
Generalization for Multiclass Classification with Overparameterized Linear Models

Vignesh Subramanian

Department of Electrical Engineering and Computer Sciences
University of California Berkeley
Berkeley, CA-94720, USA
vignesh.subramanian@eecs.berkeley.edu

Rahul Arya

Department of Electrical Engineering and Computer Sciences
University of California Berkeley
Berkeley, CA-94720, USA
rahularya@berkeley.edu

Anant Sahai

Department of Electrical Engineering and Computer Sciences
University of California Berkeley
Berkeley, CA-94720, USA
sahai@eecs.berkeley.edu

Abstract

Via an overparameterized linear model with Gaussian features, we provide conditions for good generalization for multiclass classification of minimum-norm interpolating solutions in an asymptotic setting where both the number of underlying features and the number of classes scale with the number of training points. The survival/contamination analysis framework for understanding the behavior of overparameterized learning problems is adapted to this setting, revealing that multiclass classification qualitatively behaves like binary classification in that, as long as there are not too many classes (made precise in the paper), it is possible to generalize well even in some settings where the corresponding regression tasks would not generalize. Besides various technical challenges, it turns out that the key difference from the binary classification setting is that there are relatively fewer positive training examples of each class in the multiclass setting as the number of classes increases, making the multiclass problem “harder” than the binary one.

1 Introduction

Multiclass classification on standardized datasets is where the current deep-learning revolution really made the community take notice with previously unattainable levels of performance. Contemporary systems have demonstrated tremendous success at these tasks, typically using gigantic models with parameters that vastly exceed the (also large) number of data points used to train these models. In defiance of traditional statistical wisdom regarding overfitting, these big models can be trained to achieve zero training error even with noisy labels, but still generalize well in practice [84, 28].

To better understand this empirical phenomenon, one line of work uses appropriate high-dimensional linear models for regression problems to show how benign fitting of noise in training data is possible [31, 52, 4, 9, 55]. Essentially, the model must have enough “non-preferred” degrees of freedom to be able to absorb the training noise without contaminating predictions by too much. Simultaneously, there has to be enough of a preference for degrees of freedom that can capture the true pattern to enable it to survive the learning procedure and be well represented in the final learned model.

A subsequent line of work studies binary classification [56, 13, 76] and shows that binary classification can generalize well beyond what can be proved by classical margin-based bounds [3] and there exist regimes where binary classification can even succeed in generalizing where regression fails — less preference is required for the degrees of freedom that capture the true pattern [56]. Very recently, the generalization of multiclass classification in similar models was studied in Wang et al. [77] but the analysis was limited to a fixed finite number of classes. In practice, we see that larger datasets often come with more classes and are tackled with even bigger models and so it is important to see what happens to generalization when everything scales together. To have a crisply understandable approach that allows everything to scale, this paper also adopts the bi-level covariance model with Gaussian features that is used in Muthukumar et al. [55, 56], Wang et al. [75], Wang et al. [77].

To understand classification, we must understand the role of training loss functions in determining what is learned. Empirical evidence shows that least-squares can yield classification performance competitive to cross-entropy minimization [64, 35, 11]. Muthukumar et al. [56], Hsu et al. [33] show that indeed with sufficient overparameterization, the support vector machine (SVM) solution, which also arises from minimizing the logistic loss using gradient descent [68, 36], is identical to that obtained by the minimum-norm interpolation (MNI) of binary labels — what would be obtained by gradient descent while minimizing the squared loss. A similar equivalence¹ holds for different variations of multiclass SVMs and the MNI of one-hot-encoded labels [77]. Consequently, this paper focuses on the MNI approach to overparameterized learning for multiclass classification.

2 Our contributions

Our study provides an asymptotic analysis of the error of the minimum-norm interpolating classifier for the multiclass classification problem with weighted Gaussian features. We consider an overparameterized setting using a bi-level feature weighting model where the number of features, classes, favored features, and the feature weights themselves all scale with the number of training points. Under this model, Theorem 5.1 provides sufficient conditions for good generalization in the form of a region in which as the number of training points increase, the number of classes grows slowly enough, the total number of features (i.e. level of overparameterization) grows fast enough, the number of favored features grows slowly enough, and the amount of favoring of those favored features is sufficient to allow for asymptotic generalization. We assume that our labels are generated noiselessly based on which of the first k features is the largest.²

To prove our main result, Theorem 5.1, we present a novel typicality-style argument featuring the feature margin (gap between the largest and second-largest feature) for computing sufficient conditions for correct classification utilizing the signal-processing inspired concepts of survival and contamination from Muthukumar et al. [55, 56] and leveraging the random-matrix analysis tools sharpened in Bartlett et al. [4]. The survival concept relates to the shrinkage induced by the regularizing effect of having lots of features in the context of min-norm interpolation — survival captures what is left of the true pattern after shrinkage. Contamination reflects the consequence of overparameterization when training via optimization: in addition to the true pattern, there is an infinite family of other³ false patterns (aliases) that also happen to explain the limited training data, and the optimizer ends up hedging its bet across the true pattern and these other competing false explanations. The learned false patterns contaminate the predictions on test points, and this can be quantified by the relevant standard deviation.

The key is analyzing what happens with multiclass training data where there are relatively fewer positive examples of each class, and where the training data for a particular class is not independent of the features corresponding to other classes. The analysis shows that as a result of having fewer positive exemplars for a class relative to the total size of the training data, the survival drops by a factor of k (the number of classes), while the contamination only drops by a factor of $\frac{1}{k}$. As in binary classification, the ratio of the relevant survival to contamination terms plays the role of the effective signal-to-noise ratio and shows up as a key quantity in our error analysis (Equation (22) from

¹For an interesting alternative perspective on this equivalence as an indication of a potential bug instead of as a promising feature, see Shamir [67].

²This assumption is without loss of generality for the bi-level model as long as the classes are defined by orthogonal directions as in Wang et al. [77].

³This is related to what is called the challenge of “underspecification” in ML [19], and this in turn is also one aspect of the challenge of covariate shifts [73].

Section 5.1). When this ratio grows asymptotically to ∞ , multiclass classification generalizes well. To the best of our knowledge, this is the first work that quantifies this effect of fewer informative samples per class and in what sense that makes multiclass classification harder than binary classification. The closest related work ([77]) only considers multiclass classification in the fixed finite class setting and consequently, doesn't compute exact dependencies on the number of classes k . We provide a more detailed comparison of our work with Wang et al. [77] and Muthukumar et al. [56] in Appendix H of the Supplemental material.

3 Related Work

The present work is situated within a larger stream of theoretical research trying to understand why overparameterized learning works and its limits. The limited page budget here forces brevity, but we recommend the recent surveys Bartlett et al. [5], Belkin [6], Dar et al. [20] for further context.

Classically, by either operating in the underparameterized regime or by performing explicit regularization, we can force the training procedure to average out the harmful effects of training noise and thereby hope to obtain good generalization. The present cycle of seeking a deeper understanding began after it was observed that modern deep networks were overparameterized, capable of memorizing noise, and yet still generalized well, even when they were trained without explicit regularization [59, 84]. Experiments in Geiger et al. [28], Belkin et al. [8] observed a double-descent behavior of the generalization error where in addition to the traditional U-shaped curve in the underparameterized regime, the error decreases in the overparameterized regime as we increase the number of model parameters. This double descent phenomenon is not unique to deep learning models and was replicated for kernel learning [7]. Further, the good generalization performance in the overparameterized regime cannot be explained by traditional worst-case generalization bounds based on Rademacher complexity or VC-dimension since the models have the capacity to fit purely random labels. Overparameterized models must therefore have some fortuitous combination of the model architecture with the training algorithm that leads us to a particular solution that generalizes well.

To understand the phenomenon better, several works study the simpler setting of overparameterized linear regression. The minimum- ℓ_2 norm⁴ interpolator is of particular interest since gradient descent on the squared loss has an implicit⁵ bias towards this solution in the overparameterized regime [24] and has been studied extensively. (An incomplete list is Hastie et al. [31], Mei and Montanari [52], Bartlett et al. [4], Belkin et al. [9], Muthukumar et al. [55], Bibas et al. [10], Kobak et al. [41], Wu and Xu [81], Richards et al. [63].) To generalize well, the underlying feature family must satisfy a balance between having a few important directions that sufficiently favor the true pattern, and a large number of unimportant directions that can absorb the noise in a harmless manner.

3.1 High dimensional binary classification

Both concurrently with and subsequent to the wave of analyses on overparameterized regression, researchers turned their attention to binary classification. A line of work poses the overparameterized binary classification problem as an optimization problem and analyzes it directly to obtain precise asymptotic behaviours of the generalization error [22, 66, 37, 69, 54, 38, 70]. The key technical tool employed in these works is the Convex Gaussian Min-max Theorem and the resultant error formulas involve solutions to a system of non-linear equations that typically do not admit closed-form expressions. The generalization error of the max-margin SVM has also been analyzed directly by studying the iterates of gradient descent in [13] and leveraging the implicit regularization perspective of optimization algorithms.

However, although the above works did significantly enhance our understanding of binary classification in the overparameterized regime, a fundamental question was not answered: "Is classification easier than regression?" While the classification task is easier than the regression task at test time (regression requires us to correctly predict a real value while binary classification requires us to only predict its sign correctly), the training data for classification is less informative than that for regression

⁴The minimum- ℓ_1 norm interpolator has also been studied in Muthukumar et al. [55], Mitra [53], Li and Wei [50], Wang et al. [75] and while sparsity-seeking behavior helps preserve the true signal (if the true pattern indeed depends only on a few features), it poses a challenge for the harmless absorption of noise since the desired averaging behaviour is not achieved fully [55].

⁵In fact, there is an important complementary literature that brings out the implicit regularization performed by training methods, especially variants of gradient descent and stochastic gradient descent, and how the underlying architecture of the model shapes this implicit regularization [30, 68, 36, 80, 57, 2, 82].

since the labels are also binary. As described earlier, this question was answered in Muthukumar et al. [56], by exhibiting an asymptotic regime where binary classification error goes to zero, but the regression error does not. This was shown using Gaussian features with a bi-level covariance model. It turns out that the level of anisotropy (favoring of true features) required to perform regression correctly is significantly higher than that required for binary classification.

The key to the result in Muthukumar et al. [56] was the signal-processing inspired survival/contamination framework introduced in Muthukumar et al. [55] as a reconceptualization of the “effective ranks” perspective of Bartlett et al. [4]. For binary classification to succeed, what matters is that the survival exceed the contamination so that the sign of the prediction remains correct. Meanwhile, regression is harder since for regression to succeed, the survival must also tend to 1.

3.2 Multiclass classification and the role of training loss function

There is a large classical body of work on multiclass classification algorithms [79, 12, 23, 18, 46], with further works giving computationally efficient algorithms for extreme multiclass problems with a huge number of classes [15, 83, 62]. Numerous theoretical works investigate the consistency of classifiers [85, 60, 61, 71, 14]. Finite-sample analysis of the generalization error in multiclass classification problems in the underparameterized regime has been studied in Koltchinskii and Panchenko [42], Guermeur [29], Allwein et al. [1], Li et al. [49], Cortes et al. [16], Lei et al. [47], Maurer [51], Lei et al. [48], Kuznetsov et al. [44, 45] and includes both data dependent bounds using Rademacher complexity, Gaussian complexity and covering numbers as well as data-independent bounds using the VC dimension. Recent work [72] leverages the Convex Gaussian Min-max Theorem to precisely characterize the asymptotic behaviour of the least-squares classifier in underparameterized multiclass classification.

