

$$\tilde{\mathbf{b}}_\ell = \mathbb{E}[Y_\ell] = \mathbb{E}[V_\ell],$$

924 we now have to optimize

$$\min_{\alpha_0 \geq 0, \alpha} \sqrt{\alpha_0^2 + \mathbb{E} \left[ (\mathbf{a}^T \mathbf{g} - Y_\ell)^2 \right]} - (\mathbb{E}[V_\ell])^2 - \alpha_0 \sqrt{\gamma}. \quad (\text{H.13})$$

925 By direct differentiation and first-order optimality, we compute the optimal values as follows:

$$\tilde{\alpha}_j = \mathbb{E}[\mathbf{g}_j Y_\ell] = \mathbb{E}[\mathbf{g}_j V_\ell], \quad j \in [r] \quad (\text{H.14})$$

$$\tilde{\alpha}_0^2 = \frac{\gamma}{1-\gamma} \left( \text{Var}[Y_\ell] - \sum_{j=1}^r (\mathbb{E}[\mathbf{g}_j V_\ell])^2 \right) = \frac{\gamma}{1-\gamma} \left( \mathbb{E}[V_\ell] - (\mathbb{E}[V_\ell])^2 - \sum_{j=1}^r (\mathbb{E}[\mathbf{g}_j V_\ell])^2 \right) \quad (\text{H.15})$$

926 **Asymptotic Predictions.** From the analysis above, we conclude with the following limits about the  
 927 solution  $\hat{\mathbf{b}}_\ell, \hat{\mathbf{w}}_\ell$  of the PO:

$$\hat{\mathbf{b}}_\ell \xrightarrow{P} \mathbb{E}[V_\ell] \quad (\text{H.16a})$$

$$\langle \boldsymbol{\mu}_c, \hat{\mathbf{w}}_\ell \rangle \xrightarrow{P} \mathbf{e}_c^T \mathbf{V} \Sigma \mathbb{E}[\mathbf{g} V_\ell], \quad c \in [k] \quad (\text{H.16b})$$

$$\|\hat{\mathbf{w}}_\ell\|_{\ell_2}^2 \xrightarrow{P} \sum_{j=1}^r (\mathbb{E}[\mathbf{g}_j V_\ell])^2 + \frac{\gamma}{1-\gamma} \left( \mathbb{E}[V_\ell] - (\mathbb{E}[V_\ell])^2 - \sum_{j=1}^r (\mathbb{E}[\mathbf{g}_j V_\ell])^2 \right) \quad (\text{H.16c})$$

$$= \frac{\gamma}{1-\gamma} \left( \mathbb{E}[V_\ell] - (\mathbb{E}[V_\ell])^2 \right) + \frac{1-2\gamma}{1-\gamma} \sum_{j=1}^r (\mathbb{E}[\mathbf{g}_j V_\ell])^2 \quad (\text{H.16d})$$

928 Recall the notation in (4.1). Note that  $\mathbb{E}[V_\ell] = \boldsymbol{\pi}_\ell$ . Moreover, using Gaussian integration by parts  
 929 Lemma B.2, it can be shown that  $\mathbb{E}[V_\ell \mathbf{g}] = \Sigma \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \mathbf{e}_\ell$ . Using these and writing in  
 930 matrix form, we arrive at (4.4a)

### 931 H.1.2 Computing $\Sigma_{w,w}$

932 Here, we prove (4.4b). Specifically, we compute the correlations  $\langle \hat{\mathbf{w}}_\ell, \hat{\mathbf{w}}_c \rangle$ ,  $\ell \neq c \in [k]$  by following  
 933 the strategy of Section G.1.2. Specifically, in view of Lemma I.1 we need to study the following PO:

$$\min_{\mathbf{w}, \mathbf{b}} \max_{\mathbf{s}} \frac{1}{n} \left( \mathbf{s}^T \mathbf{X}^T \mathbf{w} + \mathbf{b} \mathbf{s}^T \mathbf{1}_n - \mathbf{s}^T (\mathbf{Y}_\ell + \mathbf{Y}_c) - \frac{\|\mathbf{s}\|_{\ell_2}^2}{2} \right) \quad (\text{H.17})$$

934 which is minimized by  $\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_c$ . Thus the analysis will lead us to an asymptotic formula for  
 935  $\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_c\|_{\ell_2}$ . This when combined with the formulae for  $\|\widehat{\mathbf{w}}_\ell\|_{\ell_2}$  and  $\|\widehat{\mathbf{w}}_c\|_{\ell_2}$  in (H.16d) will give  
 936 the desired.

937 The analysis of (H.17) is almost identical to the analysis of (H.1) in the previous section. Specifically,  
 938 without repeating all the details for brevity, it can be shown that the AO of (H.17) converges to the  
 939 following (cf. (H.11)):

$$\mathcal{D}_\ell(\alpha_0, \boldsymbol{\alpha}, b_\ell) := \sqrt{\mathbb{E}\left[(\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_{\ell,c})^2\right]} - \alpha_0 \sqrt{\gamma}, \quad (\text{H.18})$$

940 where as before  $G_0 \sim \mathcal{N}(0, 1)$ ,  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)$ , only now (H.10) is modified to:

$$Y_{\ell,c} \sim \text{Bern}(V_c + V_\ell) \quad \text{and as before:} \quad V_\ell = \frac{e^{e_\ell^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{g}}}{\sum_{\ell'=1}^k e^{e_{\ell'}^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{g}}}. \quad (\text{H.19})$$

941 With these, it can be shown that

$$\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_c\|_{\ell_2}^2 \xrightarrow{P} \sum_{j=1}^r (\mathbb{E}[\mathbf{g}_j(V_c + V_\ell)])^2 + \frac{\gamma}{1-\gamma} \left( \mathbb{E}[V_c + V_\ell] - (\mathbb{E}[V_c + V_\ell])^2 - \sum_{j=1}^r (\mathbb{E}[\mathbf{g}_j(V_c + V_\ell)])^2 \right)$$

942 Combining this with (H.16d), we conclude that for  $\ell \neq c \in [k]$ :

$$\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_c \rangle \xrightarrow{P} \frac{1-2\gamma}{1-\gamma} \sum_{j=1}^r \mathbb{E}[\mathbf{g}_j V_c] \mathbb{E}[\mathbf{g}_j V_\ell] - \frac{\gamma}{1-\gamma} \mathbb{E}[V_c] \mathbb{E}[V_\ell]. \quad (\text{H.20})$$

943 This shows (4.4b) after applying Gaussian integration by parts and expressing it in matrix form.

## 944 H.2 Orthogonal means and equal-energy

945 Here, we use Theorem 4.2 to prove that, in contrast to the GMM, in the MLM under orthogonal and  
 946 equal-energy means: LS outperforms the averaging classifier for large enough sample sizes.

$$\boldsymbol{\pi} = \boldsymbol{\pi}_1 \mathbf{1}_k = (1/k) \mathbf{1}_k, \quad (\text{H.21})$$

$$\boldsymbol{\Pi} = (\boldsymbol{\Pi}_{11} - \boldsymbol{\Pi}_{12}) \mathbf{I}_k + \boldsymbol{\Pi}_{12} \mathbf{1}_k \mathbf{1}_k^T \quad \text{with} \quad \boldsymbol{\Pi}_{12} = \frac{1-k^2 \boldsymbol{\Pi}_{11}^2}{k(k-1)} \quad \text{and} \quad \boldsymbol{\Pi}_{11} = \mathbb{E} \left[ \frac{e^{2\mu G_1}}{(\sum_{\ell \in [k]} e^{\mu G_\ell})^2} \right].$$

947 Then,

$$\boldsymbol{\Sigma}_{\mathbf{w}, \mathbf{w}} - \boldsymbol{\Sigma}_{\mathbf{w}, \boldsymbol{\mu}} \boldsymbol{\Sigma}_{\boldsymbol{\mu}, \boldsymbol{\mu}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\mu}, \mathbf{w}}^T \xrightarrow{P} \frac{\gamma}{1-\gamma} \cdot (p \mathbf{I}_k - q \mathbf{1}_k \mathbf{1}_k^T), \quad (\text{H.22})$$

948 where we defined

$$p := \boldsymbol{\pi}_1 - \mu^2 (\boldsymbol{\pi}_1 - \boldsymbol{\Pi}_{11} + \boldsymbol{\Pi}_{12})^2 \quad \text{and} \quad q := (\boldsymbol{\pi}_1^2 + \boldsymbol{\Pi}_{12}^2 \mu^2 k - 2\mu^2 \boldsymbol{\Pi}_{12} (\boldsymbol{\pi}_1 - \boldsymbol{\Pi}_{11} + \boldsymbol{\Pi}_{12})). \quad (\text{H.23})$$

949 Thus, similar to (4.3) and with the same notation,

$$\mathbb{P}_{e, \text{LS}} \xrightarrow{P} \mathbb{P} \left\{ \arg \max_{\ell \in [k]} \left\{ \sqrt{\frac{\gamma}{1-\gamma}} \cdot (p \mathbf{I}_k - q \mathbf{1}_k \mathbf{1}_k^T)^{1/2} \cdot \tilde{\mathbf{g}} + \mu \cdot (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{g} \right\} \neq Y(\mathbf{g}) \right\}. \quad (\text{H.24})$$

950 In (H.24) (as well as in (4.3)), note that the matrices multiplying  $\tilde{\mathbf{g}}$  and  $\mathbf{g}$  have all the form of a rank  
 951 one update of a (scaled) identity matrix. It turns out that we can exploit this structure to simplify the  
 952 formulae for the test error even further. Importantly, this lets us directly compare  $\mathbb{P}_{e, \text{LS}}$  and  $\mathbb{P}_{e, \text{Avg}}$  of  
 953 the two classifiers. These are detailed in Section H.3.

## 954 H.3 Proof of Proposition 4.3

955 In (D.7) and (H.24), we showed the following limits for orthogonal means of equal-energy  $\mu > 0$ :

$$\begin{aligned} \mathbb{P}_{e, \text{Avg}} &\xrightarrow{P} \mathbb{P} \left( \arg \max_{\ell \in [k]} \left\{ \gamma \cdot \boldsymbol{\pi}_1 \cdot \mathbf{I}_k \cdot \tilde{\mathbf{g}} + \mu \left( (\boldsymbol{\pi}_1 - \boldsymbol{\Pi}_{11}) \cdot \mathbf{I}_k + \boldsymbol{\Pi}_{12} \mathbf{1}_k \mathbf{1}_k^T \right) \cdot \mathbf{g} \right\} \neq Y(\mathbf{g}) \right) \\ \mathbb{P}_{e, \text{LS}} &\xrightarrow{P} \mathbb{P} \left( \arg \max_{\ell \in [k]} \left\{ \sqrt{\frac{\gamma}{1-\gamma}} \cdot (p \mathbf{I}_k - q \mathbf{1}_k \mathbf{1}_k^T)^{1/2} \cdot \tilde{\mathbf{g}} + \mu \left( (\boldsymbol{\pi}_1 - \boldsymbol{\Pi}_{11}) \cdot \mathbf{I}_k + \boldsymbol{\Pi}_{12} \mathbf{1}_k \mathbf{1}_k^T \right) \cdot \mathbf{g} \right\} \neq Y(\mathbf{g}) \right), \end{aligned}$$

956 where  $\mathbb{P}(Y(\mathbf{g}) = \ell) = \frac{e^{\mu \mathbf{g} \ell}}{\sum_{j \in [k]} e^{\mu \mathbf{g} j}}$  and we have further used (H.21) and the notation in (H.23).

957 We compare the expression on the RHS in the above display by applying Lemma H.1 below with the  
958 following substitutions

$$\begin{aligned} \mathbf{g} &\leftarrow \tilde{\mathbf{g}}, \quad \mathbf{h} \leftarrow \mathbf{g}, \quad c(\mathbf{h}) \leftarrow Y(\mathbf{g}) \\ p_2 &\leftarrow \frac{\gamma}{1-\gamma} (\boldsymbol{\pi}_1 - \mu^2 (\boldsymbol{\pi}_1 - \mathbf{\Pi}_{11} + \mathbf{\Pi}_{12})^2), \quad q_2 \leftarrow \frac{\gamma}{1-\gamma} (\boldsymbol{\pi}_1^2 + \mathbf{\Pi}_{12}^2 \mu^2 k - 2\mu^2 \mathbf{\Pi}_{12} (\boldsymbol{\pi}_1 - \mathbf{\Pi}_{11} + \mathbf{\Pi}_{12})), \\ p_1 &\leftarrow \gamma \boldsymbol{\pi}_1, \quad q_1 \leftarrow 0. \end{aligned}$$

959 This shows that with probability 1,  $\mathbb{P}_{e, \text{LS}} < \mathbb{P}_{e, \text{Avg}}$  if and only if  $p_2 < p_1 \Leftrightarrow \gamma < \gamma_* =$   
960  $\mu^2 (\boldsymbol{\pi}_1 - \mathbf{\Pi}_{11} + \mathbf{\Pi}_{12})^2 / \boldsymbol{\pi}_1$ . To retrieve, recall that  $\boldsymbol{\pi}_1 = 1/k$  and  $k\mathbf{\Pi}_{11} + (k^2 - k)\mathbf{\Pi}_{12} = 1$ . The  
961 only thing left to prove is that  $\gamma_* < 1$ . To see this note that  $p_2 > 0$  from positive semi-definiteness of  
962 the Schur matrix in (H.23). It takes simple algebra to conclude that  $p_2 > 0 \implies \gamma_* < 1$ .

**Lemma H.1** *Let  $k \geq 2$ ,  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}_k, \mathbf{I}_k)$ ,  $\mathbf{h} \sim \mathcal{N}(\mathbf{0}_k, \mathbf{I}_k)$ , and discrete random variable  $c(\mathbf{h})$  such that  $\mathbb{P}(c(\mathbf{h}) = \ell) = e^{\mathbf{h} \ell} / \sum_{j \in [k]} e^{\mathbf{h} j}$ . Consider the function  $F : \mathbb{R}_{>0} \times \mathbb{R} \rightarrow [0, 1]$  defined as follows*

$$F(p, q) = \mathbb{P} \left( \arg \max \left\{ (p\mathbf{I}_k - q\mathbf{1}_k\mathbf{1}_k^T)^{1/2} \mathbf{g} + (\alpha\mathbf{I}_k - \beta\mathbf{1}_k\mathbf{1}_k^T)^{1/2} \mathbf{h} \right\} \neq c(\mathbf{h}) \right),$$

963 such that  $p\mathbf{I}_k - q\mathbf{1}_k\mathbf{1}_k^T > 0$  and fixed  $\alpha\mathbf{I}_k - \beta\mathbf{1}_k\mathbf{1}_k^T > 0$ . Then, the following statements are true.

964 1.  $F(p, q) = \mathbb{P} \left( \arg \max \left\{ \sqrt{p} \cdot \mathbf{g} + \sqrt{\alpha} \mathbf{h} \right\} \neq c(\mathbf{h}) \right)$ .

965 2. For  $0 < p_2 < p_1$  and any  $q_1 < \frac{p_1}{k}, q_2 < \frac{p_2}{k}$ , it holds that  $F(p_2, q_2) < F(p_1, q_1)$ .