So, how different is multiclass classification from binary classification? The test time task is more difficult and for the same total number of training points, we have fewer positive training examples from each class. Several empirical studies comparing the performances of multiclass classification via learning multiple binary classifiers have been undertaken [64, 25, 1]. The effects of the loss function while using deep nets to perform classification has also been investigated [32, 26, 43, 11, 21, 40, 35, 39]. Empirical evidence of least-squares minimization yielding competitive test classification performance to cross-entropy minimization has been presented in Rifkin and Klautau [64], Hui and Belkin [35], Bosman et al. [11].

More recently, Wang et al. [77] makes progress towards bridging the gap between empirical observations and theoretical understanding by proving that in certain overparameterized regimes the solution to a multiclass SVM problem is identical to the one obtained by minimum-norm interpolation of one-hot encoded labels (equivalently, that gradient descent on squared loss leads to the same solution as gradient descent on cross-entropy loss as a result of implicit bias of these algorithms [24, 36, 68]). In addition, Wang et al. [77] extends the analysis presented in Muthukumar et al. [56] for the binary classification problem to the multiclass problem with finitely many classes via an interesting reduction to analyzing a finite set of pairwise competitions, all of which must be won for multiclass classification to succeed. (We give further comments on the relationship of the present paper with Wang et al. [77] in Appendix H of the Supplemental material.)

4 Problem setup

We consider the multiclass classification problem with k classes. The training data consists of n pairs $\{\mathbf{x}_i, \ell_i\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^d$ are i.i.d Gaussian vectors drawn from distribution,

$$\mathbf{x}_i \sim \mathcal{N}(0, I_d). \quad (1)$$

We make the following assumption on how the labels $\ell_i \in [k]$ are generated.

Assumption 4.1. Orthogonal classes noiseless model⁶ *The class labels ℓ_i are generated based on which of the first k dimensions of a point \mathbf{x}_i has the largest value,*

$$\ell_i = \operatorname{argmax}_{m \in [k]} \mathbf{x}_i[m]. \quad (2)$$

⁶A more generic model is $\ell_i = \operatorname{argmax}_{m \in [k]} \mu_m^\top \mathbf{x}_i$ where the μ_m are unit norm orthogonal vectors. If we further assume the bi-level model(Definition 4.2) and that the vectors μ_m have no support outside of the favored features then it suffices to consider the simplified setting where μ_m are 1-sparse unit vectors like we do here, due to the indifference of minimum norm interpolation to orthogonal transformations.

We use the notation $x_i[m]$ to refer to the m^{th} element of vector \mathbf{x}_i . For clarity of exposition, we make explicit a feature weighting that transforms the training points as follows:

$$x_i^w[j] = \sqrt{\lambda_j} x_i[j] \quad \forall j \in [d]. \quad (3)$$

Here $\boldsymbol{\lambda} \in \mathbb{R}^d$ contains the squared feature weights. The feature weighting serves the role of favoring the true pattern, something that is essential for good generalization.⁷

The weighted feature matrix $\mathbf{X}^w \in \mathbb{R}^{n \times d}$ is given by,

$$\mathbf{X}^w = [\mathbf{x}_1^w \quad \dots \quad \mathbf{x}_j^w \quad \dots \quad \mathbf{x}_n^w]^> = [\sqrt{\lambda_1} \mathbf{z}_1 \quad \dots \quad \sqrt{\lambda_j} \mathbf{z}_j \quad \dots \quad \sqrt{\lambda_d} \mathbf{z}_d], \quad (4)$$

where $\mathbf{z}_j \in \mathbb{R}^d$ contains the j^{th} features from the n training points. Note that $\mathbf{z}_j \sim \mathcal{N}(0, I_n)$ are i.i.d Gaussians. We use a one-hot encoding for representing the labels as the matrix $\mathbf{Y}^{oh} \in \mathbb{R}^{n \times k}$,

$$\mathbf{Y}^{oh} = [\mathbf{y}_1^{oh} \quad \dots \quad \mathbf{y}_m^{oh} \quad \dots \quad \mathbf{y}_k^{oh}], \quad (5)$$

where,

$$y_m^{oh}[i] = \begin{cases} 1, & \text{if } \ell_i = m \\ 0, & \text{otherwise} \end{cases}. \quad (6)$$

A zero-mean variant of the encoding where we subtract the mean $\frac{1}{k}$ from each entry is denoted:

$$\mathbf{y}_m = \mathbf{y}_m^{oh} - \frac{1}{k} \mathbf{1}. \quad (7)$$

Our classifier consists of k coefficient vectors $\hat{\mathbf{f}}_m$ for $m \in [k]$ that are learned by minimum-norm interpolation of the zero-mean one-hot variants using the weighted features.⁸

$$\hat{\mathbf{f}}_m = \arg \min_{\mathbf{f}} \|\mathbf{X}^w \mathbf{f} - \mathbf{y}_m\|_2 \quad (8)$$

$$\text{s.t. } \mathbf{X}^w \mathbf{f} = \mathbf{y}_m - \frac{1}{k} \mathbf{1}. \quad (9)$$

We can express these coefficients in closed form as,

$$\hat{\mathbf{f}}_m = (\mathbf{X}^w)^> (\mathbf{X}^w (\mathbf{X}^w)^>)^{-1} \mathbf{y}_m. \quad (10)$$

On a test point $\mathbf{x}_{test} \sim \mathcal{N}(0, I_d)$ we predict a label as follows: First, we transform the test point into the weighted feature space to obtain \mathbf{x}_{test}^w where $x_{test}^w[j] = \sqrt{\lambda_j} x_{test}[j]$ for $j \in [d]$. Then we compute k scalar ‘‘scores’’ and assign the class based on the largest score as follows:

$$\hat{\ell} = \arg \max_{1 \leq m \leq k} \hat{\mathbf{f}}_m^> \mathbf{x}_{test}^w. \quad (11)$$

The true label of the test point is $\ell_{test} = \arg \max_{1 \leq m \leq k} x_{test}[m]$. A misclassification event E_{err} occurs iff

$$\arg \max_{1 \leq m \leq k} x_{test}[m] \neq \arg \max_{1 \leq m \leq k} \hat{\mathbf{f}}_m^> \mathbf{x}_{test}^w. \quad (12)$$

In our work we determine sufficient conditions under which the probability of misclassification (computed over the randomness in both the training data and test point) goes to zero in an asymptotic regime where the number of training points, number of features, number of classes and feature weights scale according to the bi-level ensemble model.

⁷Our weighted feature model is equivalent to the one used in other works (e.g. [56]) that assume that the covariates come from an anisotropic Gaussian with a covariance matrix that favors the truly important directions.

⁸The classifier learned via this method is equivalent to those obtained by other natural training methods under sufficient overparameterization [77].

Definition 4.2. (Bi-level ensemble): The bi-level ensemble is parameterized by p, q, r and t where $p > 1, 0 < r < 1, 0 < q < (p - r)$ and $0 < t < r$. Here, parameter p controls the extent of overparameterization, r determines the number of favored features, q controls the weights on favored features and t controls the number of classes. The number of features (d), number of favored features (s), number of classes (k) and feature weights ($\sqrt{\lambda_j}$) all scale with the number of training points (n) as follows:

$$d = bn^p c, s = bn^r c, a = n^{-q}, k = c_k bn^t c, \quad (13)$$

where c_k is a positive integer. The feature weights are given by,

$$\sqrt{\lambda_j} = \begin{cases} \sqrt{\frac{ad}{s}}, & 1 \leq j \leq s \\ \sqrt{\frac{(1-a)d}{d-s}}, & \text{otherwise} \end{cases}. \quad (14)$$

We provide a visualization of the bi-level model in Figure 1.

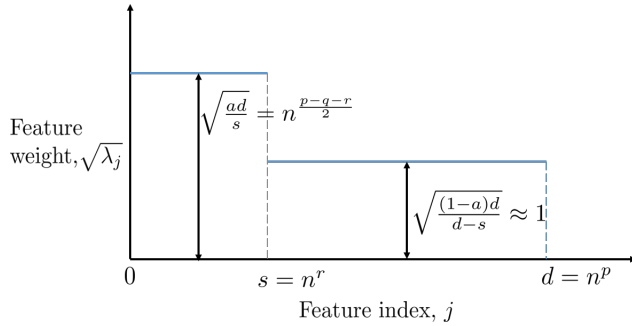


Figure 1: Bi-level feature weighting model. The first s features have a higher weight and are favored during minimum-norm interpolation. These can be thought of as the square-roots of the eigenvalues of the feature covariance matrix in a Gaussian model for the covariates as in Bartlett et al. [4].

5 Main result

Theorem 5.1. (Asymptotic classification region in the bi-level model): Under the bi-level ensemble model 4.2, when the true data generating process is 1-sparse (Assumption 4.1), the probability of misclassification $P(E_{err}) \rightarrow 0$ as $n \rightarrow \infty$ if the following conditions hold:

$$t < \min(r, 1 - r, p + 1 - 2(q + r), p - 2, 2q + r - 2) \quad (15)$$

$$q + r > 1. \quad (16)$$

Note that from Muthukumar et al. [56], the condition $q + r > 1$ corresponds to the regime where the corresponding regression does not generalize well and thus our result shows that multiclass classification can generalize in regimes where the corresponding regression problem does not. In this challenging regime, the empirical eigenstructure does not reveal the true nature of underlying features as illustrated in Appendix J.

Figure 2 visualizes the regimes by considering slices of the four dimensional scaling parameter space of p, q, r and t . (1a) and (2a) fix the value of q to 0.75 and 0.95 respectively and contrast the multiclass problem with a fixed finite number of classes ($t = 0$) to the binary classification and regression problems. From these plots we observe that if we fix p, q, t and increase r , i.e. increasing how many features are favored (and thereby favoring each of them less), we transition from the regime where both regression and binary classification work, into the regime where binary classification works but regression does not, then the regime where this paper can prove multiclass classification works and finally to the regime where neither regression nor binary classification works.

In Figure 2, subplots (1b),(1c),(2b) and (2c) each visualize a slice along the r and t (class scaling) dimensions with fixed p and q . The x axis itself in these plots corresponds to a fixed finite classes setting. From (1b) we observe that the right-hand boundary of the region where multiclass classification generalizes well contains two slopes. These slopes arise from the two conditions $t < 1 - r$ and

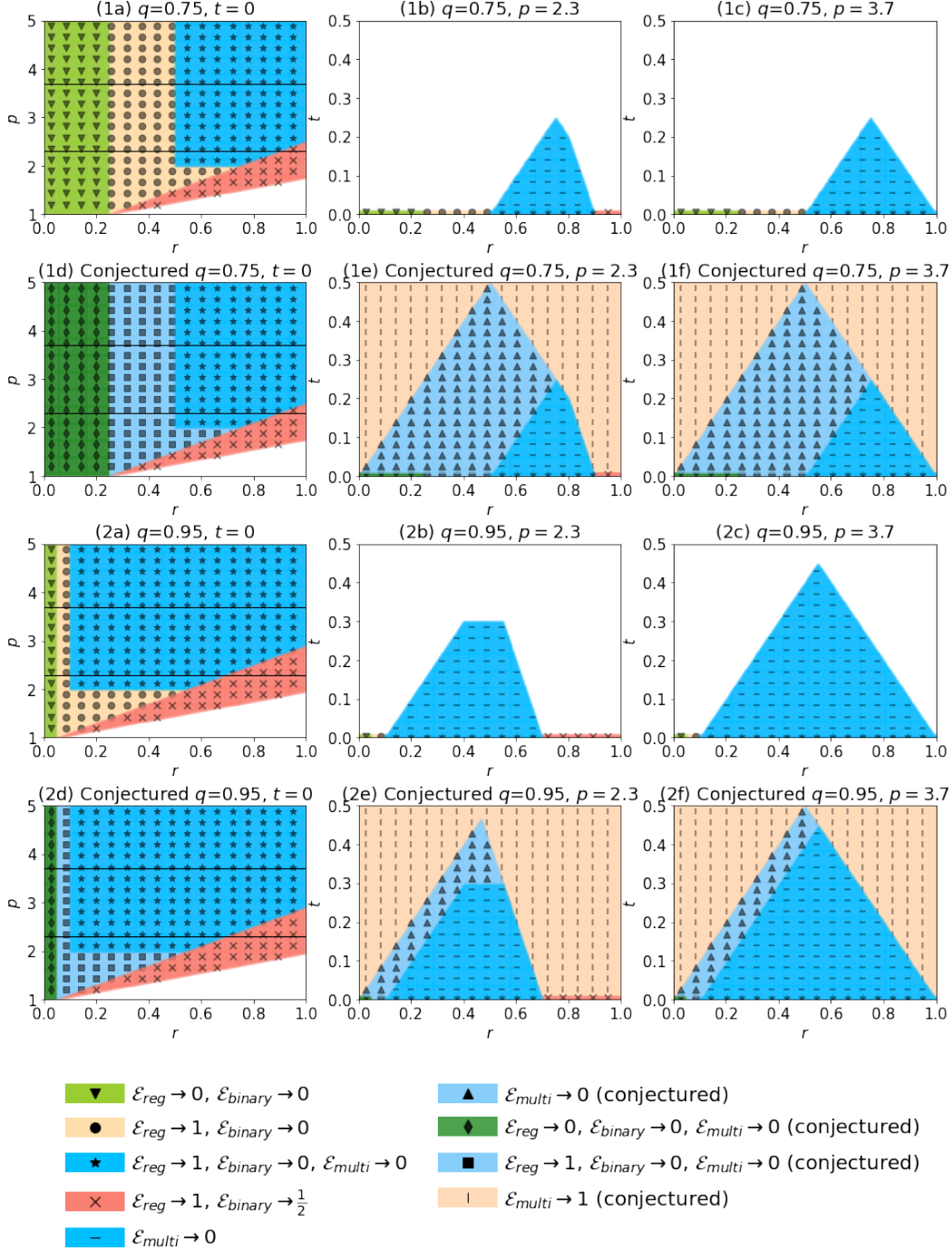


Figure 2: Visualization of the bi-level regimes in four dimensions p, q, r, t . (1a) and (2a) contrast multiclass classification with finite classes to binary classification and regression. The horizontal lines $p = 2.3$ and $p = 3.7$ correspond to the slices visualized in (1b), (1c), (2b) and (2c). The conjectured regimes are visualized in (1d), (1e), (1f), (2d), (2e) and (2f).

$t < p + 1 - 2(q + r)$ in Theorem 5.1 and are a result of either contamination from favored (but not true) features dominating or contamination from the unfavored features dominating. In (1c) we are

in the regime where binary classification works for all values of $r < 1$. However, as we increase t , eventually multiclass classification stops working.⁹

When we go from the binary problem to a multiclass problem with k classes, the survival drops by a factor of k as a consequence of having only $\frac{1}{k}$ fraction of positive training examples per class. This is because the one-hot labels we interpolate while training have fewer large values close to 1 that are able to positively correlate with the true feature vector. Having fewer positive exemplars also reduces the total energy in the training vector by a factor of k , and because of the square-root relationship of the standard deviation to the energy, the contamination only shrinks by a factor of \sqrt{k} . The overall survival/contamination ratio decreases by a factor of \sqrt{k} making the multiclass classification task more difficult.¹⁰ An interesting observation here is the amount of favoring required for good generalization is linked to the number of positive training examples per class. Indeed, if we consider a setting where the binary classification problem generalizes well, and we switch to the k class multiclass problem, then by increasing the number of training samples k fold (and thus matching the number of positive training examples per class in the multiclass case to the binary case) and keeping the number of features and feature weights constant we can generalize well for multiclass classification. (Appendix G of the Supplemental material elaborates on this phenomenon, as well as why it is somewhat surprising.)

Next, we present a brief overview of our proof that utilizes the survival/contamination analysis framework from Muthukumar et al. [56] along with a typicality-inspired argument where the feature margin (difference between largest and second largest feature) on the test point plays a key role. The complete proof is provided in Appendices B, C, D, and E of the Supplemental material.

5.1 Proof sketch

Assume without loss of generality that for the test point $\mathbf{x}_{test} \sim \mathcal{N}(0, I_d)$, the true class is α for some $\alpha \in [k]$. Let \mathbf{x}_{test}^w be the weighted version of this test point. A necessary and sufficient condition for classification error is that for some $\beta \neq \alpha, \beta \in [k]$,

$$\begin{aligned} \hat{f}_\alpha[\alpha]x_{test}^w[\alpha] + \hat{f}_\alpha[\beta]x_{test}^w[\beta] + \sum_{j \notin \{\alpha, \beta\}} \hat{f}_\alpha[j]x_{test}^w[j] &< \hat{f}_\beta[\alpha]x_{test}^w[\alpha] \\ &+ \hat{f}_\beta[\beta]x_{test}^w[\beta] + \sum_{j \notin \{\alpha, \beta\}} \hat{f}_\beta[j]x_{test}^w[j]. \end{aligned} \quad (17)$$

By converting into the unweighted feature space we obtain

$$\lambda_\alpha \hat{h}_{\alpha, \beta}[\alpha]x_{test}[\alpha] - \lambda_\beta \hat{h}_{\beta, \alpha}[\beta]x_{test}[\beta] < \sum_{j \notin \{\alpha, \beta\}} \lambda_j \hat{h}_{\beta, \alpha}[j]x_{test}[j], \quad (18)$$

where

$$\hat{h}_{\alpha, \beta}[j] = \lambda_j^{1/2}(\hat{f}_\alpha[j] - \hat{f}_\beta[j]). \quad (19)$$

Performing some algebraic manipulations and because $\lambda_\alpha = \lambda_\beta = \lambda$ since both α and β are favored features, we can rewrite this as

$$\begin{aligned} \frac{\lambda \hat{h}_{\alpha, \beta}[\alpha]}{\text{CN}_{\alpha, \beta}} \left((x_{test}[\alpha] - x_{test}[\beta]) + x_{test}[\beta] \frac{\hat{h}_{\alpha, \beta}[\alpha] - \hat{h}_{\beta, \alpha}[\beta]}{\hat{h}_{\alpha, \beta}[\alpha]} \right) \\ < \frac{1}{\text{CN}_{\alpha, \beta}} \sum_{j \notin \{\alpha, \beta\}} \lambda_j \hat{h}_{\beta, \alpha}[j]x_{test}[j], \end{aligned} \quad (20)$$

⁹To be precise, what the region actually illustrates is that our proof approach stops being able to show that multiclass classification works. In the Conclusion section, we conjecture where we believe that multiclass classification actually stops working. The conjectured regions are illustrated in (1e),(1f),(2e) and (2f).

¹⁰This is also responsible for contamination due to favored features being able to cause errors. For binary classification, because the true feature survival is constant (depending only on the level of label noise), the survival can always asymptotically overcome any contamination from other favored features [56].

where

$$\text{CN}_{\alpha,\beta} = \sqrt{\left(\sum_{j \notin \mathcal{F}_{\alpha,\beta g}} \lambda_j^2 (\hat{h}_{\beta,\alpha}[j])^2 \right)}. \quad (21)$$

We divide by $\text{CN}_{\alpha,\beta}$ to normalize the RHS above to have a standard normal distribution. Next, by removing the dependency on β , we obtain a sufficient condition for correct classification:

$$\underbrace{\frac{\min_{\beta} \lambda \hat{h}_{\alpha,\beta}[\alpha]}{\max_{\beta} \text{CN}_{\alpha,\beta}}}_{\text{SU/CN ratio}} \left(\underbrace{\min_{\beta} (x_{test}[\alpha] \quad x_{test}[\beta])}_{\text{closest feature margin}} \quad \underbrace{\max_{\beta} j x_{test}[\beta] j}_{\text{largest competing feature}} \quad \underbrace{\max_{\beta} \left| \frac{\hat{h}_{\alpha,\beta}[\alpha] \quad \hat{h}_{\beta,\alpha}[\beta]}{\hat{h}_{\alpha,\beta}[\alpha]} \right|}_{\text{survival variation}} \right) > \underbrace{\max_{\beta} \frac{1}{\text{CN}_{\alpha,\beta}} \left(\sum_{j \notin \mathcal{F}_{\alpha,\beta g}} \lambda_j \hat{h}_{\beta,\alpha}[j] x_{test}[j] \right)}_{\text{normalized contamination}}. \quad (22)$$

Here the min and max are over all competing features: $1 \leq \beta \leq k, \beta \neq \alpha$ and the sum is over all d feature indices except α and β , but we simplify the notation for convenience. We show via intermediate lemmas introduced in Appendix B of the Supplemental material that under the conditions specified in Theorem 5.1, with sufficiently high probability¹¹, the relevant survival to contamination SU/CN ratio grows at a polynomial rate n^v for some $v > 0$, the closest feature margin shrinks at a less-than-polynomial rate $1/\sqrt{\ln nk}$, the survival variation decays at a polynomial rate n^{-u} for some $u > 0$. Further, the magnitudes of the largest competing feature and the normalized contamination are no more than $\sqrt{\ln(nk)}$.

This implies that the left-hand side of Equation (22) grows at a polynomial rate n^v (ignoring logarithmic terms) and dominates the right-hand side which grows at the much slower rate $1/\sqrt{\ln nk}$. A survival/contamination ratio also plays a key role in the analysis of the binary classification problem in Muthukumar et al. [56] but in the multiclass setting, we additionally have the survival variation term and feature margin playing important roles since we are comparing different scores while predicting the class label. For correct classification, the survival/contamination ratio must be sufficiently large, the survival variation must be small enough and the feature margin must be sufficiently large.

6 Conclusion

In this work we compute sufficient conditions for good generalization of multiclass classification in a bi-level overparameterized linear model with Gaussian features. We observed that multiclass classification can generalize even when the regression problem does not generalize (for $q + r > 1$). Further, the multiclass problem is “harder” than the binary problem because we have fewer positive training examples per class. The nature of the training data complicates our analysis in the multiclass setting since the true class labels are generated by comparing k features and thus we no longer have independence of the encoded class label y with any of these features. This becomes relevant when we compute bounds on the survival and contamination quantities since the Hanson-Wright inequality [65] is no longer applicable directly on the quantities of interest as was the case for the binary classification problem in prior work [56]. As a consequence of working around this non-independence we believe that our sufficient conditions for good generalization in the regime $q + r > 1$ are loose.

Even though in our work we focus on the regime where regression does not work, $q + r > 1$, we can extend the analysis to the regime where $q + r < 1$ by grinding through the expressions for survival and contamination in this regime. Even in this regime, for multiclass training data, survival is of the order $\frac{1}{k}$ while contamination scales similarly to the regime $q + r > 1$. Thus, while it is true that for

¹¹This is where we leverage the idea of typicality-style proofs in information theory [17] to avoid unnecessarily loose union bounds that end up being dominated by the atypical behavior of quantities. In our case, by pulling the feature margin out explicitly, we can just deal with its typical behavior. Similarly, the typical behavior of the largest competing feature and the true feature is all that matters.

binary classification or a fixed number of classes, the regime where regression works is a regime where classification also works, this need not be true if there are too many classes.