966 **Proof** [of Lemma H.1] Fix any  $p > 0, q \leq \frac{p}{k}$ . Denote  $\mathbf{T} := (p\mathbf{I}_k - q\mathbf{1}_k\mathbf{1}_k^T)^{1/2}$  and  $\mathbf{S} :=$   
967  $(\alpha\mathbf{I}_k - \beta\mathbf{1}_k\mathbf{1}_k^T)^{1/2}$  for convenience. It can be checked that  $\mathbf{T} := \left( \sqrt{p}\mathbf{I}_k + \frac{\sqrt{p-qk}-\sqrt{p}}{k} \right) \mathbf{1}_k\mathbf{1}_k^T$  and  
968  $\mathbf{S} := \left( \sqrt{\alpha}\mathbf{I}_k + \frac{\sqrt{\alpha-\beta k}-\sqrt{\alpha}}{k} \right) \mathbf{1}_k\mathbf{1}_k^T$ . From these, it follows that

$$F(p, q) = \mathbb{P} \left( \arg \max \left\{ \sqrt{p} \cdot \mathbf{g} + \sqrt{\alpha} \mathbf{h} \right\} \neq c(\mathbf{h}) \right).$$

969 This shows the first statement.

970 Next, we show the second statement. Using the distribution of  $c(\mathbf{h})$  and symmetry we have the  
971 following chain of inequalities:

$$\begin{aligned} 1 - F(p, q) &= \mathbb{P} \left( \arg \max_{j \in [k]} \left\{ \sqrt{p} \cdot \mathbf{g} + \sqrt{\alpha} \mathbf{h} \right\} = c(\mathbf{h}) \right) \\ &= k \cdot \mathbb{E} \left[ \frac{e^{\mathbf{h} k}}{\sum_{j \in [k]} e^{\mathbf{h} j}} \cdot \mathbb{1} \left\{ \arg \max \left\{ \sqrt{p} \cdot \mathbf{g} + \sqrt{\alpha} \mathbf{h} \right\} = k \right\} \right] \\ &= k \cdot \mathbb{E} \left[ \frac{e^{\mathbf{h} k}}{\sum_{j \in [k]} e^{\mathbf{h} j}} \cdot \prod_{j \in [k-1]} \mathbb{1} \left\{ \sqrt{p} \cdot \mathbf{g}_j + \sqrt{\alpha} \mathbf{h}_j < \sqrt{p} \cdot \mathbf{g}_k + \sqrt{\alpha} \mathbf{h}_k \right\} \right] \\ &= k \cdot \mathbb{E} \left[ \frac{e^{\mathbf{h} k}}{\sum_{j \in [k]} e^{\mathbf{h} j}} \cdot \prod_{j \in [k-1]} \mathbb{1} \left\{ \mathbf{g}_j < \mathbf{g}_k + \frac{\sqrt{\alpha} \mathbf{h}_k - \sqrt{\alpha} \mathbf{h}_j}{\sqrt{p}} \right\} \right] \\ &= k \cdot \mathbb{E} \left[ \frac{e^{\mathbf{h} k}}{\sum_{j \in [k]} e^{\mathbf{h} j}} \cdot \prod_{j \in [k-1]} Q \left( \mathbf{g}_k + \frac{\sqrt{\alpha} \mathbf{h}_j - \sqrt{\alpha} \mathbf{h}_k}{\sqrt{p}} \right) \right] =: k \cdot G(\sqrt{p}), \quad (\text{H.25}) \end{aligned}$$

where in the last line we used the rotational symmetry of the Gaussian distribution:

$$\mathbb{P} \left\{ \mathbf{g}_j < \mathbf{g}_k + \frac{\sqrt{\alpha} \mathbf{h}_k - \sqrt{\alpha} \mathbf{h}_j}{\sqrt{p}} \mid \mathbf{h}_1, \dots, \mathbf{h}_k \right\} = \mathbb{P} \left\{ \mathbf{g}_j > \mathbf{g}_k + \frac{\sqrt{\alpha} \mathbf{h}_j - \sqrt{\alpha} \mathbf{h}_k}{\sqrt{p}} \mid \mathbf{h}_1, \dots, \mathbf{h}_k \right\},$$

972 and the fact that  $\mathbf{g}_1, \dots, \mathbf{g}_{k-1}$  are independent.

Next, we will show that the function  $\mathcal{G}(\cdot)$  defined above is strictly decreasing in  $(0, \infty)$ . Towards this goal, using  $Q'(x) = -\frac{1}{\sqrt{2\pi}}e^{-x^2/2} = -\phi(x)$  and using the shorthand

$$H_{kj} = \mathbf{h}_j - \mathbf{h}_k, \quad j \in [k],$$

973 we may compute the derivative of  $\mathcal{G}$  at any  $s > 0$  as follows:

$$\begin{aligned} \frac{d\mathcal{G}(s)}{ds} &= \sum_{i \in [k-1]} \mathbb{E} \left[ \frac{e^{\mathbf{h}_k}}{\sum_{j \in [k]} e^{\mathbf{h}_j}} \cdot \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \frac{\sqrt{\alpha} H_{ki}}{s^2} \cdot \prod_{j \neq i \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right) \right] \\ &= \sum_{i \in [k-1]} \mathbb{E} \left[ \frac{1}{\sum_{j \in [k]} e^{H_{kj}}} \cdot \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \frac{\sqrt{\alpha} H_{ki}}{s^2} \cdot \prod_{j \neq i \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right) \right] \\ &= \sum_{i \in [k-1]} \frac{\sqrt{\alpha}}{s^2} \mathbb{E} [H_{ki} \cdot \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]})], \end{aligned} \quad (\text{H.26})$$

974 where in the last line we have defined

$$\mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]}) := \frac{1}{1 + \sum_{j \in [k-1]} e^{H_{kj}}} \cdot \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \prod_{j \neq i \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right),$$

975 Next, we use Gaussian integration by parts (GIBP) to further simplify the expression in (H.26). Fix  
976 any  $i \in [k-1]$ . Then, by (GIBP):

$$\begin{aligned} A_i &:= \mathbb{E} [H_{ki} \cdot \mathcal{A}_i(H_{k1}, H_{k2}, \dots, H_{kk})] \quad (\text{H.27}) \\ &= \mathbb{E} [H_{ki}^2] \mathbb{E} \left[ \frac{d}{dH_{ki}} \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]}) \right] + \sum_{\substack{\ell \in [k-1] \\ \ell \neq i}} \mathbb{E} [H_{ki} \cdot H_{k\ell}] \mathbb{E} \left[ \frac{d}{dH_{k\ell}} \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]}) \right] \\ &= 2 \underbrace{\mathbb{E} \left[ \frac{d}{dH_{ki}} \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]}) \right]}_{\text{TermI}} + \underbrace{\sum_{\substack{\ell \in [k-1] \\ \ell \neq i}} \mathbb{E} \left[ \frac{d}{dH_{k\ell}} \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]}) \right]}_{\text{TermII}}, \end{aligned} \quad (\text{H.28})$$

where in the second line, we used the fact that  $\mathbf{h} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$  to compute

$$\mathbb{E} [H_{ki}^2] = 2, \quad \text{and} \quad \mathbb{E} [H_{ki} H_{k\ell}] = 1, \quad \ell \neq i, \quad \ell \in [k-1].$$

977 We now compute the derivatives in (H.28). First, for any  $\ell \in [k-1]$ ,  $\ell \neq i$ ,

$$\begin{aligned} \frac{d\mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]})}{dH_{k\ell}} &= -\frac{e^{H_{k\ell}}}{(1 + \sum_{j \in [k-1]} e^{H_{kj}})^2} \cdot \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \prod_{j \neq i \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right) =: \text{TermII(a)}_\ell \\ &- \frac{\sqrt{\alpha}}{s} \cdot \frac{1}{1 + \sum_{j \in [k-1]} e^{H_{kj}}} \cdot \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{k\ell}}{s}\right) \cdot \prod_{j \neq (i, \ell) \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right) =: \text{TermII(b)}_\ell \\ &= \text{TermII(a)}_\ell + \text{TermII(b)}_\ell \end{aligned} \quad (\text{H.29})$$

978 Thus,

$$\text{TermII} = \sum_{\ell \neq i \in [k-1]} \mathbb{E} [\text{TermII(a)}_\ell] + \mathbb{E} [\text{TermII(b)}_\ell] =: \mathbb{E} [\text{TermII(a)}] + N_i, \quad (\text{H.30})$$

979 where we defined

$$N_i = -\frac{\sqrt{\alpha}}{s} \cdot \mathbb{E} \left[ \frac{\phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right)}{1 + \sum_{j \in [k-1]} e^{H_{kj}}} \cdot \sum_{\ell \neq i \in [k-1]} \left\{ \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{k\ell}}{s}\right) \cdot \prod_{j \neq (i, \ell) \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right) \right\} \right] < 0. \quad (\text{H.31})$$

980 and we remark for later use that

$$\mathbb{E} [\text{TermII(a)}] = \sum_{\ell \neq i \in [k-1]} \mathbb{E} [\text{TermII(a)}_\ell] < 0. \quad (\text{H.32})$$

981 Second, it holds that

$$\begin{aligned}
\frac{d\mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]})}{dH_{ki}} &= -\frac{e^{H_{ki}}}{(1 + \sum_{j \in [k-1]} e^{H_{kj}})^2} \cdot \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \prod_{j \neq i \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right) \\
&+ \frac{d\phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right)}{dH_{ki}} \cdot \frac{1}{1 + \sum_{j \in [k-1]} e^{H_{kj}}} \cdot \prod_{j \neq i \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right) \\
&= -\frac{e^{H_{ki}}}{(1 + \sum_{j \in [k-1]} e^{H_{kj}})^2} \cdot \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \prod_{j \neq i \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right) =: \text{TermI(a)} \\
&- \frac{\sqrt{\alpha}}{s} \left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \frac{1}{1 + \sum_{j \in [k-1]} e^{H_{kj}}} \cdot \prod_{j \neq i \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right) =: \text{TermI(b)}, \\
&= \text{TermI(a)} + \text{TermI(b)},
\end{aligned} \tag{H.33}$$

982 where in the penultimate line we used the fact that  $\phi'(x) = -x\phi(x)$ . Consider the two terms in  
983 (H.33). Clearly,

$$\mathbb{E}[\text{TermI(a)}] < 0. \tag{H.34}$$

984 For the second term we observe that:

$$\mathbb{E}[\text{TermI(b)}] = -\frac{\sqrt{\alpha}}{s} \mathbb{E}\left[\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]})\right] \tag{H.35}$$

$$\begin{aligned}
&= -\frac{\alpha}{s^2} \cdot \mathbb{E}[H_{ki} \cdot \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]})] - \frac{\sqrt{\alpha}}{s} \mathbb{E}[\mathbf{g}_k \cdot \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]})] \\
&= -\frac{\alpha}{s^2} \cdot A_i - \frac{\sqrt{\alpha}}{s} \mathbb{E}[\mathbf{g}_k \cdot \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]})].
\end{aligned} \tag{H.36}$$

985 Moreover, using again GIBP,  $\mathbb{E}[\mathbf{g}_k^2] = 1$ ,  $\mathbb{E}[\mathbf{g}_k H_{kj}] = 0$ ,  $j \in [k]$  and the fact that  $\phi'(x) = -x\phi(x)$ ,

$$\begin{aligned}
\mathbb{E}[\mathbf{g}_k \cdot \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]})] &= \mathbb{E}\left[\frac{d\mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k]})}{d\mathbf{g}_k}\right] \\
&= -\mathbb{E}\left[\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \mathcal{A}_i(\mathbf{g}_k, \{H_{kj}\}_{j \in [k-1]})\right] \\
&- \sum_{\ell \neq i \in [k-1]} \mathbb{E}\left[\phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{k\ell}}{s}\right) \cdot \frac{1}{1 + \sum_{j \in [k-1]} e^{H_{kj}}} \cdot \phi\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{ki}}{s}\right) \cdot \prod_{j \neq (i, \ell) \in [k-1]} Q\left(\mathbf{g}_k + \frac{\sqrt{\alpha} H_{kj}}{s}\right)\right] \\
&= \frac{s}{\sqrt{\alpha}} \mathbb{E}[\text{TermI(b)}] + \frac{s}{\sqrt{\alpha}} N_i,
\end{aligned} \tag{H.37}$$

986 where, we have recalled (H.35) and (H.31). Using (H.37) in (H.36), we find that

$$\mathbb{E}[\text{TermI(b)}] = -\frac{\alpha}{s^2} \cdot A_i - \mathbb{E}[\text{TermI(b)}] - N_i \implies \mathbb{E}[\text{TermI(b)}] = -\frac{\alpha}{2s^2} \cdot A_i - \frac{N_i}{2}. \tag{H.38}$$

987 We are now ready to put things together:

$$\begin{aligned}
A_i &= 2 \cdot \text{TermI} + \text{TermII} \quad \text{by (H.28)} \\
&= 2 \mathbb{E}[\text{TermI(a)}] + 2 \mathbb{E}[\text{TermI(b)}] + \mathbb{E}[\text{TermII(a)}] + N_i \quad \text{by (H.30)} \\
&= 2 \mathbb{E}[\text{TermI(a)}] - \frac{\alpha}{s^2} \cdot A_i - N_i + \mathbb{E}[\text{TermII(a)}] + N_i \quad \text{by (H.38)} \\
\implies A_i &= \frac{s^2}{s^2 + \alpha} (2 \mathbb{E}[\text{TermI(a)}] + \mathbb{E}[\text{TermII(a)}]) \\
&< 0 \quad \text{by (H.34) and (H.32)}.
\end{aligned}$$

988 From this, (H.26) and (H.27), we have shown that  $\mathcal{G}$  is strictly decreasing in  $(0, \infty)$ . Recalling  
989 the definition of  $\mathcal{G}$  in (H.25), this implies that  $F(p, q)$  is strictly increasing in  $p > 0$ , as desired to  
990 complete the proof.  $\blacksquare$

991 **I Weighted LS for GMM (Proof of Theorem 3.4)**

992 **I.1 Computing  $\Sigma_{w,\mu}$**

$$\begin{aligned} \min_{\mathbf{W} \in \mathbb{R}^{k \times d}, \mathbf{b} \in \mathbb{R}^k} \frac{1}{2n} \left\| (\mathbf{W}\mathbf{X} + \mathbf{b}\mathbf{1}_n^T - \mathbf{Y}) \mathbf{D} \right\|_F^2 &= \sum_{\ell=1}^k \min_{\mathbf{w}_\ell, b_\ell} \frac{1}{2n} \left\| \mathbf{D} (\mathbf{X}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell) \right\|_{\ell_2}^2 \\ &= \sum_{\ell=1}^k \min_{\mathbf{w}_\ell, b_\ell} \frac{1}{2n} \left\| \mathbf{D} (\mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + \mathbf{Z}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell) \right\|_{\ell_2}^2 \end{aligned}$$

993 Define

$$\mathcal{L}_{PO}(\mathbf{w}_\ell, b_\ell) := \frac{1}{2n} \left\| \mathbf{D} (\mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + \mathbf{Z}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell) \right\|_{\ell_2}^2 \quad (\text{I.1})$$

994 **Identifying the AO.** To continue further note that by duality we have

$$\min_{\mathbf{w}_\ell, b_\ell} \mathcal{L}_{PO}(\mathbf{w}_\ell, b_\ell) = \min_{\mathbf{u}, \mathbf{w}_\ell, b_\ell} \max_{\mathbf{s}} \frac{1}{n} \left( \mathbf{s}^T \mathbf{D} \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + \mathbf{s}^T \mathbf{D} \mathbf{Z}^T \mathbf{w}_\ell + b_\ell \mathbf{s}^T \mathbf{D} \mathbf{1}_n - \mathbf{s}^T \mathbf{D} \mathbf{Y}_\ell - \mathbf{s}^T \mathbf{u} + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2} \right)$$

995 Note that the above is jointly convex in  $(\mathbf{u}, \mathbf{w}_\ell, b_\ell)$  and concave in  $\mathbf{s}$ . We consider the Auxiliary  
996 Optimization (AO) problem

$$\min_{\mathbf{u}, \mathbf{w}_\ell, b_\ell} \max_{\mathbf{s}} \frac{1}{n} \left( \mathbf{s}^T \mathbf{D} \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + \sigma \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g}^T \mathbf{D} \mathbf{s} + \sigma \|\mathbf{D} \mathbf{s}\|_{\ell_2} \mathbf{h}^T \mathbf{w}_\ell + b_\ell \mathbf{s}^T \mathbf{D} \mathbf{1}_n - \mathbf{s}^T \mathbf{D} \mathbf{Y}_\ell - \mathbf{s}^T \mathbf{u} + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2} \right),$$

997 where  $\mathbf{g} \in \mathbb{R}^n$  and  $\mathbf{h} \in \mathbb{R}^d$  are independent Gaussian random vectors with i.i.d.  $\mathcal{N}(0, 1)$  entries. To  
998 continue further we carry out a change of variable  $\mathbf{s} \rightarrow \mathbf{D} \mathbf{s}$  to arrive at

$$\min_{\mathbf{u}, \mathbf{w}_\ell, b_\ell} \max_{\mathbf{s}} \frac{1}{n} \left( \mathbf{s}^T \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + \sigma \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g}^T \mathbf{s} + \sigma \|\mathbf{s}\|_{\ell_2} \mathbf{h}^T \mathbf{w}_\ell + b_\ell \mathbf{s}^T \mathbf{1}_n - \mathbf{s}^T \mathbf{Y}_\ell - \mathbf{s}^T \mathbf{D}^{-1} \mathbf{u} + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2} \right),$$

999 Maximizing over the direction of  $\mathbf{s}$  and its norm  $\beta = \|\mathbf{s}\|_{\ell_2}$  we arrive at

$$\begin{aligned} \min_{\mathbf{u}, \mathbf{w}_\ell, b_\ell} \max_{\beta \geq 0} \max_{\mathbf{s}: \|\mathbf{s}\|_{\ell_2} = 1} \frac{1}{n} &\left( \beta \mathbf{s}^T \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + \sigma \beta \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g}^T \mathbf{s} + \sigma \beta \mathbf{h}^T \mathbf{w}_\ell + b_\ell \beta \mathbf{s}^T \mathbf{1}_n - \beta \mathbf{s}^T \mathbf{Y}_\ell - \beta \mathbf{s}^T \mathbf{D}^{-1} \mathbf{u} + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2} \right) \\ &= \min_{\mathbf{u}, \mathbf{w}_\ell, b_\ell} \max_{\beta \geq 0} \frac{1}{n} \left( \beta \left\| \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + \sigma \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell - \mathbf{D}^{-1} \mathbf{u} \right\|_{\ell_2} + \sigma \beta \mathbf{h}^T \mathbf{w}_\ell + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2} \right) \\ &= \min_{\mathbf{u}, b_\ell} \max_{\beta \geq 0} \min_{\mathbf{w}_\ell} \frac{1}{n} \left( \beta \left\| \mathbf{Y}^T \mathbf{M}^T \mathbf{w}_\ell + \sigma \|\mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell - \mathbf{D}^{-1} \mathbf{u} \right\|_{\ell_2} + \sigma \beta \mathbf{h}^T \mathbf{w}_\ell + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2} \right) \end{aligned}$$

1000 To continue, consider the singular value decomposition

$$\mathbf{M} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = [\mathbf{u}_1 \quad \mathbf{u}_2 \quad \dots \quad \mathbf{u}_r] \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r) \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_r^T \end{bmatrix}, \quad (\text{I.2})$$

1001 with  $r := \text{rank}(\mathbf{M}) \leq k$  and define the variable  $\boldsymbol{\alpha} = \mathbf{U}^T \mathbf{w}_\ell$  and  $\boldsymbol{\alpha}_\perp = \mathbf{U}_\perp^T \mathbf{w}_\ell$  where  $\mathbf{U}_\perp$  is the  
1002 orthogonal complement of the columns of  $\mathbf{U}$ . With these definitions the above optimization problem  
1003 reduces to

$$\min_{\mathbf{u}, b_\ell} \max_{\beta \geq 0} \min_{\boldsymbol{\alpha}} \min_{\boldsymbol{\alpha}_\perp} \frac{1}{n} \left( \beta \left\| \mathbf{Y}^T \mathbf{V} \mathbf{\Sigma} \boldsymbol{\alpha} + \sigma \sqrt{\|\boldsymbol{\alpha}\|_{\ell_2}^2 + \|\boldsymbol{\alpha}_\perp\|_{\ell_2}^2} \mathbf{g} + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell - \mathbf{D}^{-1} \mathbf{u} \right\|_{\ell_2} + \sigma \beta \mathbf{h}^T \mathbf{U} \boldsymbol{\alpha} + \sigma \beta \mathbf{h}^T \mathbf{U}_\perp \boldsymbol{\alpha}_\perp + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2} \right)$$

1004 Decomposing the optimization over  $\boldsymbol{\alpha}_\perp$  in terms of its direction and norm  $\alpha_0 = \|\boldsymbol{\alpha}_\perp\|_{\ell_2}$  we arrive at

$$\min_{\mathbf{u}, b_\ell} \max_{\beta \geq 0} \min_{\boldsymbol{\alpha}} \min_{\alpha_0 \geq 0} \frac{1}{n} \left( \beta \left\| \mathbf{Y}^T \mathbf{V} \mathbf{\Sigma} \boldsymbol{\alpha} + \sigma \sqrt{\|\boldsymbol{\alpha}\|_{\ell_2}^2 + \alpha_0^2} \mathbf{g} + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell - \mathbf{D}^{-1} \mathbf{u} \right\|_{\ell_2} + \sigma \beta \mathbf{h}^T \mathbf{U} \boldsymbol{\alpha} - \sigma \alpha_0 \beta \left\| \mathbf{U}_\perp^T \mathbf{h} \right\|_{\ell_2} + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2} \right)$$



1005 Since  $U^T \mathbf{h}$  is  $r \leq k$  dimensional in our asymptotic regime the  $\mathbf{h}^T U \boldsymbol{\alpha}$  term can be ignored. Also  
 1006 replacing  $\beta$  with  $\beta/\sqrt{n}$  we thus arrive at

$$\begin{aligned}
 & \min_{\mathbf{u}, \mathbf{b}_\ell} \max_{\beta \geq 0} \min_{\boldsymbol{\alpha}} \min_{\alpha_0 \geq 0} \frac{\beta}{\sqrt{n}} \left\| \mathbf{Y}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} + \sigma \sqrt{\|\boldsymbol{\alpha}\|_{\ell_2}^2 + \alpha_0^2} \mathbf{g} + \mathbf{b}_\ell \mathbf{1}_n - \mathbf{Y}_\ell - \mathbf{D}^{-1} \mathbf{u} \right\|_{\ell_2} - \frac{1}{\sqrt{n}} \sigma \alpha_0 \beta \|\mathbf{U}_\perp^T \mathbf{h}\|_{\ell_2} + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2n} \\
 &= \min_{\mathbf{u}, \mathbf{b}_\ell} \max_{\beta \geq 0} \min_{\boldsymbol{\alpha}} \min_{\alpha_0 \geq 0} \min_{\tau \geq 0} \frac{\beta}{2n\tau} \left\| \mathbf{Y}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} + \sigma \sqrt{\|\boldsymbol{\alpha}\|_{\ell_2}^2 + \alpha_0^2} \mathbf{g} + \mathbf{b}_\ell \mathbf{1}_n - \mathbf{Y}_\ell - \mathbf{D}^{-1} \mathbf{u} \right\|_{\ell_2}^2 + \frac{\beta\tau}{2} \\
 &\quad - \frac{1}{\sqrt{n}} \sigma \alpha_0 \beta \|\mathbf{U}_\perp^T \mathbf{h}\|_{\ell_2} + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2n} \\
 &= \min_{\mathbf{b}_\ell} \max_{\beta \geq 0} \min_{\boldsymbol{\alpha}} \min_{\alpha_0 \geq 0} \min_{\tau \geq 0} \min_{\mathbf{u}} \frac{\beta}{2n\tau} \left\| \mathbf{Y}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} + \sigma \sqrt{\|\boldsymbol{\alpha}\|_{\ell_2}^2 + \alpha_0^2} \mathbf{g} + \mathbf{b}_\ell \mathbf{1}_n - \mathbf{Y}_\ell - \mathbf{D}^{-1} \mathbf{u} \right\|_{\ell_2}^2 + \frac{\beta\tau}{2} \\
 &\quad - \frac{1}{\sqrt{n}} \sigma \alpha_0 \beta \|\mathbf{U}_\perp^T \mathbf{h}\|_{\ell_2} + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2n}
 \end{aligned}$$

1007 Setting the derivative with respect to  $\mathbf{u}$  to zero we arrive at

$$\mathbf{u} = \frac{\beta}{\tau} \mathbf{D}^{-1} \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \left( \mathbf{Y}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} + \sigma \sqrt{\|\boldsymbol{\alpha}\|_{\ell_2}^2 + \alpha_0^2} \mathbf{g} + \mathbf{b}_\ell \mathbf{1}_n - \mathbf{Y}_\ell \right)$$

1008 Plugging the latter into the above the AO simplifies to

$$\min_{\mathbf{b}_\ell} \max_{\beta \geq 0} \min_{\boldsymbol{\alpha}} \min_{\alpha_0 \geq 0} \min_{\tau \geq 0} \frac{\beta}{2\tau n} \text{trace} \left( \mathbf{t}^T \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \mathbf{t} \right) - \frac{1}{\sqrt{n}} \sigma \alpha_0 \beta \|\mathbf{U}_\perp^T \mathbf{h}\|_{\ell_2} + \frac{\beta\tau}{2}$$

1009 where

$$\mathbf{t} := \mathbf{Y}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} + \sigma \sqrt{\|\boldsymbol{\alpha}\|_{\ell_2}^2 + \alpha_0^2} \mathbf{g} + \mathbf{b}_\ell \mathbf{1}_n - \mathbf{Y}_\ell = \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + \sigma \sqrt{\|\boldsymbol{\alpha}\|_{\ell_2}^2 + \alpha_0^2} \mathbf{g} + \mathbf{b}_\ell \mathbf{1}_n$$

1010 To continue note that in our asymptotic regime we have

$$\frac{1}{\sqrt{n}} \|\mathbf{U}_\perp^T \mathbf{h}\|_{\ell_2} \xrightarrow{P} \sqrt{\gamma}$$

1011 and the cross terms can be ignored so that in an asymptotic sense

$$\begin{aligned}
 & \frac{1}{n} \text{trace} \left( \mathbf{t}^T \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \mathbf{t} \right) \\
 &= \frac{1}{n} \sigma^2 \left( \|\boldsymbol{\alpha}\|_{\ell_2}^2 + \alpha_0^2 \right) \text{trace} \left( \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \right) \\
 &\quad + \frac{1}{n} \text{trace} \left( \left( \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + \mathbf{b}_\ell \mathbf{1}_n \right)^T \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \left( \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + \mathbf{b}_\ell \mathbf{1}_n \right) \right)
 \end{aligned}$$

1012 Therefore we arrive at

$$\begin{aligned}
 & \min_{\mathbf{b}_\ell} \max_{\beta \geq 0} \min_{\boldsymbol{\alpha}} \min_{\tau \geq 0} \min_{\alpha_0 \geq 0} \frac{\beta}{2\tau n} \sigma^2 \left( \|\boldsymbol{\alpha}\|_{\ell_2}^2 + \alpha_0^2 \right) \text{trace} \left( \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \right) - \sigma \alpha_0 \beta \sqrt{\gamma} + \frac{\beta\tau}{2} \\
 &\quad + \frac{\beta}{2\tau n} \text{trace} \left( \left( \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + \mathbf{b}_\ell \mathbf{1}_n \right)^T \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \left( \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + \mathbf{b}_\ell \mathbf{1}_n \right) \right)
 \end{aligned}$$

1013 which can be rewritten in the form

$$\begin{aligned}
 & \min_{\mathbf{b}_\ell} \max_{\beta \geq 0} \min_{\boldsymbol{\alpha}} \min_{\tau \geq 0} \frac{\beta}{2\tau n} \sigma^2 \|\boldsymbol{\alpha}\|_{\ell_2}^2 \text{trace} \left( \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \right) + \frac{\beta\tau}{2} \\
 &\quad + \frac{\beta}{2\tau n} \text{trace} \left( \left( \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + \mathbf{b}_\ell \mathbf{1}_n \right)^T \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \left( \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + \mathbf{b}_\ell \mathbf{1}_n \right) \right) \\
 &\quad + \frac{\beta}{2\tau n} \sigma^2 \alpha_0^2 \text{trace} \left( \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \right) - \alpha_0 \sigma \beta \sqrt{\gamma}
 \end{aligned}$$

1014 To continue further we shall assume  $\mathbf{D} = \text{diag}(\mathbf{Y}^T \boldsymbol{\omega})$ . Note that in this case

$$\frac{1}{n} \text{trace} \left( \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \right) = \frac{1}{n} \sum_{i=1}^n \frac{(\mathbf{y}_i^T \boldsymbol{\omega})^2}{(\mathbf{y}_i^T \boldsymbol{\omega})^2 + \frac{\beta}{\tau}} = \sum_{\ell=1}^k \frac{n_\ell}{n} \frac{\omega_\ell^2}{\omega_\ell^2 + \frac{\beta}{\tau}} \xrightarrow{P} \sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{\omega_\ell^2 + \frac{\beta}{\tau}}$$

1015 Also

$$\begin{aligned} & \frac{1}{n} \text{trace} \left( \left( \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + b_\ell \mathbf{1}_n \right)^T \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \left( \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + b_\ell \mathbf{1}_n \right) \right) \\ &= \frac{1}{n} \text{trace} \left( (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell)^T \mathbf{Y} \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) \right) \\ & \quad + \frac{2}{n} b_\ell \mathbf{1}_n^T \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \mathbf{Y}^T (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + \frac{b_\ell^2}{n} \text{trace} \left( \left( \mathbf{I} + \frac{\beta}{\tau} \mathbf{D}^{-2} \right)^{-1} \right) \\ & \xrightarrow{P} \text{trace} \left( (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell)^T \text{diag} \left( \frac{\pi_1 \omega_1^2}{\omega_1^2 + \frac{\beta}{\tau}}, \frac{\pi_2 \omega_2^2}{\omega_2^2 + \frac{\beta}{\tau}}, \dots, \frac{\pi_k \omega_k^2}{\omega_k^2 + \frac{\beta}{\tau}} \right) (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) \right) \\ & \quad + 2b_\ell \left[ \frac{\pi_1 \omega_1^2}{\omega_1^2 + \frac{\beta}{\tau}} \quad \frac{\pi_2 \omega_2^2}{\omega_2^2 + \frac{\beta}{\tau}} \quad \dots \quad \frac{\pi_k \omega_k^2}{\omega_k^2 + \frac{\beta}{\tau}} \right] (\mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\alpha} - \mathbf{e}_\ell) + b_\ell^2 \left( \sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{\omega_\ell^2 + \frac{\beta}{\tau}} \right) \end{aligned}$$

1016 Next define

$$\begin{aligned} \mathbf{A}(\eta) &:= \begin{bmatrix} \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}(\eta)) \mathbf{I} + \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\pi}) \text{diag}(\boldsymbol{\nu}(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\nu}(\eta)) \boldsymbol{\pi} \\ \boldsymbol{\pi}^T \text{diag}(\boldsymbol{\nu}(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\pi}^T \boldsymbol{\nu}(\eta) \end{bmatrix} \\ \mathbf{c}_\ell &:= \begin{bmatrix} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{e}_\ell \\ 1 \end{bmatrix} \end{aligned} \quad (I.3)$$

1017 where

$$\boldsymbol{\nu}(\eta) = \frac{1}{\gamma} \begin{bmatrix} \frac{\omega_1^2}{\omega_1^2 + \eta} \\ \frac{\omega_2^2}{\omega_2^2 + \eta} \\ \dots \\ \frac{\omega_k^2}{\omega_k^2 + \eta} \end{bmatrix} \quad (I.4)$$

1018 We thus arrive at

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \min_{b_\ell} \max_{\alpha_0 \geq 0} \max_{\beta \geq 0} \min_{\tau \geq 0} & \frac{\gamma \beta}{2\tau} \left( \boldsymbol{\pi}_\ell \boldsymbol{\nu}_\ell \left( \frac{\beta}{\tau} \right) + [\boldsymbol{\alpha}^T \quad b_\ell] \mathbf{A} \left( \frac{\beta}{\tau} \right) \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - 2\pi_\ell \boldsymbol{\nu}_\ell \left( \frac{\beta}{\tau} \right) \mathbf{c}_\ell^T \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} \right) \\ & + \frac{\gamma \beta}{2\tau} \sigma^2 \left( \boldsymbol{\pi}^T \boldsymbol{\nu} \left( \frac{\beta}{\tau} \right) \right) \alpha_0^2 - \alpha_0 \sigma \beta \sqrt{\gamma} + \frac{\beta \tau}{2} \end{aligned}$$

1019 Setting the derivative of the above with respect to  $\alpha_0$  to zero we arrive at

$$\frac{\gamma \beta}{\tau} \sigma^2 \left( \boldsymbol{\pi}^T \boldsymbol{\nu} \left( \frac{\beta}{\tau} \right) \right) \alpha_0 - \sigma \beta \sqrt{\gamma} = 0 \quad \Rightarrow \quad \alpha_0 = \frac{\tau}{\sigma \sqrt{\gamma} \left( \boldsymbol{\pi}^T \boldsymbol{\nu} \left( \frac{\beta}{\tau} \right) \right)}$$

1020 Note that the above objective has the form

$$f \left( \frac{\beta}{\tau} \right) - \alpha_0 \sigma \beta \sqrt{\gamma} + \frac{\beta \tau}{2}$$

1021 with

$$f(\eta) := \frac{\eta \gamma}{2} \left( \boldsymbol{\pi}_\ell \boldsymbol{\nu}_\ell(\eta) + [\boldsymbol{\alpha}^T \quad b_\ell] \mathbf{A}(\eta) \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - 2\pi_\ell \boldsymbol{\nu}_\ell(\eta) \mathbf{c}_\ell^T \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} \right) + \frac{\gamma \eta}{2} \sigma^2 \left( \boldsymbol{\pi}^T \boldsymbol{\nu}(\eta) \right) \alpha_0^2.$$

1022 Thus, the derivatives with respect to  $\beta$  and  $\tau$  to zero we have

$$\frac{1}{\tau} f' \left( \frac{\beta}{\tau} \right) - \alpha_0 \sigma \sqrt{\gamma} + \frac{\tau}{2} = 0 \quad \Rightarrow \quad f' \left( \frac{\beta}{\tau} \right) - \alpha_0 \sigma \sqrt{\gamma} \tau + \frac{\tau^2}{2} = 0 \quad \Rightarrow \quad f' \left( \frac{\beta}{\tau} \right) = \tau^2 \left( \frac{1}{\boldsymbol{\pi}^T \boldsymbol{\nu} \left( \frac{\beta}{\tau} \right)} - \frac{1}{2} \right)$$

1023 and

$$-\frac{\beta}{\tau^2} f' \left( \frac{\beta}{\tau} \right) + \frac{\beta}{2} = 0 \quad \Rightarrow \quad \tau^2 = 2f' \left( \frac{\beta}{\tau} \right)$$

1024 Combining the latter two we conclude that  $\boldsymbol{\pi}^T \boldsymbol{\nu} \left( \frac{\beta}{\tau} \right) = 1$ . Thus,  $\eta = \frac{\beta}{\tau}$  is the solution to  $\boldsymbol{\pi}^T \boldsymbol{\nu}(\eta) = 1$ .