We conjecture that the following is a set of necessary and sufficient conditions for asymptotically good generalization (We elaborate on this in Appendix F in the Supplemental material):

Conjecture 6.1. (Conjectured bi-level regions): *Under the bi-level ensemble model 4.2, when the true data generating process is 1-sparse (Assumption 4.1), as $n \rightarrow \infty$, the probability of misclassification event $P(E_{err})$ behaves as follows:*

$$P(E_{err}) \rightarrow \begin{cases} 0, & \text{if } t < \min(r, 1 - r, p + 1 - 2 \max(1, q + r)) \\ 1, & \text{if } t > \min(r, 1 - r, p + 1 - 2 \max(1, q + r)) \end{cases}. \quad (23)$$

The conjectured regions are visualized in (1d),(1e),(1f),(2d),(2e) and (2f) in Figure 2. Subfigures (1d) and (2d) illustrate that we believe multiclass classification with finitely many classes works if binary classification works. Further, comparing (1e) to (2e) when we increase q , the conjectured parameter region where multiclass classification works shrinks since we decrease the amount of favoring of true features. Interestingly, the nature of the looseness in our approach is such that our proof technique is able to recover a larger fraction of the conjectured region for larger q which intuitively is a result of less favoring leading to stronger concentration of certain random quantities. Tightening the potential looseness in our analysis and proving the converse result by computing sufficient conditions for poor generalization of multiclass classification are interesting avenues of future work.

Further, although the present analysis focuses on solutions that exactly interpolate the training data, we can extend our results to account for additional ridge regularization by viewing ridge regularization as minimum-norm interpolation using augmented contamination-free features as in the Appendix of Muthukumar et al. [55] and computing bounds leveraging tools from Tsigler and Bartlett [74]. Our assumption of the strict bi-level weighting model is largely to simplify the calculations and by substituting terms appropriately in our lemmas from Appendix B in the Supplemental material, it should be possible to compute results for other weighting models. Finally, exploring the new phenomena that can be encountered as we go beyond the 1-sparse noiseless model is an exciting direction for future work.

Acknowledgments and Disclosure of Funding

We are grateful to our earlier collaborators Vidya Muthukumar, Misha Belkin, Daniel Hsu, and Adhyayan Narang. In addition, we want to thank the students and course staff for the Fall 2020 iteration of Berkeley’s CS189/289A machine learning courses, where we had adapted ideas from Muthukumar et al. [55, 56] in teaching the foundations of modern machine learning — the need for the present paper became more clear during this process.

We gratefully acknowledge the support from ML4Wireless center member companies and NSF grants AST-2132700 and AST-2037852 for making this research possible.

References

- [1] Erin Allwein, Robert E. Schapire, and Yoram Singer. Reducing multiclass to binary: A unifying approach for margin classifiers. *Journal of Machine Learning Research*, 1:113–141, 2000.
- [2] Navid Azizan, Sahin Lale, and Babak Hassibi. A study of generalization of stochastic mirror descent algorithms on overparameterized nonlinear models. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3132–3136. IEEE, 2020.
- [3] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [4] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- [5] Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint. *Acta numerica*, 30:87–201, 2021.

- [6] Mikhail Belkin. Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation. Acta Numerica, 30:203–248, 2021.
- [7] Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand kernel learning. ICML, 2018.
- [8] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. Proceedings of the National Academy of Sciences, 116(32):15849–15854, 2019.
- [9] Mikhail Belkin, Daniel Hsu, and Ji Xu. Two models of double descent for weak features. SIAM Journal on Mathematics of Data Science, 2(4):1167–1180, 2020.
- [10] Koby Bibas, Yaniv Fogel, and Meir Feder. A new look at an old problem: A universal learning approach to linear regression. CoRR, abs/1905.04708, 2019. URL <http://arxiv.org/abs/1905.04708>.
- [11] Anna Bosman, Andries Engelbrecht, and Mardé Helbig. Visualising basins of attraction for the cross-entropy and the squared error neural network loss functions. Neurocomputing, 400, 03 2020. doi: 10.1016/j.neucom.2020.02.113.
- [12] Erin J. Bredensteiner and Kristin P. Bennett. Multicategory classification by support vector machines. Computational Optimization and Applications, 12, 1999. doi: 10.1023/A:1008663629662.
- [13] Niladri S Chatterji and Philip M Long. Finite-sample analysis of interpolating linear classifiers in the overparameterized regime. Journal of Machine Learning Research, 22(129):1–30, 2021.
- [14] Di-Rong Chen and Tao Sun. Consistency of multiclass empirical risk minimization methods based on convex loss. Journal of Machine Learning Research, 7:2435–2447, dec 2006. ISSN 1532-4435.
- [15] Anna Choromanska, Alekh Agarwal, and John Langford. Extreme multi class classification. In NIPS Workshop: eXtreme Classification, submitted, volume 1, pages 2–1, 2013.
- [16] Corinna Cortes, Vitaly Kuznetsov, Mehryar Mohri, and Scott Yang. Structured prediction theory based on factor graph complexity. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/535ab76633d94208236a2e829ea6d888-Paper.pdf>.
- [17] Thomas M. Cover and Joy A. Thomas. Elements of Information Theory 2nd Edition (Wiley Series in Telecommunications and Signal Processing). Wiley-Interscience, July 2006. ISBN 0471241954.
- [18] Koby Crammer, Yoram Singer, Nello Cristianini, John Shawe-taylor, and Bob Williamson. On the algorithmic implementation of multiclass kernel-based vector machines. Journal of Machine Learning Research, 2:265–292, 2001.
- [19] Alexander D’Amour, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, Christina Chen, Jonathan Deaton, Jacob Eisenstein, Matthew D Hoffman, et al. Under-specification presents challenges for credibility in modern machine learning. arXiv preprint arXiv:2011.03395, 2020.
- [20] Yehuda Dar, Vidya Muthukumar, and Richard G Baraniuk. A farewell to the bias-variance tradeoff? an overview of the theory of overparameterized machine learning. arXiv preprint arXiv:2109.02355, 2021.
- [21] Ahmet Demirkaya, Jiasi Chen, and Samet Oymak. Exploring the role of loss functions in multiclass classification. In 2020 54th Annual Conference on Information Sciences and Systems (CISS), pages 1–5, 2020. doi: 10.1109/CISS48834.2020.1570627167.

- [22] Zeyu Deng, Abla Kammoun, and Christos Thrampoulidis. A model of double descent for high-dimensional binary linear classification. Information and Inference: A Journal of the IMA, 04 2021. ISSN 2049-8772. doi: 10.1093/imaiai/iaab002. URL <https://doi.org/10.1093/imaiai/iaab002>. iaab002.
- [23] Thomas G Dietterich and Ghulum Bakiri. Solving multiclass learning problems via error-correcting output codes. Journal of Artificial Intelligence Research, 2(1):263–286, 1994. ISSN 1076-9757.
- [24] Heinz Werner Engl, Martin Hanke, and Andreas Neubauer. Regularization of inverse problems, volume 375. Springer Science & Business Media, 1996.
- [25] Johannes Fürnkranz. Round robin classification. Journal of Machine Learning Research, 2: 721–747, 2002.
- [26] Krzysztof Gajowniczek, Leszek Chmielewski, Arkadiusz Orłowski, and Tomasz Ząbkowski. Generalized entropy cost function in neural networks. In International Conference on Artificial Neural Networks, pages 128–136, 10 2017. ISBN 978-3-319-68611-0. doi: 10.1007/978-3-319-68612-7_15.
- [27] Robert G Gallager. Information theory and reliable communication, volume 588. Springer, 1968.
- [28] Mario Geiger, Stefano Spigler, Stéphane d’Ascoli, Levent Sagun, Marco Baity-Jesi, Giulio Biroli, and Matthieu Wyart. Jamming transition as a paradigm to understand the loss landscape of deep neural networks. Physical Review E, 100(1):012115, 2019.
- [29] Yann Guermeur. Combining Discriminant Models with New Multi-Class SVMs. Pattern Anal. Appl., 5:168–179, 06 2002. doi: 10.1007/s100440200015.
- [30] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. In International Conference on Machine Learning, pages 1832–1841, 2018.
- [31] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. arXiv preprint arXiv:1903.08560, 2019.
- [32] Le Hou, Chen-Ping Yu, and Dimitris Samaras. Squared Earth Mover’s Distance-based Loss for Training Deep Neural Networks. arXiv e-prints, art. arXiv:1611.05916, November 2016.
- [33] Daniel Hsu, Vidya Muthukumar, and Ji Xu. On the proliferation of support vectors in high dimensions. In International Conference on Artificial Intelligence and Statistics, pages 91–99. PMLR, 2021.
- [34] Iosif Pinelis (https://mathoverflow.net/users/36721/iosif_pinelis). Concentration and anti-concentration of gap between largest and second largest value in gaussian iid sample. MathOverflow. URL <https://mathoverflow.net/q/379688>. URL:<https://mathoverflow.net/q/379688> (version: 2020-12-25).
- [35] Like Hui and Mikhail Belkin. Evaluation of Neural Architectures Trained with Square Loss vs Cross-Entropy in Classification Tasks. arXiv e-prints, art. arXiv:2006.07322, June 2020.
- [36] Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In Conference on Learning Theory, pages 1772–1798, 2019.
- [37] Abla Kammoun and Mohamed-Slim Alouini Fellow. On the precise error analysis of support vector machines. IEEE Open Journal of Signal Processing, 2:99–118, 2021.
- [38] Ganesh Ramachandra Kini and Christos Thrampoulidis. Analytic study of double descent in binary classification: The impact of loss. In 2020 IEEE International Symposium on Information Theory (ISIT), pages 2527–2532. IEEE, 2020.

- [39] Ganesh Ramachandra Kini, Orestis Paraskevas, Samet Oymak, and Christos Thrampoulidis. Label-imbalanced and group-sensitive classification under overparameterization. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, Advances in Neural Information Processing Systems, volume 34, pages 18970–18983. Curran Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/file/9dfcf16f0adbc5e2a55ef02db36bac7f-Paper.pdf>.
- [40] Doug M. Kline and Victor L. Berardi. Revisiting squared-error and cross-entropy functions for training neural network classifiers. Neural Computing & Applications, 14:310–318, 2005.
- [41] Dmitry Kobak, Jonathan Lomond, and Benoit Sanchez. The optimal ridge penalty for real-world high-dimensional data can be zero or negative due to the implicit ridge regularization. Journal of Machine Learning Research, 21:169–1, 2020.
- [42] Vladimir Koltchinskii and Dmitry Panchenko. Empirical margin distributions and bounding the generalization error of combined classifiers. The Annals of Statistics, 30(1):1–50, 2002. ISSN 00905364. URL <http://www.jstor.org/stable/2700001>.
- [43] Himanshu Kumar and P. Shanti Sastry. Robust loss functions for learning multi-class classifiers. 2018 IEEE International Conference on Systems, Man, and Cybernetics (SMC), pages 687–692, 2018.
- [44] Vitaly Kuznetsov, Mehryar Mohri, and Umar Syed. Multi-class deep boosting. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, editors, Advances in Neural Information Processing Systems, volume 27. Curran Associates, Inc., 2014. URL <https://proceedings.neurips.cc/paper/2014/file/7bb060764a818184ebb1cc0d43d382aa-Paper.pdf>.
- [45] Vitaly Kuznetsov, Mehryar Mohri, and Umar Syed. Rademacher complexity margin bounds for learning with a large number of classes. In ICML Workshop on Extreme Classification: Learning with a Very Large Number of Labels, 2015.
- [46] Yoonkyung Lee, Yi Lin, and Grace Wahba. Multicategory support vector machines: Theory and application to the classification of microarray data and satellite radiance data. Journal of the American Statistical Association, 99(465):67–81, 2004. doi: 10.1198/016214504000000098. URL <https://doi.org/10.1198/016214504000000098>.
- [47] Yunwen Lei, Urun Dogan, Alexander Binder, and Marius Kloft. Multi-class SVMs: From Tighter Data-Dependent Generalization Bounds to Novel Algorithms. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 28. Curran Associates, Inc., 2015. URL <https://proceedings.neurips.cc/paper/2015/file/3a029f04d76d32e79367c4b3255dda4d-Paper.pdf>.
- [48] Yunwen Lei, Ürün Dogan, Ding-Xuan Zhou, and Marius Kloft. Data-dependent generalization bounds for multi-class classification. IEEE Transactions on Information Theory, 65(5):2995–3021, 2019. doi: 10.1109/TIT.2019.2893916.
- [49] Jian Li, Yong Liu, Rong Yin, Hua Zhang, Lizhong Ding, and Weiping Wang. Multi-class learning: From theory to algorithm. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/1141938ba2c2b13f5505d7c424ebae5f-Paper.pdf>.
- [50] Yue Li and Yuting Wei. Minimum ℓ_1 -norm interpolators: Precise asymptotics and multiple descent. arXiv preprint arXiv:2110.09502, 2021.
- [51] Andreas Maurer. A vector-contraction inequality for Rademacher complexities. In Hans Ulrich Simon Ronald Ortner and Sandra Zilles, editors, Algorithmic Learning Theory, pages 3–17. Springer International Publishing, 2016.
- [52] Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve. arXiv preprint arXiv:1908.05355, 2019.