1025 To calculate  $\tau$  and hence  $\alpha_0$  we calculate  $f'$  which is equal to

$$\begin{aligned} f'(\eta) &= \frac{\gamma}{2} \left( \pi_\ell \boldsymbol{\nu}_\ell(\eta) + [\boldsymbol{\alpha}^T \quad b_\ell] \mathbf{A}(\eta) \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - 2\pi_\ell \boldsymbol{\nu}_\ell(\eta) \mathbf{c}_\ell^T \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} \right) + \frac{\gamma}{2} \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}(\eta)) \alpha_0^2 + \frac{\gamma\eta}{2} \sigma^2 \alpha_0^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta)) \\ &\quad + \frac{\gamma\eta}{2} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + [\boldsymbol{\alpha}^T \quad b_\ell] \mathbf{A}'(\eta) \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - 2\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell^T \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} \right) \end{aligned}$$

1026 Where

$$\begin{aligned} \boldsymbol{\nu}'(\eta) &= -\frac{1}{\gamma} \begin{bmatrix} \frac{\omega_1^2}{(\omega_1^2 + \eta)^2} \\ \frac{\omega_2^2}{(\omega_2^2 + \eta)^2} \\ \dots \\ \frac{\omega_k^2}{(\omega_k^2 + \eta)^2} \end{bmatrix} \\ \mathbf{A}'(\eta) &:= \begin{bmatrix} \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta)) \mathbf{I} + \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\pi}) \text{diag}(\boldsymbol{\nu}'(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\nu}'(\eta)) \boldsymbol{\pi} \\ \boldsymbol{\pi}^T \text{diag}(\boldsymbol{\nu}'(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta) \end{bmatrix} \\ \mathbf{c}_\ell &:= \begin{bmatrix} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{e}_\ell \\ 1 \end{bmatrix} \end{aligned}$$

1027 Now note that at the optimal point we have

$$f'(\eta) = \frac{\tau^2}{2} = \frac{\gamma}{2} \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}(\eta)) \alpha_0^2$$

1028 Thus from the above we can conclude that

$$\begin{aligned} \alpha_0^2 &= -\frac{1}{\eta \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta))} \left( \pi_\ell \boldsymbol{\nu}_\ell(\eta) + [\boldsymbol{\alpha}^T \quad b_\ell] \mathbf{A}(\eta) \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - 2\pi_\ell \boldsymbol{\nu}_\ell(\eta) \mathbf{c}_\ell^T \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} \right) \\ &\quad - \frac{1}{\sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta))} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + [\boldsymbol{\alpha}^T \quad b_\ell] \mathbf{A}'(\eta) \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - 2\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell^T \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} \right) \end{aligned}$$

1029 Thus the AO optimization problem reduces to

$$\min_{\mathbf{b}_\ell} \min_{\boldsymbol{\alpha}} \frac{\eta\gamma}{2} \left( \pi_\ell \boldsymbol{\nu}_\ell + [\boldsymbol{\alpha}^T \quad b_\ell] \mathbf{A} \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - 2\pi_\ell \boldsymbol{\nu}_\ell \mathbf{c}_\ell^T \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} \right)$$

1030 where  $\eta$  is the solution to

$$\sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{\omega_\ell^2 + \eta} = \gamma$$

1031 and

$$\begin{aligned} \mathbf{A}(\eta) &:= \begin{bmatrix} \sigma^2 \mathbf{I}_r + \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu}) \\ (\boldsymbol{\pi} \odot \boldsymbol{\nu})^T \mathbf{V} \boldsymbol{\Sigma} & 1 \end{bmatrix} \\ \mathbf{c}_\ell &:= \begin{bmatrix} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{e}_\ell \\ 1 \end{bmatrix} \end{aligned} \tag{I.5}$$

1032 where

$$\boldsymbol{\nu} := \frac{1}{\gamma} \begin{bmatrix} \frac{\omega_1^2}{\omega_1^2 + \eta} \\ \frac{\omega_2^2}{\omega_2^2 + \eta} \\ \dots \\ \frac{\omega_k^2}{\omega_k^2 + \eta} \end{bmatrix}$$

1033 First, note that the matrix  $\mathbf{A}$  is positive definite. This can be checked by computing the Schur  
1034 complement of  $\mathbf{A}$ :

$$\mathbf{\Delta} := \sigma^2 \mathbf{I}_r + \mathbf{\Sigma} \mathbf{V}^T \mathbf{P} \mathbf{V} \mathbf{\Sigma} := \sigma^2 \mathbf{I}_r + \mathbf{\Sigma} \mathbf{V}^T \left( \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}) - (\boldsymbol{\pi} \odot \boldsymbol{\nu})(\boldsymbol{\pi} \odot \boldsymbol{\nu})^T \right) \mathbf{V} \mathbf{\Sigma} > \mathbf{0}_{r \times r}. \quad (\text{I.6})$$

1035 Positive definiteness above holds because  $\mathbf{P} := \left( \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}) - (\boldsymbol{\pi} \odot \boldsymbol{\nu})(\boldsymbol{\pi} \odot \boldsymbol{\nu})^T \right) \geq \mathbf{0}_{k \times k}$ . Thus  
1036 the objective is a strictly convex quadratic and is jointly convex in its arguments. We proceed by  
1037 minimizing the objective over  $(\boldsymbol{\alpha}, b_\ell)$  which is equal to

$$\begin{bmatrix} \hat{\boldsymbol{\alpha}} \\ \hat{b}_\ell \end{bmatrix} = \pi_\ell \nu_\ell \mathbf{A}^{-1} \mathbf{c}_\ell = \pi_\ell \nu_\ell \begin{bmatrix} -\mathbf{\Delta}^{-1} \mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \\ 1 + (\boldsymbol{\pi} \odot \boldsymbol{\nu})^T \mathbf{V} \mathbf{\Sigma} \mathbf{\Delta}^{-1} \mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \end{bmatrix} \quad (\text{I.7})$$

1038 Thus, the minimum value attained is

$$\begin{aligned} & -\pi_\ell^2 \nu_\ell^2 \left[ -(\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell)^T \mathbf{V} \mathbf{\Sigma} \quad 1 \right] \begin{bmatrix} \mathbf{\Delta}^{-1} & \mathbf{0} \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} -\mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \\ 1 \end{bmatrix} \\ & = -\pi_\ell^2 \nu_\ell^2 \left( 1 + (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell)^T \mathbf{V} \mathbf{\Sigma} \mathbf{\Delta}^{-1} \mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \right). \end{aligned} \quad (\text{I.8})$$

1039 Thus the objective reduces to

$$\frac{\eta \gamma}{2} \left( \pi_\ell \nu_\ell (1 - \pi_\ell \nu_\ell) - \pi_\ell^2 \nu_\ell^2 (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell)^T \mathbf{V} \mathbf{\Sigma} \mathbf{\Delta}^{-1} \mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \right). \quad (\text{I.9})$$

1040 Therefore,

$$\begin{aligned} \alpha_0^2 &= -\frac{1}{\eta \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta))} \left( \pi_\ell \nu_\ell (1 - \pi_\ell \nu_\ell) - \pi_\ell^2 \nu_\ell^2 (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell)^T \mathbf{V} \mathbf{\Sigma} \mathbf{\Delta}^{-1} \mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \right) \\ &\quad - \frac{1}{\sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta))} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + [\boldsymbol{\alpha}^T \quad b_\ell] \mathbf{A}'(\eta) \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - 2\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell^T \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} \right) \end{aligned}$$

1041 where

$$\begin{bmatrix} \hat{\boldsymbol{\alpha}} \\ \hat{b}_\ell \end{bmatrix} = \pi_\ell \nu_\ell \begin{bmatrix} -\mathbf{\Delta}^{-1} \mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \\ 1 + (\boldsymbol{\pi} \odot \boldsymbol{\nu})^T \mathbf{V} \mathbf{\Sigma} \mathbf{\Delta}^{-1} \mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \end{bmatrix} \quad (\text{I.10})$$

1042 and

$$\begin{aligned} \boldsymbol{\nu}'(\eta) &= -\frac{1}{\gamma} \begin{bmatrix} \frac{\omega_1^2}{(\omega_1^2 + \eta)^2} \\ \frac{\omega_2^2}{(\omega_2^2 + \eta)^2} \\ \dots \\ \frac{\omega_k^2}{(\omega_k^2 + \eta)^2} \end{bmatrix} \\ \mathbf{A}'(\eta) &:= \begin{bmatrix} \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta)) \mathbf{I} + \mathbf{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\pi}) \text{diag}(\boldsymbol{\nu}'(\eta)) \mathbf{V} \mathbf{\Sigma} & \mathbf{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\nu}'(\eta)) \boldsymbol{\pi} \\ \boldsymbol{\pi}^T \text{diag}(\boldsymbol{\nu}'(\eta)) \mathbf{V} \mathbf{\Sigma} & \boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta) \end{bmatrix} \end{aligned}$$

1043 To continue note that

$$\begin{aligned} & \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + [\boldsymbol{\alpha}^T \quad b_\ell] \mathbf{A}'(\eta) \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} - 2\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell^T \begin{bmatrix} \boldsymbol{\alpha} \\ b_\ell \end{bmatrix} \\ & = \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_\ell^2 \nu_\ell^2 \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{A}'(\eta) \mathbf{A}^{-1} \mathbf{c}_\ell - 2\pi_\ell^2 \nu_\ell \nu_\ell \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{c}_\ell \end{aligned}$$

1044 Thus,

$$\begin{aligned} \alpha_0^2 &= -\frac{1}{\eta \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta))} \left( \pi_\ell \nu_\ell (1 - \pi_\ell \nu_\ell) - \pi_\ell^2 \nu_\ell^2 (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell)^T \mathbf{V} \mathbf{\Sigma} \mathbf{\Delta}^{-1} \mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \right) \\ &\quad - \frac{1}{\sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta))} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_\ell^2 \nu_\ell^2 \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{A}'(\eta) \mathbf{A}^{-1} \mathbf{c}_\ell - 2\pi_\ell^2 \nu_\ell \nu_\ell \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{c}_\ell \right) \\ &= \frac{\gamma}{\eta \sigma^2 \left( \sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{(\omega_\ell^2 + \eta)^2} \right)} \left( \pi_\ell \nu_\ell (1 - \pi_\ell \nu_\ell) - \pi_\ell^2 \nu_\ell^2 (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell)^T \mathbf{V} \mathbf{\Sigma} \mathbf{\Delta}^{-1} \mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \right) \\ &\quad + \frac{\gamma}{\sigma^2 \left( \sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{(\omega_\ell^2 + \eta)^2} \right)} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_\ell^2 \nu_\ell^2 \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{A}'(\eta) \mathbf{A}^{-1} \mathbf{c}_\ell - 2\pi_\ell^2 \nu_\ell \nu_\ell \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{c}_\ell \right) \\ &:= \frac{\zeta}{\sigma^2} \left( \pi_\ell \nu_\ell (1 - \pi_\ell \nu_\ell) - \pi_\ell^2 \nu_\ell^2 (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell)^T \mathbf{V} \mathbf{\Sigma} \mathbf{\Delta}^{-1} \mathbf{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \right) \\ &\quad + \frac{\zeta \eta}{\sigma^2} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_\ell^2 \nu_\ell^2 \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{A}'(\eta) \mathbf{A}^{-1} \mathbf{c}_\ell - 2\pi_\ell^2 \nu_\ell \nu_\ell \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{c}_\ell \right) \end{aligned}$$

1045 where  $\zeta := \frac{\gamma}{\eta \left( \sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{(\omega_\ell^2 + \eta)^2} \right)}$ .

**Asymptotic predictions.** First, from (I.7) the bias term converges as follows:

$$\widehat{b}_\ell \xrightarrow{P} \pi_\ell \nu_\ell \left( 1 + (\boldsymbol{\pi} \odot \boldsymbol{\nu})^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \right).$$

1046 Thus,

$$\widehat{\mathbf{b}} \xrightarrow{P} (\mathbf{I}_k - P \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T) (\boldsymbol{\pi} \odot \boldsymbol{\nu})$$

1047 Recall that  $\boldsymbol{\alpha} = \mathbf{U}^T \mathbf{w}_\ell$ . Thus, the correlations  $\langle \boldsymbol{\mu}_i, \mathbf{w}_\ell \rangle$ ,  $i \in [k]$  converge as follows:

$$\mathbf{M}^T \mathbf{w}_\ell = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{w}_\ell \xrightarrow{P} \mathbf{V} \boldsymbol{\Sigma} \widehat{\boldsymbol{\alpha}} = -\pi_\ell \nu_\ell \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell). \quad (\text{I.11})$$

1048 Here, convergence applies element-wise to the entries of the involved random vectors. Moreover, from  
 1049 the analysis above we can predict the limit of the norm  $\|\mathbf{w}_\ell\|_{\ell_2}$ . For this, note that  $\|\mathbf{w}_\ell\|_{\ell_2}^2 = \widehat{\alpha}_0^2 + \widehat{\boldsymbol{\alpha}}^T \widehat{\boldsymbol{\alpha}}$ .  
 1050 Thus,

$$\begin{aligned} \|\mathbf{w}_\ell\|_{\ell_2}^2 &\xrightarrow{P} \frac{\zeta}{\sigma^2} \pi_\ell \nu_\ell (1 - \pi_\ell \nu_\ell) + \pi_\ell^2 \nu_\ell^2 (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell)^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\zeta}{\sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) \\ &\quad + \frac{\eta \zeta}{\sigma^2} (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_\ell^2 \nu_\ell^2 \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{A}' \mathbf{A}^{-1} \mathbf{c}_\ell - 2\pi_\ell^2 \nu_\ell' \nu_\ell \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{c}_\ell) \end{aligned} \quad (\text{I.12})$$

## 1051 I.2 Computing $\Sigma_{\mathbf{w}, \mathbf{w}}$

1052 In the previous section we used the CGMT to predict the bias  $\widehat{b}_\ell$ , the correlations  $\langle \boldsymbol{\mu}_i, \widehat{\mathbf{w}}_\ell \rangle$ ,  $i \in [k]$   
 1053 and the norm  $\|\widehat{\mathbf{w}}_\ell\|_{\ell_2}$  for all  $\ell \in [k]$  members of the multi-output classifier. Here, we show how to  
 1054 compute the limits of the cross-correlations  $\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j \rangle$ ,  $\ell \neq j \in [k]$ .

1055 **Lemma I.1** For  $\ell \neq j \in [k]$  let  $\widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j$  be solutions to the least-squares minimization (2.4), i.e.,

$$\begin{aligned} &(\widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j, \widehat{b}_\ell, \widehat{b}_j) \\ &= \arg \min_{\mathbf{w}_\ell, \mathbf{w}_j, b_\ell, b_j} \left\{ \frac{1}{2n} \|\mathbf{D}(\mathbf{Y}_\ell - \mathbf{X}^T \mathbf{w}_\ell - b_\ell \mathbf{1}_n)\|_{\ell_2}^2 + \frac{1}{2n} \|\mathbf{D}(\mathbf{Y}_j - \mathbf{X}^T \mathbf{w}_j - b_j \mathbf{1}_n)\|_{\ell_2}^2 \right\}. \end{aligned}$$

1056 Denote  $\widehat{\mathbf{w}}_{\ell,j} := \widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j$  and  $\widehat{b}_{\ell,j} := \widehat{b}_\ell + \widehat{b}_j$ . Then,  $(\widehat{\mathbf{w}}_{\ell,j}, \widehat{b}_{\ell,j})$  is a minimizer in the following  
 1057 least-squares problem:

$$(\widehat{\mathbf{w}}_{\ell,j}, \widehat{b}_{\ell,j}) = \arg \min_{\mathbf{w}, b} \frac{1}{2n} \|\mathbf{D}(\mathbf{Y}_\ell + \mathbf{Y}_j - \mathbf{X}^T \mathbf{w} - b \mathbf{1}_n)\|_{\ell_2}^2 \quad (\text{I.13})$$

**Proof** Clearly the minimization in (I.13) is convex. Thus, it suffices to prove that  $\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j$  satisfies the KKT conditions. First, by optimality of  $\widehat{\mathbf{w}}_\ell$ , we have that

$$\mathbf{X} \mathbf{D}^2 (\mathbf{Y}_\ell - \mathbf{X}^T \widehat{\mathbf{w}}_\ell - \widehat{b}_\ell \mathbf{1}_n) = 0$$

Similarly, for  $\widehat{\mathbf{w}}_j$ :

$$\mathbf{X} \mathbf{D}^2 (\mathbf{Y}_j - \mathbf{X}^T \widehat{\mathbf{w}}_j - \widehat{b}_j \mathbf{1}_n) = 0.$$

Adding the equations on the above displays we find that

$$\mathbf{X} \mathbf{D}^2 (\mathbf{Y}_\ell + \mathbf{Y}_j - \mathbf{X}^T (\widehat{\mathbf{w}}_j + \widehat{\mathbf{w}}_\ell) - (\widehat{b}_j + \widehat{b}_\ell) \mathbf{1}_n) = 0.$$

1058 Recognize that this coincides with the optimality condition for (I.13). Thus, the proof is complete. ■

1059 Thanks to Lemma I.1, we can use the CGMT to characterize the limiting behavior of  $\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j\|_{\ell_2}$ .

1060 Observe that this immediately gives the limit of  $\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j \rangle$  since

$$\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j \rangle = \frac{\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j\|_{\ell_2}^2 - \|\widehat{\mathbf{w}}_\ell\|_{\ell_2}^2 - \|\widehat{\mathbf{w}}_j\|_{\ell_2}^2}{2}. \quad (\text{I.14})$$

The analysis of (I.13) is very similar to that of (I.1). In particular we use the following decomposition

$$\mathbf{w}_{\ell,j} = \sum_{i=1}^r \beta_i \mathbf{u}_i + \beta_0 \mathbf{w}_{\ell,j}^\perp,$$

1061 with  $\|\mathbf{w}_{\ell,j}^\perp\|_{\ell_2} = 1$  and  $\mathbf{U}^T \mathbf{w}_{\ell,j}^\perp = \mathbf{0}_r$ . This allows us to arrive at

$$\begin{aligned} & \min_{\boldsymbol{\beta}} \min_{\mathbf{b}_{\ell,j}} \max_{\alpha_0 \geq 0} \max_{\beta \geq 0} \min_{\tau \geq 0} \\ & \frac{\gamma\beta}{2\tau} \left( \pi_\ell \boldsymbol{\nu}_\ell \left( \frac{\beta}{\tau} \right) + \pi_j \boldsymbol{\nu}_j \left( \frac{\beta}{\tau} \right) + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A} \left( \frac{\beta}{\tau} \right) \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2 \left( \pi_\ell \boldsymbol{\nu}_\ell \left( \frac{\beta}{\tau} \right) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}_j \left( \frac{\beta}{\tau} \right) \mathbf{c}_j \right)^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} \right) \\ & + \frac{\gamma\beta}{2\tau} \sigma^2 \left( \boldsymbol{\pi}^T \boldsymbol{\nu} \left( \frac{\beta}{\tau} \right) \right) \beta_0^2 - \beta_0 \sigma \beta \sqrt{\gamma} + \frac{\beta\tau}{2} \end{aligned}$$

1062 where  $\mathbf{A}(\eta)$  and  $\mathbf{c}_\ell$  are as in (I.3) and  $\boldsymbol{\nu}(\eta)$  is as in (I.4).

1063 Setting the derivative of the above with respect to  $\alpha_0$  to zero we arrive at

$$\frac{\gamma\beta}{\tau} \sigma^2 \left( \boldsymbol{\pi}^T \boldsymbol{\nu} \left( \frac{\beta}{\tau} \right) \right) \beta_0 - \sigma \beta \sqrt{\gamma} = 0 \quad \Rightarrow \quad \beta_0 = \frac{\tau}{\sigma \sqrt{\gamma} \left( \boldsymbol{\pi}^T \boldsymbol{\nu} \left( \frac{\beta}{\tau} \right) \right)}$$

1064 Note that the above objective has the form

$$g \left( \frac{\beta}{\tau} \right) - \beta_0 \sigma \beta \sqrt{\gamma} + \frac{\beta\tau}{2}$$

1065 with

$$\begin{aligned} g(\eta) & := \frac{\eta\gamma}{2} \left( \pi_\ell \boldsymbol{\nu}_\ell(\eta) + \pi_j \boldsymbol{\nu}_j(\eta) + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A}(\eta) \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2 \left( \pi_\ell \boldsymbol{\nu}_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}_j(\eta) \mathbf{c}_j \right)^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} \right) \\ & + \frac{\gamma\eta}{2} \sigma^2 \left( \boldsymbol{\pi}^T \boldsymbol{\nu}(\eta) \right) \beta_0^2. \end{aligned}$$

1066 Thus, the derivatives with respect to  $\beta$  and  $\tau$  to zero we have

$$\frac{1}{\tau} g' \left( \frac{\beta}{\tau} \right) - \beta_0 \sigma \sqrt{\gamma} + \frac{\tau}{2} = 0 \quad \Rightarrow \quad g' \left( \frac{\beta}{\tau} \right) - \beta_0 \sigma \sqrt{\gamma} \tau + \frac{\tau^2}{2} = 0 \quad \Rightarrow \quad g' \left( \frac{\beta}{\tau} \right) = \tau^2 \left( \frac{1}{\boldsymbol{\pi}^T \boldsymbol{\nu} \left( \frac{\beta}{\tau} \right)} - \frac{1}{2} \right)$$

1067 and

$$-\frac{\beta}{\tau^2} g' \left( \frac{\beta}{\tau} \right) + \frac{\beta}{2} = 0 \quad \Rightarrow \quad \tau^2 = 2g' \left( \frac{\beta}{\tau} \right)$$

1068 Combining the latter two we conclude that  $\boldsymbol{\pi}^T \boldsymbol{\nu} \left( \frac{\beta}{\tau} \right) = 1$ . Thus,  $\eta = \frac{\beta}{\tau}$  is the solution to  $\boldsymbol{\pi}^T \boldsymbol{\nu}(\eta) = 1$ .

1069 To calculate  $\tau$  and hence  $\beta_0$  we calculate  $g'$  which is equal to

$$\begin{aligned} g'(\eta) & = \frac{\gamma}{2} \left( \pi_\ell \boldsymbol{\nu}_\ell(\eta) + \pi_j \boldsymbol{\nu}_j(\eta) + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A}(\eta) \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2 \left( \pi_\ell \boldsymbol{\nu}_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}_j(\eta) \mathbf{c}_j \right)^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} \right) \\ & + \frac{\gamma}{2} \sigma^2 \left( \boldsymbol{\pi}^T \boldsymbol{\nu}(\eta) \right) \beta_0^2 + \frac{\gamma\eta}{2} \sigma^2 \beta_0^2 \left( \boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta) \right) \\ & + \frac{\gamma\eta}{2} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_j \boldsymbol{\nu}'_j(\eta) + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A}'(\eta) \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2 \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j \right)^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} \right) \end{aligned}$$

1070 Where

$$\begin{aligned} \boldsymbol{\nu}'(\eta) & = -\frac{1}{\gamma} \begin{bmatrix} \frac{\omega_1^2}{(\omega_1^2 + \eta)^2} \\ \frac{\omega_2^2}{(\omega_2^2 + \eta)^2} \\ \dots \\ \frac{\omega_k^2}{(\omega_k^2 + \eta)^2} \end{bmatrix} \\ \mathbf{A}'(\eta) & := \begin{bmatrix} \sigma^2 \left( \boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta) \right) \mathbf{I} + \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\pi}) \text{diag}(\boldsymbol{\nu}'(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\nu}'(\eta)) \boldsymbol{\pi} \\ \boldsymbol{\pi}^T \text{diag}(\boldsymbol{\nu}'(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta) \end{bmatrix} \\ \mathbf{c}_\ell & := \begin{bmatrix} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{e}_\ell \\ 1 \end{bmatrix} \end{aligned}$$

1071 Now note that at the optimal point we have

$$g'(\eta) = \frac{\tau^2}{2} = \frac{\gamma}{2} \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}(\eta)) \beta_0^2$$

1072 Thus from the above we can conclude that

$$\begin{aligned} \beta_0^2 = & -\frac{1}{\eta \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta))} \left( \pi_\ell \boldsymbol{\nu}_\ell(\eta) + \pi_j \boldsymbol{\nu}_j(\eta) + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A}(\eta) \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2(\pi_\ell \boldsymbol{\nu}_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}_j(\eta) \mathbf{c}_j)^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} \right) \\ & -\frac{1}{\sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta))} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_j \boldsymbol{\nu}'_j(\eta) + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A}'(\eta) \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2(\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j)^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} \right) \end{aligned}$$

1073 Thus the AO optimization problem reduces to

$$\min_{b_{\ell,j}} \min_{\boldsymbol{\beta}} \frac{\eta \gamma}{2} \left( \pi_\ell \boldsymbol{\nu}_\ell + \pi_j \boldsymbol{\nu}_j + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A} \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2(\pi_\ell \boldsymbol{\nu}_\ell \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}_j \mathbf{c}_j)^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} \right)$$

1074 where  $\eta$  is the solution to

$$\sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{\omega_\ell^2 + \eta} = \gamma$$

1075 Thus, similar to (L.10) we can compute the minimizer of the deterministic

$$\begin{bmatrix} \widehat{\boldsymbol{\beta}} \\ \widehat{b}_{\ell,j} \end{bmatrix} = \begin{bmatrix} -\boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\pi_\ell \boldsymbol{\nu}_\ell (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) + \pi_j \boldsymbol{\nu}_j (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_j)) \\ \pi_\ell \boldsymbol{\nu}_\ell + \pi_j \boldsymbol{\nu}_j + (\boldsymbol{\pi} \odot \boldsymbol{\nu})^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\pi_\ell \boldsymbol{\nu}_\ell (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) + \pi_j \boldsymbol{\nu}_j (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_j)) \end{bmatrix} \quad (\text{I.15})$$

1076 and

$$\boldsymbol{\nu}'(\eta) = -\frac{1}{\gamma} \begin{bmatrix} \frac{\omega_1^2}{(\omega_1^2 + \eta)^2} \\ \frac{\omega_2^2}{(\omega_2^2 + \eta)^2} \\ \dots \\ \frac{\omega_k^2}{(\omega_k^2 + \eta)^2} \end{bmatrix}$$

$$\mathbf{A}'(\eta) := \begin{bmatrix} \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta)) \mathbf{I} + \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\pi}) \text{diag}(\boldsymbol{\nu}'(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\nu}'(\eta)) \boldsymbol{\pi} \\ \boldsymbol{\pi}^T \text{diag}(\boldsymbol{\nu}'(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta) \end{bmatrix}$$

1077 To continue note that

$$\begin{aligned} & \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_j \boldsymbol{\nu}'_j(\eta) + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A}'(\eta) \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2(\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j)^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} \\ & = \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_j \boldsymbol{\nu}'_j(\eta) + (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j)^T \mathbf{A}^{-1} \mathbf{A}'(\eta) \mathbf{A}^{-1} (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j) \\ & \quad - 2(\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j)^T \mathbf{A}^{-1} (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j) \end{aligned}$$

1078 Thus,

$$\begin{aligned} \beta_0^2 = & -\frac{1}{\eta \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta))} \left( \pi_\ell \boldsymbol{\nu}_\ell(\eta) + \pi_j \boldsymbol{\nu}_j(\eta) + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A}(\eta) \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2(\pi_\ell \boldsymbol{\nu}_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}_j(\eta) \mathbf{c}_j)^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} \right) \\ & -\frac{1}{\sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta))} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_j \boldsymbol{\nu}'_j(\eta) + [\boldsymbol{\beta}^T \quad b_{\ell,j}] \mathbf{A}'(\eta) \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} - 2(\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j)^T \begin{bmatrix} \boldsymbol{\beta} \\ b_{\ell,j} \end{bmatrix} \right) \\ = & \frac{\zeta}{\sigma^2} \left( \pi_\ell \boldsymbol{\nu}_\ell + \pi_j \boldsymbol{\nu}_j - (\pi_\ell \boldsymbol{\nu}_\ell + \pi_j \boldsymbol{\nu}_j)^2 \right. \\ & \left. - (\pi_\ell \boldsymbol{\nu}_\ell (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) + \pi_j \boldsymbol{\nu}_j (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_j))^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\pi_\ell \boldsymbol{\nu}_\ell (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) + \pi_j \boldsymbol{\nu}_j (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_j)) \right) \\ & + \frac{\zeta \eta}{\sigma^2} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_j \boldsymbol{\nu}'_j(\eta) + (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j)^T \mathbf{A}^{-1} \mathbf{A}'(\eta) \mathbf{A}^{-1} (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j) \right. \\ & \left. - 2(\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j)^T \mathbf{A}^{-1} (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j) \right) \end{aligned}$$

1079 where  $\zeta := \frac{\gamma}{\eta \left( \sum_{\ell=1}^k \frac{\pi_\ell \omega_\ell^2}{(\omega_\ell^2 + \eta)^2} \right)}$ . From the CGMT, we have that  $\|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j\|_{\ell_2}^2 \xrightarrow{P} \widehat{\beta}_0^2 + \|\beta\|_{\ell_2}^2$ . Combining

1080 this with the calculations above, we conclude that

$$\begin{aligned}
& \|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_j\|_{\ell_2}^2 \\
& \xrightarrow{P} \frac{\zeta}{\sigma^2} (\pi_\ell \nu_\ell + \pi_j \nu_j) (1 - \pi_\ell \nu_\ell - \pi_j \nu_j) \\
& + (\pi_\ell \nu_\ell (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) + \pi_j \nu_j (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_j))^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\eta}{\sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T (\pi_\ell \nu_\ell (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell) + \pi_j \nu_j (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_j)) \\
& + \frac{\zeta \eta}{\sigma^2} \left( \pi_\ell \boldsymbol{\nu}'_\ell(\eta) + \pi_j \boldsymbol{\nu}'_j(\eta) + (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j)^T \mathbf{A}^{-1} \mathbf{A}'(\eta) \mathbf{A}^{-1} (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j) \right. \\
& \quad \left. - 2 (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j)^T \mathbf{A}^{-1} (\pi_\ell \boldsymbol{\nu}'_\ell(\eta) \mathbf{c}_\ell + \pi_j \boldsymbol{\nu}'_j(\eta) \mathbf{c}_j) \right) \quad (\text{I.16})
\end{aligned}$$