- [53] Partha P Mitra. Understanding overfitting peaks in generalization error: Analytical risk curves for ℓ_2 and ℓ_1 penalized interpolation. arXiv preprint arXiv:1906.03667, 2019.
- [54] Andrea Montanari, Feng Ruan, Youngtak Sohn, and Jun Yan. The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime. arXiv preprint arXiv:1911.01544, 2019.
- [55] Vidya Muthukumar, Kailas Vodrahalli, Vignesh Subramanian, and Anant Sahai. Harmless interpolation of noisy data in regression. IEEE Journal on Selected Areas in Information Theory, 1(1):67–83, 2020.
- [56] Vidya Muthukumar, Adhyayan Narang, Vignesh Subramanian, Mikhail Belkin, Daniel J. Hsu, and Anant Sahai. Classification vs regression in overparameterized regimes: Does the loss function matter? Journal of Machine Learning Research, 22:222:1–222:69, 2021.
- [57] Mor Shpigel Nacson, Jason Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese, Nathan Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. In The 22nd International Conference on Artificial Intelligence and Statistics, pages 3420–3428, 2019.
- [58] Preetum Nakkiran. More Data Can Hurt for Linear Regression: Sample-wise Double Descent. arXiv e-prints, art. arXiv:1912.07242, December 2019.
- [59] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. arXiv preprint arXiv:1412.6614, 2014.
- [60] Bernardo Ávila Pires and Csaba Szepesvári. Multiclass Classification Calibration Functions. arXiv e-prints, art. arXiv:1609.06385, September 2016.
- [61] Bernardo Ávila Pires, Mohammad Ghavamzadeh, and Csaba Szepesvári. Cost-sensitive multiclass classification risk bounds. In Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28, pages 1391–1399, 2013.
- [62] Ankit Singh Rawat, Jiecao Chen, Felix Xinnan X Yu, Ananda Theertha Suresh, and Sanjiv Kumar. Sampled softmax with random fourier features. Advances in Neural Information Processing Systems, 32, 2019.
- [63] Dominic Richards, Jaouad Mourtada, and Lorenzo Rosasco. Asymptotics of ridge (less) regression under general source condition. In International Conference on Artificial Intelligence and Statistics, pages 3889–3897. PMLR, 2021.
- [64] Ryan Rifkin and Aldebaro Klautau. In defense of one-vs-all classification. Journal of Machine Learning Research, 5:101–141, 12 2004.
- [65] Mark Rudelson and Roman Vershynin. Hanson-Wright inequality and sub-Gaussian concentration. Electronic Communications in Probability, 18:1–9, 2013.
- [66] Fariborz Salehi, Ehsan Abbasi, and Babak Hassibi. The impact of regularization on high-dimensional logistic regression. Advances in Neural Information Processing Systems, 32, 2019.
- [67] Ohad Shamir. The implicit bias of benign overfitting. arXiv preprint arXiv:2201.11489, 2022.
- [68] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. Journal of Machine Learning Research, 19(1):2822–2878, 2018.
- [69] Hossein Taheri, Ramtin Pedarsani, and Christos Thrampoulidis. Sharp asymptotics and optimal performance for inference in binary models. In International Conference on Artificial Intelligence and Statistics, pages 3739–3749. PMLR, 2020.
- [70] Hossein Taheri, Ramtin Pedarsani, and Christos Thrampoulidis. Fundamental limits of ridge-regularized empirical risk minimization in high dimensions. In International Conference on Artificial Intelligence and Statistics, pages 2773–2781. PMLR, 2021.

- [71] Ambuj Tewari and Peter Bartlett. On the consistency of multiclass classification methods. *Journal of Machine Learning Research*, 8:143–157, 01 2005.
- [72] Christos Thrampoulidis, Samet Oymak, and Mahdi Soltanolkotabi. Theoretical insights into multiclass classification: A high-dimensional asymptotic view. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 8907–8920. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/6547884cea64550284728eb26b0947ef-Paper.pdf>.
- [73] Nilesh Tripuraneni, Ben Adlam, and Jeffrey Pennington. Overparameterization improves robustness to covariate shift in high dimensions. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 13883–13897. Curran Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/file/73fed7fd472e502d8908794430511f4d-Paper.pdf>.
- [74] Alexander Tsigler and Peter L Bartlett. Benign overfitting in ridge regression. *arXiv preprint arXiv:2009.14286*, 2020.
- [75] Guillaume Wang, Konstantin Donhauser, and Fanny Yang. Tight bounds for minimum ℓ_1 -norm interpolation of noisy data. *arXiv preprint arXiv:2111.05987*, 2021.
- [76] Ke Wang and Christos Thrampoulidis. Benign overfitting in binary classification of Gaussian mixtures. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4030–4034. IEEE, 2021.
- [77] Ke Wang, Vidya Muthukumar, and Christos Thrampoulidis. Benign Overfitting in Multiclass Classification: All Roads Lead to Interpolation. *arXiv e-prints*, art. arXiv:2106.10865, June 2021.
- [78] Weichen Wang and Jianqing Fan. Asymptotics of empirical eigenstructure for high dimensional spiked covariance. *Annals of statistics*, 45(3):1342, 2017.
- [79] Jason Weston and Chris Watkins. Multi-class support vector machines. Technical report, 1998.
- [80] Blake Woodworth, Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Kernel and deep regimes in overparametrized models. *arXiv preprint arXiv:1906.05827*, 2019.
- [81] Denny Wu and Ji Xu. On the optimal weighted ℓ_2 regularization in overparameterized linear regression. *Advances in Neural Information Processing Systems*, 33:10112–10123, 2020.
- [82] Jingfeng Wu, Difan Zou, Vladimir Braverman, and Quanquan Gu. Direction matters: On the implicit bias of stochastic gradient descent with moderate learning rate. *arXiv preprint arXiv:2011.02538*, 2020.
- [83] Ian En-Hsu Yen, Xiangru Huang, Pradeep Ravikumar, Kai Zhong, and Inderjit Dhillon. Pd-sparse: A primal and dual sparse approach to extreme multiclass and multilabel classification. In *International conference on machine learning*, pages 3069–3077. PMLR, 2016.
- [84] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- [85] Tong Zhang. Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5:1225–1251, 2004.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) In the conclusion as well as in the statements of the theorems.

- (c) Did you discuss any potential negative societal impacts of your work? [No] Because this is a purely theoretical paper elucidating the foundations of overparameterized learning. The negative societal impacts are largely a matter of opinion whether theoretical insights from idealized models are useful. Because they are useful in teaching, we believe this is a positive social impact. However, some might believe that theoretical results for key ideas can end up having a gatekeeping effect in teaching that disadvantages populations with less access to advanced mathematics, and their mere existence therefore empowers would-be gatekeepers. However, we have attempted to mitigate this risk by building this theory together with more intuitive explanations that, while still mathematical, we believe reduce the barrier to accessing the core insights. Since this debate is orthogonal to the specific claims in this paper, we didn't feel that it needed to be addressed in the paper itself.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
- (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes] But in the Supplemental material since they are too long to fit.
3. If you ran experiments...
- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A] Not really relevant since there are no empirical experiments, only plots of the regions defined by the theorems — which are themselves given by simple linear inequalities.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A] Not applicable since no plots are random in any way.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A] Not applicable since everything could be done by hand.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? [Yes] We cited the papers that inspired us and whose results we are building upon.
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

The appendix is organized as follows. In Appendix A we provide a table of notations used throughout the paper. Appendix B provides an overall proof for Theorem 5.1 by introducing some intermediate lemmas and assuming they hold. Appendix C introduces some key tools that we need and Appendix D leverages those tools to build towards a proof of these intermediate lemmas by introducing some helper results that are needed to deal with the key challenge posed by multiclass training data. Appendix E actually proves the intermediate lemmas used in Appendix B and completes the proof. Appendix F discusses the potential looseness in our analysis and describes how we obtained Conjecture 6.1. Appendix G elaborates on the effect of fewer number of positive training examples per class in the multiclass setting and investigates an alternative setting where the total number of positive training examples per class is kept constant while we increase the number of classes. In Appendix H we provide a more detailed comparison of our work with Wang et al. [77] and Muthukumar et al. [56]. Appendix I complements the theoretical and asymptotic proofs in our paper with an empirical evaluation of relevant quantities using simulated data. Appendix J illustrates the underlying challenge provided by the regime in which regression does not generalize — namely that the empirical eigenstructure does not reveal the true nature of the underlying features. Finally, in Appendix K we provide a heuristic derivation/conjecture of our main results using an asymptotic linear-algebraic perspective.

Throughout the appendix we will assume that n is large enough for asymptotic behavior to kick in. We also will introduce various universal positive constants, indexed as c_i . These constants are all independent of n , and constants with the same index are to be treated as equal throughout this Appendix.

A Notation

We summarize the notation used in the problem setup (as well as some terms defined later) as follows:

B Proof of Theorem 5.1

We restate Theorem 5.1, our main result, here for convenience:

Theorem 5.1. (*Asymptotic classification region in the bi-level model*): *Under the bi-level ensemble model 4.2, when the true data generating process is 1-sparse (Assumption 4.1), the probability of misclassification $P(E_{err}) \rightarrow 0$ as $n \rightarrow \infty$ if the following conditions hold:*

$$\begin{aligned} t &< \min(r, 1 - r, p + 1 - 2(q + r), p - 2, 2q + r - 2) & (15) \\ q + r &> 1. & (16) \end{aligned}$$

Our proof that utilizes the survival/contamination analysis framework from Muthukumar et al. [56] along with a typicality-inspired argument where the feature margin (difference between largest and second largest feature) on the test point plays a key role.

Assume without loss of generality that for the test point $\mathbf{x}_{test} \sim \mathcal{N}(0, I_d)$, the true class is α for some $\alpha \in [k]$. Let \mathbf{x}_{test}^w be the weighted version of this test point. A necessary and sufficient condition for classification error is that for some $\beta \neq \alpha, \beta \in [k]$, the score associated with class β is higher than the score associated with class α . Pulling out the key terms associated with the α and β weighted features, we get:

$$\begin{aligned} \hat{f}_\alpha[\alpha]x_{test}^w[\alpha] + \hat{f}_\alpha[\beta]x_{test}^w[\beta] + \sum_{j \notin \{\alpha, \beta\}} \hat{f}_\alpha[j]x_{test}^w[j] &< \hat{f}_\beta[\alpha]x_{test}^w[\alpha] \\ &+ \hat{f}_\beta[\beta]x_{test}^w[\beta] + \sum_{j \notin \{\alpha, \beta\}} \hat{f}_\beta[j]x_{test}^w[j] \end{aligned} \quad (24)$$

$$\Rightarrow (\hat{f}_\alpha[\alpha] - \hat{f}_\beta[\alpha])x_{test}^w[\alpha] - (\hat{f}_\beta[\beta] - \hat{f}_\alpha[\beta])x_{test}^w[\beta] < \sum_{j \notin \{\alpha, \beta\}} (\hat{f}_\beta[j] - \hat{f}_\alpha[j])x_{test}^w[j]. \quad (25)$$

$$(26)$$

Note that $\sum_{j \notin \{\alpha, \beta\}}$ refers to the sum over all feature indices 1 to d excluding α and β .