1081 Finally, using (I.16) and (I.12) in (I.14) it follows that

$$\begin{aligned}
& \langle \mathbf{w}_\ell, \mathbf{w}_j \rangle \xrightarrow{P} \\
& \pi_\ell \nu_\ell \pi_j \nu_j \left( -\frac{\zeta}{\sigma^2} + (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell)^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\zeta}{\sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_j) \right) \\
& + \frac{\zeta \eta}{\sigma^2} \pi_\ell \nu'_\ell \pi_j \nu'_j \left( \mathbf{c}_j^T \mathbf{A}^{-1} \mathbf{A}' \mathbf{A}^{-1} \mathbf{c}_\ell - 2 \mathbf{c}_j^T \mathbf{A}^{-1} \mathbf{c}_\ell \right) \\
& = \pi_\ell \nu_\ell \pi_j \nu_j \left( -\frac{\zeta}{\sigma^2} + (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_\ell)^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\zeta}{\sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T (\boldsymbol{\pi} \odot \boldsymbol{\nu} - \mathbf{e}_j) \right) \\
& + \frac{\zeta \eta}{\sigma^2} \pi_\ell \nu'_\ell \pi_j \nu'_j \left( \mathbf{e}_j^T \left[ \begin{array}{c} \boldsymbol{\Sigma} \mathbf{V}^T \\ \mathbf{1}^T \end{array} \right]^T \mathbf{A}^{-1} \mathbf{A}' \mathbf{A}^{-1} \left[ \begin{array}{c} \boldsymbol{\Sigma} \mathbf{V}^T \\ \mathbf{1}^T \end{array} \right] \mathbf{e}_\ell - 2 \mathbf{e}_j^T \left[ \begin{array}{c} \boldsymbol{\Sigma} \mathbf{V}^T \\ \mathbf{1}^T \end{array} \right]^T \mathbf{A}^{-1} \left[ \begin{array}{c} \boldsymbol{\Sigma} \mathbf{V}^T \\ \mathbf{1}^T \end{array} \right] \mathbf{e}_\ell \right) \quad (\text{I.17})
\end{aligned}$$

1082 Putting everything together we arrive at

$$\boldsymbol{\Sigma}_{\mathbf{w}, \mathbf{w}} \xrightarrow{P} \frac{\zeta}{\sigma^2} \mathbf{P} + \mathbf{P} \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\zeta}{\sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{P} + \frac{\zeta \eta}{\sigma^2} \mathbf{M}$$

1083 where

$$\mathbf{Q} := \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}') + \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}') \left[ \begin{array}{c} \boldsymbol{\Sigma} \mathbf{V}^T \\ \mathbf{1}^T \end{array} \right]^T (\mathbf{A}^{-1} \mathbf{A}' \mathbf{A}^{-1} - 2 \mathbf{A}^{-1}) \left[ \begin{array}{c} \boldsymbol{\Sigma} \mathbf{V}^T \\ \mathbf{1}^T \end{array} \right] \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}')$$

1084 and as mentioned earlier

$$\begin{aligned}
\mathbf{A}' & := \mathbf{A}'(\eta) := \begin{bmatrix} \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta)) \mathbf{I} + \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\pi}) \text{diag}(\boldsymbol{\nu}'(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\nu}'(\eta)) \boldsymbol{\pi} \\ \boldsymbol{\pi}^T \text{diag}(\boldsymbol{\nu}'(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\pi}^T \boldsymbol{\nu}'(\eta) \end{bmatrix} \\
\mathbf{A} & := \mathbf{A}(\eta) := \begin{bmatrix} \sigma^2 (\boldsymbol{\pi}^T \boldsymbol{\nu}(\eta)) \mathbf{I} + \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\pi}) \text{diag}(\boldsymbol{\nu}(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\boldsymbol{\nu}(\eta)) \boldsymbol{\pi} \\ \boldsymbol{\pi}^T \text{diag}(\boldsymbol{\nu}(\eta)) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\pi}^T \boldsymbol{\nu}(\eta) \end{bmatrix}
\end{aligned}$$

1085 Let us end by simplifying  $\mathbf{M}$  to this aim define  $\widetilde{\boldsymbol{\pi}}' = \boldsymbol{\pi} \odot \boldsymbol{\nu}'$  note that

$$\begin{aligned}
\mathbf{A}^{-1} \left[ \begin{array}{c} \boldsymbol{\Sigma} \mathbf{V}^T \\ \mathbf{1}^T \end{array} \right] \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}') & = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\widetilde{\boldsymbol{\pi}}'^T \mathbf{V} \boldsymbol{\Sigma} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Delta}^{-1} & \mathbf{0} \\ \mathbf{0}^T & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & -\boldsymbol{\Sigma} \mathbf{V}^T \widetilde{\boldsymbol{\pi}}' \\ \mathbf{0}^T & \mathbf{1} \end{bmatrix} \left[ \begin{array}{c} \boldsymbol{\Sigma} \mathbf{V}^T \\ \mathbf{1}^T \end{array} \right] \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}') \\
& = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\widetilde{\boldsymbol{\pi}}'^T \mathbf{V} \boldsymbol{\Sigma} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Delta}^{-1} & \mathbf{0} \\ \mathbf{0}^T & \mathbf{1} \end{bmatrix} \left[ \begin{array}{c} \boldsymbol{\Sigma} \mathbf{V}^T (\mathbf{I} - \widetilde{\boldsymbol{\pi}}' \mathbf{1}^T) \\ \mathbf{1}^T \end{array} \right] \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}') \\
& = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\widetilde{\boldsymbol{\pi}}'^T \mathbf{V} \boldsymbol{\Sigma} & \mathbf{1} \end{bmatrix} \left[ \begin{array}{c} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\mathbf{I} - \widetilde{\boldsymbol{\pi}}' \mathbf{1}^T) \\ \mathbf{1}^T \end{array} \right] \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}') \\
& = \left[ \begin{array}{c} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\mathbf{I} - \widetilde{\boldsymbol{\pi}}' \mathbf{1}^T) \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}') \\ -\widetilde{\boldsymbol{\pi}}'^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\mathbf{I} - \widetilde{\boldsymbol{\pi}}' \mathbf{1}^T) + \mathbf{1}^T \end{array} \right] \text{diag}(\boldsymbol{\pi} \odot \boldsymbol{\nu}') \\
& = \left[ \begin{array}{c} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\widetilde{\boldsymbol{\pi}}') - \widetilde{\boldsymbol{\pi}}' \boldsymbol{\pi}'^T) \\ -\widetilde{\boldsymbol{\pi}}'^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\widetilde{\boldsymbol{\pi}}') - \widetilde{\boldsymbol{\pi}}' \boldsymbol{\pi}'^T) + \widetilde{\boldsymbol{\pi}}'^T \end{array} \right]
\end{aligned}$$



1086 Using the above we arrive at

$$\begin{aligned} \mathbf{Q} &= \text{diag}(\tilde{\boldsymbol{\pi}}') \\ &+ \left[ \begin{array}{c} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\tilde{\boldsymbol{\pi}}') - \tilde{\boldsymbol{\pi}} \tilde{\boldsymbol{\pi}}'^T) \\ -\tilde{\boldsymbol{\pi}}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\tilde{\boldsymbol{\pi}}') - \tilde{\boldsymbol{\pi}} \tilde{\boldsymbol{\pi}}'^T) + \tilde{\boldsymbol{\pi}}'^T \end{array} \right]^T (\mathbf{A}' - 2\mathbf{A}) \left[ \begin{array}{c} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\tilde{\boldsymbol{\pi}}') - \tilde{\boldsymbol{\pi}} \tilde{\boldsymbol{\pi}}'^T) \\ -\tilde{\boldsymbol{\pi}}^T \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\tilde{\boldsymbol{\pi}}') - \tilde{\boldsymbol{\pi}} \tilde{\boldsymbol{\pi}}'^T) + \tilde{\boldsymbol{\pi}}'^T \end{array} \right] \end{aligned}$$

1087 where

$$\begin{aligned} \mathbf{A}' &= \begin{bmatrix} \sigma^2 (\tilde{\boldsymbol{\pi}}'^T \mathbf{1}) \mathbf{I} + \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\tilde{\boldsymbol{\pi}}') \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\tilde{\boldsymbol{\pi}}') \\ \tilde{\boldsymbol{\pi}}'^T \mathbf{V} \boldsymbol{\Sigma} & \tilde{\boldsymbol{\pi}}'^T \mathbf{1} \end{bmatrix} \\ \mathbf{A} &= \begin{bmatrix} \sigma^2 \mathbf{I} + \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\tilde{\boldsymbol{\pi}}) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\Sigma} \mathbf{V}^T \text{diag}(\tilde{\boldsymbol{\pi}}) \\ \tilde{\boldsymbol{\pi}}^T \mathbf{V} \boldsymbol{\Sigma} & 1 \end{bmatrix} \end{aligned}$$

1088 Using the above the cross correlation matrix is given by

$$\boldsymbol{\Sigma}_{\mathbf{w}, \mathbf{w}} \xrightarrow{P} \frac{\zeta}{\sigma^2} \mathbf{P} + \mathbf{P} \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \left( \boldsymbol{\Delta}^{-1} - \frac{\zeta}{\sigma^2} \mathbf{I}_r \right) \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{P} + \frac{\zeta \eta}{\sigma^2} \mathbf{Q}$$

## 1089 J Weighted LS for MLM (Proof of Theorem 4.4)

1090 Let  $\mathbf{D} := \mathbf{D}^{(n)} := \text{diag}(d_1, \dots, d_n)$  be a diagonal matrix with non-zero diagonal entries. In particular,  
1091 assume that the entries of  $\mathbf{D}$  are distributed  $D_i \stackrel{iid}{\sim} D$  where the random variable  $D$  may depend on  
1092 the entries of the matrix of response variables  $\mathbf{Y}$ . Here, we focus on the following setting:

$$\mathbf{D} = \sum_{j \in [k]} \text{diag}(\omega_j \mathbf{Y}_\ell), \quad \omega_j \geq 0, \quad j \in [k]. \quad (\text{J.1})$$

1093 Specifically, for (J.1), we have  $D_i \stackrel{iid}{\sim} D$  with  $D = \omega_\ell Y_\ell + \sum_{i \neq \ell \in [k]} \omega_i Y_i$ , where for all  $c \in [k]$ :

$$\mathbb{P}([Y_1, Y_2, \dots, Y_k]^T = \mathbf{e}_c) = V_c = \frac{e^{\mathbf{e}_c^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{g}}}{\sum_{\ell'=1}^k e^{\mathbf{e}_{\ell'}^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{g}}}. \quad (\text{J.2})$$

1094 With these, we consider the weighted least-squares (WLS) solution for all  $\ell \in [k]$ :

$$(\widehat{\mathbf{w}}_\ell, \widehat{b}) = \arg \min_{\mathbf{w}, b} \mathcal{L}_{PO}(\mathbf{w}, b) := \frac{1}{2n} \|\mathbf{D}(\mathbf{X}^T \mathbf{w} + b \mathbf{1}_n - \mathbf{Y}_\ell)\|_{\ell_2}^2,$$

1095 where  $\mathbf{D}$  is as in (J.1). In fact, it is convenient to rewrite the above as follows:

$$(\widehat{\mathbf{w}}_\ell, \widehat{b}) = \arg \min_{\mathbf{w}, b, \mathbf{u}} \max_{\mathbf{s}} \frac{1}{n} \left( \mathbf{s}^T \mathbf{D} \mathbf{X}^T \mathbf{w} + b \mathbf{s}^T \mathbf{D} \mathbf{1}_n \mathbf{s}^T \mathbf{D} \mathbf{Y}_\ell - \mathbf{s}^T \mathbf{u} + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2} \right). \quad (\text{J.3})$$

1096 **Identifying the AO.** The PO in (J.3) is very similar to (H.1). In particular, following step by step  
1097 the same decomposition trick as in Section G.1.1, it can be shown that the AO corresponding to (J.3)  
1098 becomes (cf. (H.6))

$$\min_{\mathbf{w}_\ell, b_\ell, \mathbf{u}} \max_{\mathbf{s}} \frac{1}{n} \left( \|\mathbf{P}^\perp \mathbf{w}_\ell\|_{\ell_2} \mathbf{g}^T \mathbf{D} \mathbf{s} + \|\mathbf{D} \mathbf{s}\|_{\ell_2} \mathbf{h}^T \mathbf{P}^\perp \mathbf{w}_\ell + \mathbf{s}^T \mathbf{D} \widetilde{\mathbf{G}}^T \mathbf{U}^T \mathbf{w}_\ell + b_\ell \mathbf{s}^T \mathbf{D} \mathbf{1}_n - \mathbf{s}^T \mathbf{D} \mathbf{Y}_\ell - \mathbf{u}^T \mathbf{s} + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2} \right),$$

1099 where we use the same notation as in Section G.1.1 for  $\mathbf{P}^\perp, \mathbf{U}, \widetilde{\mathbf{G}}, \mathbf{g}$  and  $\mathbf{h}$ . Recall also the relation  
1100 of  $\mathbf{Y}_\ell$  to  $\widetilde{\mathbf{G}}$  in (H.5).

1101 **Scalarization of the AO.** We start the process of simplifying the AO by setting  $\beta := \|\mathbf{D} \mathbf{s}\|_{\ell_2} / \sqrt{n}$   
1102 and optimizing over the direction of  $\mathbf{D} \mathbf{s}$  to equivalently write the AO as

$$\min_{\mathbf{w}_\ell, b_\ell, \mathbf{u}} \max_{\beta \geq 0} \frac{1}{\sqrt{n}} \left( \beta \left( \|\mathbf{P}^\perp \mathbf{w}_\ell\|_{\ell_2} \mathbf{g} + \widetilde{\mathbf{G}}^T \mathbf{U}^T \mathbf{w}_\ell + b_\ell \mathbf{1}_n - \mathbf{Y}_\ell - \mathbf{D}^{-1} \mathbf{u} \right)_{\ell_2} + \beta \mathbf{h}^T \mathbf{P}^\perp \mathbf{w}_\ell \right) + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2n}, \quad (\text{J.4})$$

Next, focus on the minimization over  $\mathbf{w}_\ell$ . Let us denote

$$\mathbf{a} := \mathbf{U}^T \mathbf{w}_\ell \quad \text{and} \quad \alpha_0 = \|\mathbf{P}^\perp \mathbf{w}_\ell\|_{\ell_2}.$$

1103 Notice that  $\mathbf{a} \perp \mathbf{P}^\perp \mathbf{w}_\ell$  and thus the orthogonal decomposition  $\mathbf{w}_\ell = \mathbf{U} \mathbf{a} + \mathbf{P}^\perp \mathbf{w}_\ell$ . With this  
 1104 observation, note that the optimal direction of  $\mathbf{P}^T \mathbf{w}_\ell$  in (J.4) aligns with  $\mathbf{P}^T \mathbf{h}$  for all values of  $\beta$ .  
 1105 Therefore, (J.4) reduces to

$$\min_{\mathbf{a}, \alpha_0 \geq 0, \mathbf{b}_\ell, \mathbf{u}} \max_{\beta \geq 0} \frac{1}{\sqrt{n}} \left( \beta \|\alpha_0 \mathbf{g} + \tilde{\mathbf{G}}^T \mathbf{a} + \mathbf{b}_\ell \mathbf{1}_n - \mathbf{Y}_\ell - \mathbf{D}^{-1} \mathbf{u}\|_{\ell_2} - \beta \alpha_0 \|\mathbf{P}^\perp \mathbf{h}\|_{\ell_2} \right) + \frac{\|\mathbf{u}\|_{\ell_2}^2}{2n}, \quad (\text{J.5})$$

Continuing let us denote  $\mathbf{t} := \alpha_0 \mathbf{g} + \tilde{\mathbf{G}}^T \mathbf{a} + \mathbf{b}_\ell \mathbf{1}_n - \mathbf{Y}_\ell$  for convenience and rewrite  $\|\mathbf{t} - \mathbf{D}^{-1} \mathbf{u}\|_{\ell_2}$  as follows

$$\frac{\|\mathbf{t} - \mathbf{D}^{-1} \mathbf{u}\|_{\ell_2}}{\sqrt{n}} = \min_{\tau > 0} \frac{\tau}{2} + \frac{\|\mathbf{t} - \mathbf{D}^{-1} \mathbf{u}\|_{\ell_2}^2}{2\tau n}.$$

1106 Note that the resulting minimization is convex in  $\mathbf{u}$  and concave in  $\beta$ . Also, by considering the  
 1107 bounded AO (such that  $\beta$  is bounded; see [15, Sec. A]), we can flip the order of min-max and optimize  
 1108 over  $\mathbf{u}$  first. In particular,  $\mathbf{u}$  minimizes the following strictly convex quadratic

$$\min_{\mathbf{u}} \left\{ \frac{1}{n} \left( \frac{\beta}{2\tau} \|\mathbf{D}^{-1} \mathbf{u}\|_{\ell_2} + \frac{1}{2} \|\mathbf{u}\|_{\ell_2}^2 - \frac{\beta}{\tau} \mathbf{t}^T \mathbf{D}^{-1} \mathbf{u} \right) = \frac{1}{2n} \mathbf{u}^T \left( \frac{\beta}{\tau} \mathbf{D}^{-2} + \mathbf{I}_n \right) \mathbf{u} - \frac{\beta}{\tau n} \mathbf{t}^T \mathbf{D}^{-1} \mathbf{u} \right\}.$$

In particular,

$$\mathbf{u} = \frac{\beta}{\tau} \left( \frac{\beta}{\tau} \mathbf{D}^{-2} + \mathbf{I} \right)^{-1} \mathbf{D}^{-1} \mathbf{t} = \left( \mathbf{D}^{-1} + \frac{\tau}{\beta} \mathbf{D} \right)^{-1} (\alpha_0 \mathbf{g} + \tilde{\mathbf{G}}^T \mathbf{a} + \mathbf{b}_\ell \mathbf{1}_n - \mathbf{Y}_\ell)$$

1109 Putting things together, the new objective function of (J.5) becomes

$$\min_{\mathbf{a}, \alpha_0 \geq 0, \mathbf{b}_\ell, \tau > 0} \max_{\beta \geq 0} \mathcal{R}(\mathbf{a}, \alpha_0, \mathbf{b}_\ell, \tau, \beta) \quad (\text{J.6})$$

$$\text{where } \mathcal{R}(\mathbf{a}, \alpha_0, \mathbf{b}_\ell, \tau, \beta) := \frac{\beta\tau}{2n} + \frac{\beta}{2\tau n} \|\mathbf{t}\|_{\ell_2}^2 - \frac{\beta}{2\tau n} \mathbf{t}^T \left( \mathbf{I} + \frac{\tau}{\beta} \mathbf{D}^2 \right)^{-1} \mathbf{t} - \frac{\beta\alpha_0}{\sqrt{n}} \|\mathbf{P}^\perp \mathbf{h}\|_{\ell_2}.$$

1110 **Convergence of the AO** After having simplified the AO into an optimization problem over  $r + 4$   
 1111 variables, we are ready to study its asymptotic behavior. First, we argue on point-wise convergence  
 1112 of  $\mathcal{R}$  in (J.6). Fix  $\mathbf{a}, \alpha_0, \mathbf{b}_\ell, \tau$  and  $\beta$ . From the WLLN,  $\frac{1}{\sqrt{n}} \|\mathbf{P}^\perp \mathbf{h}\|_{\ell_2} \xrightarrow{P} \sqrt{\gamma}$  and as in (H.9)

$$\frac{1}{n} \|\mathbf{t}\|_{\ell_2}^2 = \frac{1}{n} \sum_{i=1}^n (\alpha_0 \mathbf{g}_i + \mathbf{a}^T \tilde{\mathbf{g}}_i + \mathbf{b}_\ell - [\mathbf{Y}_\ell]_i)^2 \xrightarrow{P} \mathbb{E} \left[ (\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2 \right],$$

1113 where the expectation is over  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)$  (with some abuse of notation) and

$$Y_\ell \sim \text{Bern}(V_\ell) \quad \text{and} \quad V_\ell = \frac{e^{\mathbf{e}_\ell^T \mathbf{V} \Sigma \mathbf{g}}}{\sum_{\ell'=1}^r e^{\mathbf{e}_{\ell'}^T \mathbf{V} \Sigma \mathbf{g}}}. \quad (\text{J.7})$$

1114 Furthermore,

$$\frac{1}{n} \mathbf{t}^T \left( \mathbf{I} + \frac{\tau}{\beta} \mathbf{D}^2 \right)^{-1} \mathbf{t} = \frac{1}{n} \sum_{i=1}^n \frac{(\alpha_0 \mathbf{g}_i + \mathbf{a}^T \tilde{\mathbf{g}}_i + \mathbf{b}_\ell - [\mathbf{Y}_\ell]_i)^2}{1 + \frac{\tau}{\beta} d_i^2} \xrightarrow{P} \mathbb{E} \left[ \frac{(\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2}{1 + \frac{\tau}{\beta} D^2} \right]$$

1115 Therefore, point-wise on  $\mathbf{a}, \alpha_0, \mathbf{b}_\ell, \tau$  and  $\beta$ , the objective  $\mathcal{R}$  of the AO converges to

$$\begin{aligned} \mathcal{D}_\ell(\alpha_0, \alpha, \mathbf{b}_\ell, \tau, \beta) &:= \frac{\beta\tau}{2} + \frac{\beta}{2\tau} \mathbb{E} \left[ (\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2 \right] - \frac{\beta}{2\tau} \mathbb{E} \left[ \frac{(\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2}{1 + \frac{\tau}{\beta} D^2} \right] - \beta \alpha_0 \sqrt{\gamma} \\ &= \frac{\beta\tau}{2} + \frac{1}{2} \mathbb{E} \left[ \frac{D^2 (\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2}{1 + \frac{\tau}{\beta} D^2} \right] - \beta \alpha_0 \sqrt{\gamma} \\ &= \frac{\beta\tau}{2} + \frac{1}{2} \mathbb{E} \left[ \frac{(\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2}{D^{-2} + (\tau/\beta)} \right] - \beta \alpha_0 \sqrt{\gamma}. \end{aligned} \quad (\text{J.8})$$

1116 We note that the function above is jointly convex in  $(\alpha_0, \alpha, \mathbf{b}_\ell, \tau)$  and concave in  $\beta$ .

1117 **J.1 Deterministic Analysis of the AO**

1118 It can be checked that the first order optimality conditions of  $\mathcal{D}_\ell(\alpha_0, \boldsymbol{\alpha}, b_\ell, \tau, \beta)$  with respect to  $\beta$   
 1119 and  $\tau > 0$  are given as follows:

$$\beta^2 = \mathbb{E} \left[ \frac{(\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2}{(D^{-2} + (\tau/\beta))^2} \right] \quad \text{or} \quad \beta = 0, \quad (\text{J.9})$$

$$\alpha_0 \sqrt{\gamma} = \frac{\tau}{2} + \frac{\tau}{2\beta^2} \cdot \mathbb{E} \left[ \frac{(\alpha_0 G_0 + \mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2}{(D^{-2} + (\tau/\beta))^2} \right]. \quad (\text{J.10})$$

1120 Thus, at optimality either  $\beta = 0$  or  $\tau = \alpha_0 \sqrt{\gamma}$ . In what follows, consider the solution  $\tau = \alpha_0 \sqrt{\gamma}$ . We  
 1121 will show that this leads to the true saddle point of  $\mathcal{D}$ .

1122 Moreover, by denoting  $\eta := \frac{\beta}{\tau}$  and recalling from (J.7) that  $Y_\ell = \text{Bern}(V_\ell)$ , we can express  
 1123  $\mathcal{D}_\ell(\alpha_0, \boldsymbol{\alpha}, b_\ell, \tau, \beta)$  as follows

$$\frac{\beta\tau}{2} + \frac{\alpha_0^2}{2} \mathbb{E} \left[ \frac{1}{D^{-2} + 1/\eta} \right] - \beta\alpha_0\sqrt{\gamma} + \frac{1}{2} [\mathbf{a}^T \quad \mathbf{b}_\ell] \cdot \mathbf{A}(\eta) \cdot \begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix} - \mathbf{c}_\ell^T(\eta) \begin{bmatrix} \mathbf{a}^T \\ \mathbf{b}_\ell \end{bmatrix} + \frac{1}{2} \mathbb{E} \left[ \frac{Y_\ell^2}{D^{-2} + 1/\eta} \right],$$

1124 where

$$\mathbf{A}(\eta) := \begin{bmatrix} \mathbb{E} \left[ \frac{\mathbf{g}\mathbf{g}^T}{D^{-2} + 1/\eta} \right] & \mathbb{E} \left[ \frac{\mathbf{g}}{D^{-2} + 1/\eta} \right] \\ \mathbb{E} \left[ \frac{\mathbf{g}^T}{D^{-2} + 1/\eta} \right] & \mathbb{E} \left[ \frac{1}{D^{-2} + 1/\eta} \right] \end{bmatrix} \quad (\text{J.11a})$$

$$\mathbf{c}_\ell(\eta) := \left[ \mathbb{E} \left[ \frac{\mathbf{g}Y_\ell}{D^{-2} + 1/\eta} \right] \right] \quad (\text{J.11b})$$

1125 we have the following first-order optimality conditions for  $\alpha_0$ ,  $\mathbf{a}$  and  $\mathbf{b}_\ell$ :

$$\begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix} = \mathbf{A}^{-1}(\eta) \cdot \mathbf{c}_\ell(\eta) \quad (\text{J.12})$$

$$\alpha_0 = \beta\sqrt{\gamma} / \mathbb{E} \left[ \frac{1}{D^{-2} + 1/\eta} \right]. \quad (\text{J.13})$$

1126 Rearranging (J.13) and using  $\tau = \alpha_0 \sqrt{\gamma}$  gives the following equation for  $\eta$ :

$$\frac{\alpha_0 \sqrt{\gamma}}{\beta} \mathbb{E} \left[ \frac{1}{D^{-2} + 1/\eta} \right] = \gamma \stackrel{\tau = \alpha_0 \sqrt{\gamma}}{\implies} \mathbb{E} \left[ \frac{1/\eta}{D^{-2} + 1/\eta} \right] = \gamma. \quad (\text{J.14})$$

1127 Thus, the optimal values of  $\mathbf{a}$  and  $\mathbf{b}_\ell$  are found by (J.12) for  $\eta$  the positive solution of the equation in  
 1128 (J.14). To solve for  $\alpha_0$ , we combine (J.13) and (J.9) which leads to

$$\alpha_0^2 \left( \gamma \eta^2 - \mathbb{E} \left[ \left( \frac{1}{D^{-2} + 1/\eta} \right)^2 \right] \right) = \mathbb{E} \left[ \frac{(\mathbf{a}^T \mathbf{g} + \mathbf{b}_\ell - Y_\ell)^2}{(D^{-2} + 1/\eta)^2} \right], \quad (\text{J.15})$$

1129 where we have also used the RHS of (J.14). Next, we specialize these findings to the special structures  
 1130 of the weighting matrix  $\mathbf{D}$  (J.1), respectively.

1131 **Applying weighting (J.1).** Assume (J.1) holds. In this case, Equation (J.14) that determines the  
 1132 value of  $\eta > 0$  becomes

$$F(\eta) := \sum_{i \in [k]} \frac{\pi_i \omega_i^2}{\omega_i^2 + \eta} = \gamma, \quad (\text{J.16})$$

1133 where we have recalled the notation in (4.1)  $\pi_i := \mathbb{E}[V_i] > 0$ ,  $i \in [k]$ . It can be easily checked by  
 1134 direct differentiation that  $\eta \mapsto F$  is strictly decreasing in  $(0, \infty)$ . Also, using  $\sum_{i \in [k]} \pi_i = 1$  the range  
 1135 of  $F$  in  $(0, \infty)$  is  $(0, 1)$ . Thus, it follows that (J.16) has a unique solution for all  $\gamma \in (0, 1)$ .

1136 Also, in this case we can write (J.11) in the following more convenient form:

$$\mathbf{A}(\eta) := \sum_{i \in [k]} \left( \frac{\omega_i^2 \eta}{\omega_i^2 + \eta} \right) \underbrace{\mathbb{E} \left[ \begin{bmatrix} \mathbf{g} \\ 1 \end{bmatrix} [\mathbf{g}^T \ 1] V_i \right]}_{=: \tilde{\mathbf{A}}_i} \quad (\text{J.17})$$

$$\mathbf{c}_\ell(\eta) := \left( \frac{\omega_\ell^2 \eta}{\omega_\ell^2 + \eta} \right) \underbrace{\mathbb{E} \left[ \begin{bmatrix} \mathbf{g} \\ 1 \end{bmatrix} V_\ell \right]}_{=: \tilde{\mathbf{c}}_\ell}. \quad (\text{J.18})$$

1137 For convenience let us define vectors  $\boldsymbol{\nu} := \boldsymbol{\nu}(\eta)$ ,  $\tilde{\boldsymbol{\pi}} = \tilde{\boldsymbol{\pi}}(\eta) \in \mathbb{R}^k$  with entries:

$$\tilde{\pi}_\ell := \pi_i \left( \frac{1}{\gamma} \cdot \frac{\omega_i^2}{\omega_i^2 + \eta} \right) =: \pi_i \cdot \nu_i \quad (\text{J.19})$$

1138 Because of (J.16), notice that  $\tilde{\boldsymbol{\pi}}$  is a probability vector, i.e.  $\tilde{\boldsymbol{\pi}}^T \mathbf{1}_k = 1$ . Moreover, for  $\eta$  satisfying  
1139 (J.16), it also holds that  $\boldsymbol{\pi}^T \boldsymbol{\nu} = 1$ .

1140 With the notation above, it holds

$$\begin{aligned} \mathbf{A}(\eta) &= \gamma \cdot \eta \cdot \begin{bmatrix} \sum_{i \in [k]} \nu_i \cdot \mathbb{E}[V_i \mathbf{g} \mathbf{g}^T] & \sum_{i \in [k]} \nu_i \cdot \mathbb{E}[V_i \mathbf{g}] \\ \sum_{i \in [k]} \nu_i \cdot \mathbb{E}[V_i \mathbf{g}^T] & 1 \end{bmatrix} \\ &= \gamma \cdot \eta \cdot \begin{bmatrix} \sum_{i \in [k]} \nu_i \cdot \mathbb{E}[V_i \mathbf{g} \mathbf{g}^T] & \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \boldsymbol{\nu} \\ \boldsymbol{\nu}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V} \boldsymbol{\Sigma} & 1 \end{bmatrix} \end{aligned} \quad (\text{J.20})$$

$$= \gamma \cdot \eta \cdot \begin{bmatrix} \mathbb{E}[(\boldsymbol{\nu}^T \mathbf{v}) \mathbf{g} \mathbf{g}^T] & \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \boldsymbol{\nu} \\ \boldsymbol{\nu}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V} \boldsymbol{\Sigma} & 1 \end{bmatrix} \quad (\text{J.21})$$

$$\begin{aligned} \mathbf{c}_\ell(\eta) &= \gamma \cdot \eta \cdot \begin{bmatrix} \nu_\ell \mathbb{E}[V_\ell \mathbf{g}] \\ \tilde{\pi}_\ell \end{bmatrix} \\ &= \gamma \cdot \eta \cdot \begin{bmatrix} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \nu_\ell \mathbf{e}_\ell \\ \tilde{\pi}_\ell \end{bmatrix} \end{aligned} \quad (\text{J.22})$$

where we have also used the fact that  $\mathbb{E}[V_i \mathbf{g}] = \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{e}_i$ ,  $i \in [k]$  and recalled the notation

$$\mathbf{v} = [V_1, \dots, V_k]^T.$$

1141 Using (J.30) and (J.31), we conclude from (J.12) the following expressions for  $\mathbf{a}$  and  $\mathbf{b}$ :

$$\mathbf{a} = \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \cdot \boldsymbol{\nu}_\ell \cdot (\mathbf{e}_\ell - \boldsymbol{\pi}_\ell \boldsymbol{\nu}), \quad (\text{J.23})$$

$$\mathbf{b}_\ell = \tilde{\boldsymbol{\pi}}_\ell - \boldsymbol{\nu}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \cdot \boldsymbol{\nu}_\ell \cdot (\mathbf{e}_\ell - \boldsymbol{\pi}_\ell \boldsymbol{\nu}), \quad (\text{J.24})$$

1142 where we defined

$$\boldsymbol{\Delta} = \mathbb{E}[(\boldsymbol{\nu}^T \mathbf{v}) \mathbf{g} \mathbf{g}^T] - \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \boldsymbol{\nu} \boldsymbol{\nu}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V} \boldsymbol{\Sigma} > \mathbf{0}_{r \times r}. \quad (\text{J.25})$$