Table 1: Notation

Symbol	Definition	Dimension	Source
k	Number of classes	Scalar	Sec. 4
n	Number of training points	Scalar	Sec. 4
d	Dimension of each point — the total number of features	Scalar	Sec. 4
s	The number of favored features	Scalar	Def. 4.2
p	Parameter controlling overparameterization ($d = n^p$)	Scalar	Def. 4.2
r	Parameter controlling the number of favored features ($s = n^r$)	Scalar	Def. 4.2
a	Parameter controlling the favored weights ($a = n^{-a}$)	Scalar	Def. 4.2
t	Parameter controlling the number of classes ($k = c_k n^t$)	Scalar	Def. 4.2
c_k	The number of classes when $t = 0$ ($k = c_k n^t$)	Scalar	Def. 4.2
λ_j	Squared weight of the j th feature	Scalar	Def. 4.2
\mathbf{x}_i	i th training point (unweighted)	Length- n vector	Eqn. 1
ℓ_i	Class label of i th training point	Scalar	Eqn. 2
\mathbf{x}_i^w	i th training point (weighted)	Length- n vector	Eqn. 3
\mathbf{X}^w	Weighted feature matrix	$(n \times d)$ -matrix	Eqn. 4
\mathbf{z}_j	The collected j th features of all training points	Length- n vector	Eqn. 4
\mathbf{y}_m^{oh}	One-hot encoding of all the training points for label m	Length- n vector	Eqn. 6
\mathbf{Y}^{oh}	One-hot label matrix	$(n \times k)$ -matrix	Eqn. 5
\mathbf{y}_m	Zero-mean encoding of the training points for label m	Length- n vector	Eqn. 7
$\hat{\mathbf{f}}_m$	Learned coefficients for label m using min-norm interpolation	Length- d vector	Eqn. 10
\mathbf{x}_{test}	A single test point	Length- d vector	Sec. 4
\mathbf{x}_{test}^w	A single weighted test point	Length- d vector	Sec. 4
\mathbf{A}	$\mathbf{A} = \mathbf{X}^w (\mathbf{X}^w)^>$	$(n \times n)$ -matrix	Eqn. 38
$\mu_i(\mathbf{A})$	The i th eigenvalue of matrix \mathbf{A} , sorted in descending order	Scalar	App. B
$\mathbf{\Lambda}$	Matrix of squared feature weights: $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$	$(d \times d)$ -matrix	App. B
$\hat{h}_{\alpha,\beta}$	Relative survival $\hat{h}_{\alpha,\beta}[j] = \lambda_j^{-1/2}(\hat{f}_\alpha[j] - \hat{f}_\beta[j])$	Length- d vector	Eqn. 28
$\text{CN}_{\alpha,\beta}$	Normalizing factor $\text{CN}_{\alpha,\beta} = \sqrt{\left(\sum_{j \notin \{\alpha,\beta\}} \lambda_j^2 (\hat{h}_{\beta,\alpha}[j])^2\right)}$	Scalar	Eqn. 34
k_{ψ_2}	The sub-Gaussian norm of a scalar random variable	Scalar	Eqn. 74
$\bar{\mu}$	Center of the eigenvalue bounds for \mathbf{A}^{-1} , $\bar{\mu} = \frac{1}{\sum_j \lambda_j}$	Scalar	Eqn. 41
	Deviation term in eigenvalue bounds for \mathbf{A}	Scalar	Eqn. 77
Δ_μ	Deviation term in eigenvalue bounds for \mathbf{A}^{-1}	Scalar	Eqn. 43

By converting into the unweighted feature space we obtain,

$$\lambda_\alpha \hat{h}_{\alpha,\beta}[\alpha] x_{test}[\alpha] - \lambda_\beta \hat{h}_{\beta,\alpha}[\beta] x_{test}[\beta] < \sum_{j \notin \{\alpha,\beta\}} \lambda_j \hat{h}_{\beta,\alpha}[j] x_{test}[j], \quad (27)$$

where we introduce the short-hand notation,

$$\hat{h}_{\alpha,\beta}[j] = \lambda_j^{-1/2}(\hat{f}_\alpha[j] - \hat{f}_\beta[j]) \quad (28)$$

$$\hat{h}_{\beta,\alpha}[j] = \lambda_j^{-1/2}(\hat{f}_\beta[j] - \hat{f}_\alpha[j]). \quad (29)$$

Since both α and β are favored feature indices, by leveraging the definition of the bi-level model and denoting $\lambda_\alpha = \lambda_\beta = \lambda$, we get

$$\lambda \left(\hat{h}_{\alpha,\beta}[\alpha] x_{test}[\alpha] - \hat{h}_{\beta,\alpha}[\beta] x_{test}[\beta] \right) < \sum_{j \notin \{\alpha,\beta\}} \lambda_j \hat{h}_{\beta,\alpha}[j] x_{test}[j]. \quad (30)$$

Next, we perform some algebraic manipulations,

$$\lambda \left(\widehat{h}_{\alpha,\beta}[\alpha] x_{test}[\alpha] \quad \widehat{h}_{\beta,\alpha}[\beta] x_{test}[\beta] \right) < \sum_{j \notin \mathcal{F}_{\alpha,\beta g}} \lambda_j \widehat{h}_{\beta,\alpha}[j] x_{test}[j] \quad (31)$$

$$\Rightarrow \lambda \widehat{h}_{\alpha,\beta}[\alpha] (x_{test}[\alpha] \quad x_{test}[\beta]) + \lambda x_{test}[\beta] (\widehat{h}_{\alpha,\beta}[\alpha] \quad \widehat{h}_{\beta,\alpha}[\beta]) < \sum_{j \notin \mathcal{F}_{\alpha,\beta g}} \lambda_j \widehat{h}_{\beta,\alpha}[j] x_{test}[j] \quad (32)$$

$$\Rightarrow \lambda \widehat{h}_{\alpha,\beta}[\alpha] \left((x_{test}[\alpha] \quad x_{test}[\beta]) + x_{test}[\beta] \frac{\widehat{h}_{\alpha,\beta}[\alpha] \quad \widehat{h}_{\beta,\alpha}[\beta]}{\widehat{h}_{\alpha,\beta}[\alpha]} \right) < \sum_{j \notin \mathcal{F}_{\alpha,\beta g}} \lambda_j \widehat{h}_{\beta,\alpha}[j] x_{test}[j]. \quad (33)$$

We divide both sides by the quantity $\text{CN}_{\alpha,\beta}$ defined as,

$$\text{CN}_{\alpha,\beta} = \sqrt{\left(\sum_{j \notin \mathcal{F}_{\alpha,\beta g}} \lambda_j^2 (\widehat{h}_{\beta,\alpha}[j])^2 \right)}. \quad (34)$$

This normalizes the RHS of (33) to have a standard normal distribution. Thus, the necessary and sufficient condition for a misclassification error is for some $\beta \notin \alpha, \beta \geq [k]$,

$$\frac{\lambda \widehat{h}_{\alpha,\beta}[\alpha]}{\text{CN}_{\alpha,\beta}} \left((x_{test}[\alpha] \quad x_{test}[\beta]) + x_{test}[\beta] \frac{\widehat{h}_{\alpha,\beta}[\alpha] \quad \widehat{h}_{\beta,\alpha}[\beta]}{\widehat{h}_{\alpha,\beta}[\alpha]} \right) < \frac{1}{\text{CN}_{\alpha,\beta}} \sum_{j \notin \mathcal{F}_{\alpha,\beta g}} \lambda_j \widehat{h}_{\beta,\alpha}[j] x_{test}[j]. \quad (35)$$

A sufficient condition for correct classification can then be obtained by ensuring that the smallest potential value of the LHS is still greater than the value of the RHS for all values of β . Thus, we obtain a sufficient condition for correct classification by appropriately minimizing or maximizing quantities over competing feature indices $\beta \notin \alpha, \beta \geq [k]$ (for notational convenience we simply denote this as \min_{β} or \max_{β}).

$$\begin{aligned} & \underbrace{\frac{\min_{\beta} \lambda \widehat{h}_{\alpha,\beta}[\alpha]}{\max_{\beta} \text{CN}_{\alpha,\beta}}}_{\text{SU/CN ratio}} \left(\underbrace{\min_{\beta} (x_{test}[\alpha] \quad x_{test}[\beta])}_{\text{closest feature margin}} \quad \underbrace{\max_{\beta} j x_{test}[\beta] j}_{\text{largest competing feature}} \quad \underbrace{\max_{\beta} \left| \frac{\widehat{h}_{\alpha,\beta}[\alpha] \quad \widehat{h}_{\beta,\alpha}[\beta]}{\widehat{h}_{\alpha,\beta}[\alpha]} \right|}_{\text{survival variation}} \right) \\ & > \underbrace{\max_{\beta} \frac{1}{\text{CN}_{\alpha,\beta}} \left(\sum_{j \notin \mathcal{F}_{\alpha,\beta g}} \lambda_j \widehat{h}_{\beta,\alpha}[j] x_{test}[j] \right)}_{\text{normalized contamination}}. \end{aligned} \quad (36)$$

We will show that under the conditions specified in Theorem 5.1, with sufficiently high probability, the relevant survival to contamination SU/CN ratio grows at a polynomial rate n^v for some $v > 0$, the closest feature margin shrinks at a less-than-polynomial rate $1/\ln nk$, and the survival variation decays at a polynomial rate n^{-u} for some $u > 0$. Further, the magnitudes of the largest competing feature and the normalized contamination are no more than $2\sqrt{\ln(nk)}$. Here, we leverage the idea of typicality-style proofs in information theory [17] to avoid unnecessarily loose union bounds that end up being dominated by the atypical behavior of quantities. In our case, by pulling the feature margin out explicitly, we can just deal with its typical behavior. Similarly, the typical behavior of the largest competing feature and the true feature is all that matters. Before we proceed with the rest of our proof we remind the reader of a few important definitions.

Recall from (10) that our learned feature coefficients are

$$\hat{\mathbf{f}}_m = (\mathbf{X}^w)^{\succ} (\mathbf{X}^w (\mathbf{X}^w)^{\succ})^{-1} \mathbf{y}_m. \quad (37)$$

Let

$$\mathbf{A} = \mathbf{X}^w (\mathbf{X}^w)^{\top}. \quad (38)$$

Then we can express our learned coefficients as

$$\hat{f}_m[j] = \sqrt{\lambda_j} \mathbf{z}_j^{\top} \mathbf{A}^{-1} \mathbf{y}_m, \quad (39)$$

where $\mathbf{z}_j \in \mathbb{R}^n$ contains the j^{th} features of all n training points. The rows of \mathbf{X}^w are i.i.d. Gaussians with covariance matrix $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$. Let $\mu_1(\mathbf{A})$ denote the largest eigenvalue and $\mu_n(\mathbf{A})$ denote the smallest eigenvalue of \mathbf{A} respectively, with $\mu_i(\mathbf{A})$ being the i -th largest eigenvalue of \mathbf{A} .

Next, we state a useful lemma adapted from Bartlett et al. [4]¹² that bounds the eigenvalues of \mathbf{A}^{-1} . Subsequent lemmas will utilize these eigenvalue bounds.