1143 Finally, we show how to compute  $\alpha_0$  using (J.15). The RHS in (J.15) can be computed as

$$\begin{aligned} & \sum_{i \neq \ell \in [k]} \frac{[\mathbf{a}^T \ \mathbf{b}_\ell] \tilde{\mathbf{A}}_i \begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix}}{(\omega_i^{-2} + 1/\eta)^2} + \frac{[\mathbf{a}^T \ \mathbf{b}_\ell] \tilde{\mathbf{A}}_\ell \begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix} - 2[\mathbf{a}^T \ \mathbf{b}_\ell] \tilde{\mathbf{c}}_\ell + \boldsymbol{\pi}_\ell}{(\omega_\ell^{-2} + 1/\eta)^2} \\ &= \eta^2 \cdot \gamma^2 \cdot \left\{ [\mathbf{a}^T \ \mathbf{b}_\ell] \left( \sum_{i \in [k]} \nu_i^2 \tilde{\mathbf{A}}_i \right) \begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix} - 2[\mathbf{a}^T \ \mathbf{b}_\ell] \boldsymbol{\nu}_\ell^2 \tilde{\mathbf{c}}_\ell + \boldsymbol{\pi}_\ell \boldsymbol{\nu}_\ell^2 \right\}, \end{aligned}$$

where  $\mathbf{a}$ ,  $\mathbf{b}_\ell$  are as in (J.23) and (J.24). Also, note that

$$\mathbb{E} \left[ \left( \frac{1}{D^{-2} + 1/\eta} \right)^2 \right] = \eta^2 \sum_{i \in [k]} \frac{\pi_i \omega_i^4}{(\omega_i^2 + \eta)^2} = \eta^2 \cdot \gamma^2 \cdot \boldsymbol{\pi}^T \text{diag}(\boldsymbol{\nu}) \boldsymbol{\nu} = \eta^2 \cdot \gamma^2 \cdot \tilde{\boldsymbol{\pi}}^T \boldsymbol{\nu}.$$

1144 Put together, we have the following expression for  $\alpha_0$ :

$$\begin{aligned}\alpha_0^2 &= \frac{1}{(1/\gamma - \tilde{\pi}^T \boldsymbol{\nu})} \cdot \left\{ [\mathbf{a}^T \quad \mathbf{b}_\ell] \left( \sum_{i \in [k]} \nu_i^2 \tilde{\mathbf{A}}_i \right) \begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix} - 2 [\mathbf{a}^T \quad \mathbf{b}_\ell] \boldsymbol{\nu}_\ell^2 \tilde{\mathbf{c}}_\ell + \pi_\ell \boldsymbol{\nu}_\ell^2 \right\} \\ &= \frac{1}{(1/\gamma - \tilde{\pi}^T \boldsymbol{\nu})} \cdot \left\{ [\mathbf{a}^T \quad \mathbf{b}_\ell] \mathbf{A}' \begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix} - 2 [\mathbf{a}^T \quad \mathbf{b}_\ell] \left[ \frac{\boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \boldsymbol{\nu}_\ell^2 \mathbf{e}_\ell}{\tilde{\pi}_\ell \cdot \boldsymbol{\nu}_\ell} \right] + \tilde{\pi}_\ell \cdot \boldsymbol{\nu}_\ell \right\},\end{aligned}\quad (\text{J.26})$$

1145 where  $\mathbf{a}, \mathbf{b}_\ell$  are as in (J.23), (J.24) and we have also defined

$$\mathbf{A}' = \begin{bmatrix} \mathbb{E}[(\boldsymbol{\nu}^T \text{diag}(\boldsymbol{\nu}) \mathbf{v}) \mathbf{g} \mathbf{g}^T] & \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \text{diag}(\boldsymbol{\nu}) \boldsymbol{\nu} \\ \boldsymbol{\nu}^T \text{diag}(\boldsymbol{\nu}) (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V} \boldsymbol{\Sigma} & \boldsymbol{\nu}^T \text{diag}(\boldsymbol{\nu}) \boldsymbol{\pi} \end{bmatrix}. \quad (\text{J.27})$$

1146 **Asymptotic Predictions.** Writing (J.24) in vector form we find that

$$\widehat{\mathbf{b}} \xrightarrow{P} \tilde{\boldsymbol{\pi}} - \text{diag}(\boldsymbol{\nu}) (\mathbf{I}_k - \boldsymbol{\pi} \boldsymbol{\nu}^T) (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \boldsymbol{\nu}. \quad (\text{J.28})$$

1147 Also, recalling that  $\mathbf{e}_\ell^T \boldsymbol{\Sigma} \mathbf{w}, \boldsymbol{\mu} = \widehat{\mathbf{w}}_\ell \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T \xrightarrow{P} \mathbf{a} \boldsymbol{\Sigma} \mathbf{V}^T$  and using (J.23):

$$\boldsymbol{\Sigma} \mathbf{w}, \boldsymbol{\mu} \xrightarrow{P} \text{diag}(\boldsymbol{\nu}) (\mathbf{I}_k - \boldsymbol{\pi} \boldsymbol{\nu}^T) (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V} \boldsymbol{\Sigma} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T \quad (\text{J.29})$$

1148 Finally, for the magnitudes of the weight vectors, recall that  $\|\widehat{\mathbf{w}}_\ell\|_{\ell_2}^2 \xrightarrow{P} \|\mathbf{a}\|_{\ell_2}^2 + \alpha_0^2$ . Thus, to find  
1149 the limiting values of the norms, we can combine (J.26) and (J.23)-(J.24). For convenience, we  
1150 summarize the final expression here. Define the following<sup>1</sup>

$$\mathbf{A} := \begin{bmatrix} \mathbb{E}[(\boldsymbol{\nu}^T \mathbf{v}) \mathbf{g} \mathbf{g}^T] & \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \boldsymbol{\nu} \\ \boldsymbol{\nu}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{V} \boldsymbol{\Sigma} & 1 \end{bmatrix} \quad (\text{J.30})$$

$$\mathbf{c}_\ell = \begin{bmatrix} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \boldsymbol{\nu}_\ell \mathbf{e}_\ell \\ \tilde{\pi}_\ell \end{bmatrix} \quad (\text{J.31})$$

1151 Further recall the matrix  $\mathbf{A}'$  in (J.27).

$$\begin{aligned}\|\widehat{\mathbf{w}}_\ell\|_{\ell_2}^2 &\xrightarrow{P} \|\boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma} \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \cdot \boldsymbol{\nu}_\ell \cdot (\mathbf{e}_\ell - \boldsymbol{\pi}_\ell \boldsymbol{\nu})\|_{\ell_2}^2 \\ &\quad + \frac{1}{(1/\gamma - \tilde{\pi}^T \boldsymbol{\nu})} \cdot \left\{ \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{A}' \mathbf{A}^{-1} \mathbf{c}_\ell - 2 \boldsymbol{\nu}_\ell \mathbf{c}_\ell^T \mathbf{A}^{-1} \mathbf{c}_\ell + \tilde{\pi}_\ell \cdot \boldsymbol{\nu}_\ell \right\},\end{aligned}\quad (\text{J.32})$$

1152 **Remark J.2** Consider the special case  $\omega_i = 1$ ,  $i \in [k]$ . We show how the above recovers the  
1153 solution for (un-weighted) LS. First, note that in this case (J.16) simply gives  $\eta = \frac{1}{\gamma} - 1$ . Thus,  
1154  $\boldsymbol{\nu} = \mathbf{1}_k$  and  $\tilde{\boldsymbol{\pi}} = \boldsymbol{\pi}$ . Also, recall that  $(\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\Pi}) \mathbf{1}_k = \mathbf{0}$  and  $\mathbf{1}^T \mathbf{v} = 1$ . Thus, (J.25) simply  
1155 gives  $\boldsymbol{\Delta} = \mathbb{E}[\mathbf{g} \mathbf{g}^T] = \mathbf{I}_r$ . With these, it can be readily checked that (J.28) and (J.23) simplify to  
1156 the expressions in (4.4a). Similarly,  $\mathbf{A} = \mathbf{A}' = \mathbf{I}_{r+1}$  and (J.26) reduces in this case to (H.1.1). For  
1157 general weight coefficients, such simplifications do not seem possible and one needs to compute the  
1158 matrix  $\mathbb{E}[(\boldsymbol{\nu}^T \mathbf{v}) \mathbf{g} \mathbf{g}^T]$  that appears in the definitions of  $\boldsymbol{\Delta}$ ,  $\mathbf{A}$  and  $\mathbf{A}'$ . We note that this calculation  
1159 can be somewhat simplified by applying Gaussian integration by parts similar to lemma B.2.

### 1160 J.3 Computing cross-correlations $\boldsymbol{\Sigma}_{w,w}$

1161 In this section, we use Lemma I.1 to compute the cross-correlations  $\langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_j \rangle$ ,  $j \neq \ell \in [k]$ . Specif-  
1162 ically, the analysis of (I.13) is almost identical to the analysis of (J.3) in the previous section.  
1163 Specifically, without repeating all the details for brevity, it can be shown that the AO of (I.13) con-  
1164 verges to  $\min_{\mathbf{a}, \alpha_0 \geq 0, b_\ell, \tau > 0} \max_{\beta \geq 0} \mathcal{D}(\mathbf{a}, \alpha_0, b_\ell, \tau, \beta)$  where  $\mathcal{D}(\mathbf{a}, \alpha_0, b_\ell, \tau, \beta)$  is as in (J.8) only  
1165 with  $Y_\ell$  substituted by  $Y_{\ell,c}$ :

$$Y_{\ell,c} \sim \text{Bern}(V_c + V_\ell) \quad \text{and as before:} \quad V_i = \frac{e^{e_i^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{g}}}{\sum_{\ell'=1}^r e^{e_{\ell'}^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{g}}}, \quad i = \ell, c. \quad (\text{J.33})$$

<sup>1</sup>Note the slight abuse of notation compared to the definitions in (J.30) and (J.30). This ‘‘renaming’’ should not be confusing as the constant  $\gamma \cdot \eta$  (that is different between the two definitions) cancels in the evaluation of  $\begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix} = \mathbf{A}^{-1} \mathbf{c}_\ell$  (see (J.12)).

1166 Thus, what changes in the calculations above is in (J.18) and (J.35), where we now have instead

$$\mathbf{c}(\eta) := \left( \frac{1}{\frac{1}{\omega_\ell^2} + 1/\eta} \right) \underbrace{\mathbb{E} \left[ \begin{bmatrix} \mathbf{g} \\ 1 \end{bmatrix} V_\ell \right]}_{=:\tilde{\mathbf{c}}_\ell} + \left( \frac{1}{\frac{1}{\omega_c^2} + 1/\eta} \right) \underbrace{\mathbb{E} \left[ \begin{bmatrix} \mathbf{g} \\ 1 \end{bmatrix} V_c \right]}_{=:\tilde{\mathbf{c}}_c} \quad (\text{J.34})$$

1167 and

$$\begin{aligned} \sum_{i \neq \{\ell, c\} \in [k]} \frac{[\mathbf{a}^T \ \mathbf{b}_\ell] \tilde{\mathbf{A}}_i \begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix}}{(\omega_i^{-2} + 1/\eta)^2} + \frac{[\mathbf{a}^T \ \mathbf{b}_\ell] \tilde{\mathbf{A}}_\ell \begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix} - 2[\mathbf{a}^T \ \mathbf{b}_\ell] \tilde{\mathbf{c}}_\ell + \pi_\ell}{(\omega_\ell^{-2} + \eta)^2} \\ + \frac{[\mathbf{a}^T \ \mathbf{b}_\ell] \tilde{\mathbf{A}}_c \begin{bmatrix} \mathbf{a} \\ \mathbf{b}_\ell \end{bmatrix} - 2[\mathbf{a}^T \ \mathbf{b}_\ell] \tilde{\mathbf{c}}_c + \pi_c}{(\omega_c^{-2} + 1/\eta)^2}, \end{aligned} \quad (\text{J.35})$$

1168 respectively. With these and following mutatis-mutandis the steps and the notation in the previous  
1169 section, we find the following asymptotic expression for the magnitude of  $\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_c$ :

$$\begin{aligned} \|\widehat{\mathbf{w}}_\ell + \widehat{\mathbf{w}}_c\|_{\ell_2}^2 \xrightarrow{P} \|\Delta^{-1} \Sigma \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \cdot (\boldsymbol{\nu}_\ell \cdot (\mathbf{e}_\ell - \boldsymbol{\pi}_\ell \boldsymbol{\nu}) + \boldsymbol{\nu}_c \cdot (\mathbf{e}_c - \boldsymbol{\pi}_c \boldsymbol{\nu}))\|_{\ell_2}^2 \\ + \frac{1}{(1/\gamma - \tilde{\boldsymbol{\pi}}^T \boldsymbol{\nu})} \cdot \left\{ (\mathbf{c}_\ell + \mathbf{c}_c)^T \mathbf{A}^{-1} \mathbf{A}' \mathbf{A}^{-1} (\mathbf{c}_\ell + \mathbf{c}_c) - 2(\boldsymbol{\nu}_\ell \mathbf{c}_\ell + \boldsymbol{\nu}_c \mathbf{c}_c)^T \mathbf{A}^{-1} (\boldsymbol{\nu}_\ell \mathbf{c}_\ell + \boldsymbol{\nu}_c \mathbf{c}_c) + \tilde{\boldsymbol{\pi}}_\ell \cdot \boldsymbol{\nu}_\ell + \tilde{\boldsymbol{\pi}}_c \cdot \boldsymbol{\nu}_c \right\}, \end{aligned} \quad (\text{J.36})$$

1170 We may now combine this with (J.32) to conclude with the following asymptotic limits for the  
1171 cross-correlations for all  $\ell \neq c \in [k]$ :

$$\begin{aligned} \langle \widehat{\mathbf{w}}_\ell, \widehat{\mathbf{w}}_c \rangle \xrightarrow{P} \boldsymbol{\nu}_c (\mathbf{e}_c - \boldsymbol{\pi}_c \boldsymbol{\nu})^T (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \mathbf{V} \Sigma \Delta^{-2} \Sigma \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \boldsymbol{\nu}_\ell (\mathbf{e}_\ell - \boldsymbol{\pi}_\ell \boldsymbol{\nu}) \\ + \frac{1}{(1/\gamma - \tilde{\boldsymbol{\pi}}^T \boldsymbol{\nu})} \cdot \left\{ \mathbf{c}_c^T \mathbf{A}^{-1} \mathbf{A}' \mathbf{A}^{-1} \mathbf{c}_\ell - \boldsymbol{\nu}_c \boldsymbol{\nu}_\ell \mathbf{c}_c^T \mathbf{A}^{-1} \mathbf{c}_\ell \right\}. \end{aligned} \quad (\text{J.37})$$

$$\begin{aligned} = \boldsymbol{\nu}_c (\mathbf{e}_c - \boldsymbol{\pi}_c \boldsymbol{\nu})^T (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \mathbf{V} \Sigma \Delta^{-2} \Sigma \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \boldsymbol{\nu}_\ell (\mathbf{e}_\ell - \boldsymbol{\pi}_\ell \boldsymbol{\nu}) \\ + \frac{1}{(1/\gamma - \tilde{\boldsymbol{\pi}}^T \boldsymbol{\nu})} \cdot \left\{ \mathbf{c}_c^T (\mathbf{A}^{-1} \mathbf{A}' \mathbf{A}^{-1} - \boldsymbol{\nu}_c \boldsymbol{\nu}_\ell \mathbf{A}^{-1}) \mathbf{c}_\ell \right\}. \end{aligned} \quad (\text{J.38})$$

1172 In matrix form, we have

$$\begin{aligned} \boldsymbol{\Sigma}_{w,w} \xrightarrow{P} \text{diag}(\boldsymbol{\nu}) (\mathbf{I}_k - \boldsymbol{\pi} \boldsymbol{\nu}^T) (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \mathbf{V} \Sigma \Delta^{-2} \Sigma \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) (\mathbf{I}_k - \boldsymbol{\pi} \boldsymbol{\nu}^T) \text{diag}(\boldsymbol{\nu}) \\ + \frac{1}{(1/\gamma - \tilde{\boldsymbol{\pi}}^T \boldsymbol{\nu})} \left\{ [\text{diag}(\boldsymbol{\nu}) (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \mathbf{V} \Sigma \quad \tilde{\boldsymbol{\pi}}] \mathbf{A}^{-1} \mathbf{A}' \mathbf{A}^{-1} \begin{bmatrix} \text{diag}(\boldsymbol{\nu}) \Sigma \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \\ \tilde{\boldsymbol{\pi}}^T \end{bmatrix} \right. \\ \left. + \text{diag}(\boldsymbol{\nu}) [\text{diag}(\boldsymbol{\nu}) (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \mathbf{V} \Sigma \quad \tilde{\boldsymbol{\pi}}] \mathbf{A}^{-1} \begin{bmatrix} \text{diag}(\boldsymbol{\nu}) \Sigma \mathbf{V}^T (\text{diag}(\boldsymbol{\pi}) - \mathbf{\Pi}) \\ \tilde{\boldsymbol{\pi}}^T \end{bmatrix} \text{diag}(\boldsymbol{\nu}) \right\} \end{aligned} \quad (\text{J.39})$$