Lemma B.1. (Eigenvalue bounds on \mathbf{A}^{-1} adapted from Bartlett et al. [4]): If $\mathbf{\Lambda}$ is such that $\sum_j \lambda_j \geq \bar{\mu}$, then with probability at least $(1 - 2e^{-n})$,

$$\bar{\mu} - \Delta_{\mu} \leq \mu_n(\mathbf{A}^{-1}) \leq \mu_1(\mathbf{A}^{-1}) \leq \bar{\mu} + \Delta_{\mu}, \quad (40)$$

where,

$$\bar{\mu} = \frac{1}{\sum_j \lambda_j} \quad (41)$$

$$= \frac{32}{9} \left(\lambda_1 (1 + \ln 9) n + \sqrt{(1 + \ln 9) n \sum_j \lambda_j^2} \right) \quad (42)$$

$$\Delta_{\mu} = \bar{\mu} \left(\frac{1}{\sum_j \lambda_j} + \Theta \left(\frac{1}{\sum_j \lambda_j} \right)^2 \right). \quad (43)$$

Further this implies that with probability at least $(1 - 2e^{-n})$,

$$|\mu_i(\mathbf{A}^{-1}) - \bar{\mu}| \leq \Delta_{\mu} \quad (44)$$

for all $i \in [n]$.

The subsequent lemmas bound the feature margin, survival, contamination and survival variation terms, utilizing tools from [4] and building on results from [56].

Lemma B.2. (Lower bound on the closest feature margin as $k \rightarrow \infty$): For any constant $\varepsilon > 0$, there exists a constant θ such that, for sufficiently large k with probability at least $(1 - \varepsilon)$,

$$\min_{\beta:1} \min_{\beta \neq \alpha} (x_{\text{test}}[\alpha] - x_{\text{test}}[\beta]) \geq \frac{\theta}{\sqrt{2 \ln(k)}}. \quad (45)$$

Here, α is fixed and corresponds to the index of the true class — i.e. α corresponds to the index of the maximum feature among the first k features.

Lemma B.3. (Lower bound on the closest feature margin when k is constant): If $k = c_k$ for some fixed constant c_k , for any constant $\varepsilon > 0$, there exists a constant $\varepsilon^{\theta} > 0$ such that

$$\Pr \left(\min_{\beta, \gamma:1} \min_{\beta \neq \gamma} |x_{\text{test}}[\beta] - x_{\text{test}}[\gamma]| \geq \varepsilon^{\theta} \right) \geq 1 - \varepsilon. \quad (46)$$

Thus, with probability at least $(1 - \varepsilon)$,

$$\min_{\beta:1} \min_{\beta \neq \alpha} (x_{\text{test}}[\alpha] - x_{\text{test}}[\beta]) \geq \varepsilon^{\theta}. \quad (47)$$

Here, α is fixed and corresponds to the index of the true class — i.e. α corresponds to the index of the maximum feature among the first k features.

¹²More precisely this lemma appeared in the first version of this work at <https://arxiv.org/pdf/1906.11300v1.pdf>. In subsequent versions the authors use a slightly weaker version of this result since it is sufficient for their purposes.

Lemma B.4. (Lower bound on relative survival of true feature): For any fixed $\beta \geq [k]$, $\beta \notin \alpha$, with $\lambda_\alpha = \lambda_\beta = \lambda$ we have with probability at least $(1 - 5/(nk))$,

$$\widehat{h}_{\alpha,\beta}[\alpha] \geq \lambda \left(c_{10} \bar{\mu} \frac{n}{k} \sqrt{\ln(k)} - c_9 (\bar{\mu} \frac{p}{n} \sqrt{\ln(nk)} + \Delta_\mu \frac{n}{k}) \right), \quad (48)$$

for universal positive constants c_9 and c_{10} .

By substituting the asymptotic behavior of parameters from our bi-level ensemble model we get the following corollary:

Corollary B.4.1. Under the bi-level ensemble model 4.2, for any fixed $\beta \geq [k]$, $\beta \notin \alpha$, $\lambda_\alpha = \lambda_\beta = \lambda$ if $t < 1/2$, $t < 2(q + r - 1)$ and $1 < q + r < (p + 1)/2$, with probability at least $(1 - 5/(nk))$,

$$\lambda \widehat{h}_{\alpha,\beta}[\alpha] \geq c_{12} n^{1 - q - r - t} \sqrt{\ln(k)}, \quad (49)$$

for universal positive constant c_{12} .

Lemma B.5. (Upper bound on contamination): For any fixed $\beta \geq [k]$, $\beta \notin \alpha$, with probability at least $(1 - 7/(nk))$,

$$\text{CN}_{\alpha,\beta} \leq c_7 \left(\bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(ndk)} + \Delta_\mu \frac{n}{k} \sqrt{\sum \lambda_j^2} \right), \quad (50)$$

for universal positive constant c_7 .

As before, for our bi-level ensemble model we have the corollary:

Corollary B.5.1. Under the bi-level model 4.2, in the regime $1 < q + r < (p + 1)/2$, with probability at least $(1 - 7/(nk))$,

$$\text{CN}_{\alpha,\beta} \leq c_{13} n^{(1 - t - p)/2 + \max(0, 3/2 - q - r) + \max(0, p/2 - q - r/2)} \sqrt{\ln(ndk)}, \quad (51)$$

for universal positive constant c_{13} .

Lemma B.6. (Upper bound on survival variance): For any fixed competing feature $\beta \geq [k]$, $\beta \notin \alpha$ with $\lambda_\alpha = \lambda_\beta$, we have with probability at least $(1 - 15/(nk))$,

$$\frac{\widehat{h}_{\alpha,\beta}[\alpha] - \widehat{h}_{\beta,\alpha}[\beta]}{\widehat{h}_{\alpha,\beta}[\alpha]} \leq \frac{2c_9 (\bar{\mu} \frac{p}{n} \sqrt{\ln(nk)} + \Delta_\mu \frac{n}{k})}{c_{10} \bar{\mu} \frac{n}{k} \sqrt{\ln(k)} - c_9 (\bar{\mu} \frac{p}{n} \sqrt{\ln(nk)} + \Delta_\mu \frac{n}{k})}, \quad (52)$$

for universal positive constants c_9 and c_{10} .

As before, we can also obtain the asymptotic bound:

Corollary B.6.1. Under the bi-level ensemble model 4.2, for any fixed $\beta \geq [k]$, $\beta \notin \alpha$, if $t < 1/2$, $t < 2(q + r - 1)$, and $1 < q + r < (p + 1)/2$, with probability at least $(1 - 15/(nk))$,

$$\frac{\widehat{h}_{\alpha,\beta}[\alpha] - \widehat{h}_{\beta,\alpha}[\beta]}{\widehat{h}_{\alpha,\beta}[\alpha]} < n^{-u}, \quad (53)$$

for large enough n for some fixed $u > 0$.

Next, we assume that the lemmas and corollaries stated above are true and complete the proof for Theorem 5.1. We provide proofs for these lemmas in Appendices C, D and E.

Assume we are in the regime where $t < 1/2$, $t < 2(q + r - 1)$, and $1 < q + r < (p + 1)/2$, so all our corollaries above hold. Denote the misclassification event as E_{err} and let $\varepsilon > 0$ be an arbitrarily chosen constant.

Substitute Corollaries B.4.1, B.5.1, and B.6.1 into (22), applying them on all $1 \leq \beta \notin \alpha \leq k$. They hold with probability at least $(1 - 5/(nk))$, $(1 - 7/(nk))$, and $(1 - 15/(nk))$ respectively for a given test point and choice of β . So by the union bound across the three bounds and all $k - 1$ choices of β , with probability at most $27/n$, one of these corollaries will not hold for our test point for some β . Let this failure event be denoted E_1 .

In the case when E_1 does not occur, misclassification occurs only if

$$\frac{c_{12} \sqrt{\ln(k)}}{c_7 \sqrt{\ln(ndk)}} n^v \left(\min_{\beta} (x_{test}[\alpha] - x_{test}[\beta]) - \max_{\beta} |x_{test}[\beta] - n^{-u}| \right) < \max_{\beta} Z^{(\beta)}, \quad (54)$$

where we define the exponent

$$v = 1 - q - r - t - (1 - t - p)/2 - \max\left(0, \frac{3}{2} - q - r\right) - \max\left(0, \frac{p}{2} - q - \frac{q}{2}\right) \quad (55)$$

$$= \frac{p+1}{2} - q - r - \frac{t}{2} - \max\left(0, \frac{3}{2} - q - r, \frac{p}{2} - q - \frac{r}{2}, \frac{3}{2} - 2q - \frac{3r}{2}\right), \quad (56)$$

and

$$Z^{(\beta)} = \frac{1}{\text{CN}_{\alpha, \beta}} \left(\sum_{j \notin \mathcal{F}_{\alpha, \beta g}} \lambda_j \hat{h}_{\beta, \alpha}[j] x_{test}[j] \right). \quad (57)$$

For each class β , observe that we have $Z^{(\beta)} \sim \mathcal{N}(0, 1)$.¹³ Thus, by the Gaussian tail bound, for each β with probability at least $(1 - 1/(nk))$,

$$Z^{(\beta)} < \sqrt{2 \ln(nk)}. \quad (58)$$

So by the union bound over all k classes β , with probability at least $(1 - 1/n)$,

$$\max_{\beta} Z^{(\beta)} < \sqrt{2 \ln(nk)}. \quad (59)$$

Let the failure event where this is not the case be E_2 .

An identical argument shows that with probability at least $(1 - 2/n)$, $\max_{\beta} j x_{test}[\beta] j < \sqrt{2 \ln(nk)}$. Let E_3 be the failure event where this is not the case.

From Lemma B.2, we know with probability $1 - \varepsilon$ that, if $t > 0$, then for sufficiently large n (and so sufficiently large k)

$$\min_{\beta} (x_{test}[\alpha] - x_{test}[\beta]) > \frac{\theta}{\sqrt{2 \ln(k)}}. \quad (60)$$

If $t = 0$ and $k = c_k$, then Lemma B.3 states that, with probability $1 - \varepsilon$,

$$\Pr \left(\min_{\beta \in \gamma} j x_{test}[\beta] j - x_{test}[\gamma] j > \varepsilon^{\theta} \right) < 1 - \varepsilon, \quad (61)$$

for some constant ε^{θ} . Let the ε -probability event of the appropriate margin bound (depending on whether $t = 0$ or $t > 0$) being violated be the error event E_4 .

Assuming E_1, E_2, E_3 , and E_4 all do not take place, misclassification can only occur if

$$\frac{c_{12} \sqrt{\ln(k)}}{c_7 \sqrt{\ln(ndk)}} n^v \left(\min \left(1 - \varepsilon, \frac{\theta}{\sqrt{2 \ln(k)}} \right) - \sqrt{2 \ln(nk)} n^{-u} \right) < \sqrt{2 \ln(nk)}. \quad (62)$$

Clearly, if $v > 0$, then (for sufficiently large n) misclassification becomes asymptotically impossible (except via the specified error events), since the LHS of the above grows asymptotically faster than the RHS.

The union bound shows that the probability of any of E_1, E_2, E_3, E_4 occurring tends to ε as $n \rightarrow \infty$ (since the probability of the first three tend to zero). So in the regime where

$$t < \frac{1}{2} \quad (63)$$

$$t < 2(q + r - 1) \quad (64)$$

$$q + r > 1 \quad (65)$$

$$\frac{p+1}{2} > q + r + \frac{t}{2} + \max\left(0, \frac{3}{2} - q - r\right) + \max\left(0, \frac{p}{2} - q - \frac{r}{2}\right), \quad (66)$$

¹³To be precise, here we can think of fixing the training data and looking purely at the randomness arising from the features in the test point. The resulting $Z^{(\beta)}$ is a standard normal. Since we are using the union bound in our proof finally, this is sufficient for our purposes.

the probability of misclassification tends to ε for sufficiently large n , for any $\varepsilon > 0$.

Consolidation of the above bounds produces the conditions ¹⁴

$$t < \min(1 - r, p + 1 - 2(q + r), p - 2, 2q + r - 2) \quad (71)$$

$$q + r > 1. \quad (72)$$

Finally, note that the condition $t < r$ comes from the definition of the bi-level model (4.2). This condition simply states that for good generalization we must favor all the features used to determine classes. Since the analysis above holds for any ε , we see that within this regime the probability of misclassification must approach zero in the limit. This completes the proof. Note that while we show that probability of misclassification goes to zero, we do not show it to do so at any particular rate, because the result from Lemma B.2 does not specify the rate of convergence.

C Useful results from elsewhere that we need

This section collects results that are used in our proof, but which come from elsewhere or are lightly adapted to our purposes.

Hanson-Wright inequality [65]: Let \mathbf{z} be a random vector composed of i.i.d. random variables that are zero mean and with sub-Gaussian norm at most K . The sub-Gaussian norm $k\xi k_{\psi_2}$ of a random variable ξ is defined as in Rudelson and Vershynin [65],

$$k\xi k_{\psi_2} = \inf_{K>0} K \quad (73)$$

$$\text{s.t. } \mathbb{E} \exp(\xi^2/K^2) \leq 2. \quad (74)$$

Then, there exists universal constant $c > 0$ such that for any positive semi-definite matrix M and for every $t \geq 0$, we have

$$\Pr [|\mathbf{z}^T M \mathbf{z} - \mathbb{E}[\mathbf{z}^T M \mathbf{z}]| > t] \leq 2 \exp \left\{ -c \min \left\{ \frac{t^2}{K^4 \text{tr} M}, \frac{t}{K^2 \text{tr} M} \right\} \right\} \quad (75)$$

The next result bounds the eigenvalues of the $n \times n$ matrix $\mathbf{A} = \mathbf{X}^w (\mathbf{X}^w)^T$, where recall that the rows of \mathbf{X}^w are i.i.d. Gaussians with covariance matrix $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$. Let $\mu_1(\mathbf{A})$ denote the largest eigenvalue and $\mu_n(\mathbf{A})$ denote the smallest eigenvalue of \mathbf{A} respectively.

From Bartlett et al. [4], we have the following result

Lemma C.1. *With probability at least $(1 - 2e^{-n})$, the eigenvalues of \mathbf{A} satisfy:*

$$\sum_j \lambda_j - \mu_n(\mathbf{A}) \leq \mu_1(\mathbf{A}) - \sum_j \lambda_j, \quad (76)$$

where,

$$= \frac{32}{9} \left(\lambda_1 (1 + \ln 9) n + \sqrt{(1 + \ln 9) n \sum_j \lambda_j^2} \right). \quad (77)$$

¹⁴We can simplify (66) as follows:

$$\frac{p+1}{2} > q+r + \frac{t}{2} \Rightarrow t < p+1 - 2(q+r) \quad (67)$$

$$\frac{p+1}{2} > q+r + \frac{t}{2} + \frac{3}{2} - q - r \Rightarrow t < p - 2 \quad (68)$$

$$\frac{p+1}{2} > q+r + \frac{t}{2} + \frac{p}{2} - q - \frac{r}{2} \Rightarrow t < 1 - r \quad (69)$$

$$\frac{p+1}{2} > q+r + \frac{t}{2} + \frac{3}{2} - q - r + \frac{p}{2} - q - \frac{r}{2} \Rightarrow t < 2q+r - 2. \quad (70)$$

Then we note that $t < \min(r, 1 - r) \Rightarrow t < 1/2$.

Next, as stated previously in Lemma B.1 we will use this result to obtain bounds on the eigenvalues of \mathbf{A}^{-1} assuming that $\mathbf{\Lambda}$ is such that $\sum_j \lambda_j$.¹⁵

Lemma B.1. (Eigenvalue bounds on \mathbf{A}^{-1} adapted from Bartlett et al. [4]):

If $\mathbf{\Lambda}$ is such that $\sum_j \lambda_j$, then with probability at least $(1 - 2e^{-n})$,

$$\bar{\mu} - \Delta_\mu \leq \mu_n(\mathbf{A}^{-1}) \leq \mu_1(\mathbf{A}^{-1}) \leq \bar{\mu} + \Delta_\mu, \quad (40)$$

where,

$$\bar{\mu} = \frac{1}{\sum_j \lambda_j} \quad (41)$$

$$= \frac{32}{9} \left(\lambda_1(1 + \ln 9)n + \sqrt{(1 + \ln 9)n \sum_j \lambda_j^2} \right) \quad (42)$$

$$\Delta_\mu = \bar{\mu} \left(\frac{1}{\sum_j \lambda_j} + \Theta \left(\frac{1}{\sum_j \lambda_j} \right)^2 \right). \quad (43)$$

Further this implies that with probability at least $(1 - 2e^{-n})$,

$$|\mu_i(\mathbf{A}^{-1}) - \bar{\mu}| \leq \Delta_\mu \quad (44)$$

for all $i \in [n]$.

Proof. Let $S = \sum_j \lambda_j$.

$$\frac{1}{S + \Delta_\mu} = \frac{1}{S} \left(1 + \frac{\Delta_\mu}{S} \right)^{-1} \quad (78)$$

$$= \frac{1}{S} \left(1 - \frac{\Delta_\mu}{S} + \Theta \left(\frac{\Delta_\mu}{S} \right)^2 \right) \quad (79)$$

$$= \bar{\mu} - \Delta_\mu, \quad (80)$$

and analogously $(S - \Delta_\mu)^{-1} = \bar{\mu} + \Delta_\mu$. Taking reciprocals of everything in the inequality 76, and since the eigenvalues of \mathbf{A} and \mathbf{A}^{-1} are reciprocals of each other, the desired result follows. \square

As a Corollary of Lemma B.1:

Corollary C.1.1. (Asymptotic eigenvalue bounds on \mathbf{A}^{-1}) Considering the asymptotic scaling of the model parameters from the bi-level model (Definition 4.2), in the regime $1 < q + r < (1 + p)/2$,

$$\bar{\mu} = n^{-p} \quad (81)$$

$$\Delta_\mu = c_4 n^{1-p-q-r} \bar{\mu}, \quad (82)$$

where $\bar{\mu}$ and Δ_μ are defined as in Lemma B.1, and c_4 is a universal constant.

Proof. From the asymptotic scaling of the λ_j from (13) and (14), we see that (from the definition provided in Lemma B.1)

$$\bar{\mu} = \frac{1}{\sum_j \lambda_j} \quad (83)$$

$$= \frac{1}{n^r n^{p-q-r} + (n^p - n^r)(1 - n^q) n^p / (n^p - n^r)} \quad (84)$$

$$= \frac{1}{n^{p-q} + n^p - n^{p-q}} \quad (85)$$

$$= n^{-p}. \quad (86)$$

¹⁵Note that in the regime $q + r < 1$ (where regression works [56]), we do not have $\sum_j \lambda_j$ and in such scenarios we cannot simply rely on eigenvalue bounds and need to use other techniques in the proof.

Next, we have that

$$= \frac{32}{9} \left(\lambda_1(1 + \ln 9)n + \sqrt{(1 + \ln 9)n \sum_j \lambda_j^2} \right) \quad (87)$$

$$c_1 n^{1+p-q-r} + c_2 \sqrt{n(n^r n^{2p-2q-2r} + (n^p - n^r))} \quad (88)$$

$$c_1 n^{1+p-q-r} + c_2 \sqrt{n^{1+2p-2q-r} + n^{1+p}} \quad (89)$$

for constants c_1 and c_2 ,

The second term is of the order $n^{\max((1-r)/2+p-q, (1+p)/2)}$. Thus, in the regime $q+r < (1+p)/2$, and since $r < 1$ we have $1+p-q-r > (1-r)/2+p-q$ and $1+p-q-r > (1+p)/2$ and the first term dominates.

Thus, $c_3 n^{1+p-q-r}$ for some constant c_3 and sufficiently large n .

Observe that since $q+r > 1$, $\sum_j \lambda_j = n^p$. Thus, we can substitute into our relation for Δ_μ from Lemma B.1, to see that

$$\Delta_\mu = \bar{\mu} \left(\frac{1}{\sum_j \lambda_j} + \Theta \left(\frac{1}{\sum_j \lambda_j} \right)^2 \right) \quad (90)$$

$$n^{-p} ((c_3 n^{1+p-q-r})(n^{-p}) + \Theta((c_3 n^{1+p-q-r})^2 (n^{-p})^2)) \quad (91)$$

$$= n^{-p} (c_3 n^{1-q-r} + \Theta(c_3 n^{2(1-q-r)})). \quad (92)$$

In the regime where $q+r > 1$, the first term in the sum dominates the second, giving us,

$$\Delta_\mu = c_4 n^{1-p-q-r} \quad (93)$$

for some constant c_4 and sufficiently large n . This completes the proof. \square

Finally, in this section, we restate well-known bounds concerning Gaussian random variables.

Lemma C.2. *Chi-squared tail bound:*

Let $\mathbf{z} \sim \mathcal{N}(0, I_n)$. For any $\delta \in (0, 1)$, with probability at least $(1 - 2e^{-n\delta^2})$ we have:

$$n(1 - \delta) \leq \|\mathbf{z}\|^2 \leq n(1 + \delta). \quad (94)$$

From bounds on the expectation of the maximum of k Gaussians:

Lemma C.3. Let $\mathbf{z}_\alpha = \max_{1 \leq j \leq k} \mathbf{z}_j$ where $\mathbf{z}_j \sim \mathcal{N}(0, 1)$. Then,

$$\frac{1}{\pi \ln 2} \leq \frac{\rho}{\ln k} \leq \mathbb{E}[\mathbf{z}_\alpha] \leq \frac{\rho}{2} \leq \frac{\rho}{\ln k}. \quad (95)$$

D Utility Bounds

The big technical challenge in moving from binary classification (as studied in Muthukumar et al. [56]) to multiclass classification has to do with the nature of the training data. Whereas for binary classification one could change coordinates so that the binary labels only depended on a single Gaussian random variable and were independent of all other directions of Gaussian variation in the covariates, no such change of coordinates exists for multiclass labels. The one-hot-style encoding of the labels fundamentally depends on the realizations of all k of the Gaussian random variables representing each of the k classes. This means that we can no longer simply leverage independence to simplify the analysis and certain clever approaches used to invoke Hanson-Wright are no longer available to us. However, the need remains to appropriately bound quadratic forms of the form $|\mathbf{z}_j^\top \mathbf{A}^{-1} \Delta y|$ both for the cases when j represents a feature that is not dominant in the computation of Δy as well as in cases where j represents a feature that is dominant in Δy . To be able to control such quantities in the absence of the independence we could leverage in the binary case, this section of the Appendix derives two lemmas which can be viewed as helper bounds. These bounds will later be used to bound the various quantities from (22). Because our focus is on the asymptotic scaling, we will use c_i to denote the appropriate global constants.

In the subsequent lemmas, $\bar{\mu}$ and Δ_μ are defined as in the bounds on the eigenvalues of \mathbf{A}^{-1} from Lemma B.1.

The following lemma is used to upper-bound the contamination term $\text{CN}_{\alpha,\beta}$ in Lemma B.5:

Lemma D.1. *Let $\Delta\mathbf{y} = \mathbf{y}_\alpha - \mathbf{y}_\beta$. Let α, β , and j be distinct. Then, with probability at least $(1 - 7/(ndk))$, we have,*

$$|\mathbf{z}_j^\top \mathbf{A}^{-1} \Delta\mathbf{y}| \leq c_7(\bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(ndk)} + \Delta_\mu \sqrt{n/k}), \quad (96)$$

for some constant c_7 .

This next lemma is used to bound the numerator of the survival variation term from (22):

Lemma D.2. *Let $\Delta\mathbf{y} = \mathbf{y}_\alpha - \mathbf{y}_\beta$. With probability at least $(1 - 5/(nk))$, we have each of*

$$\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta\mathbf{y} \leq \bar{\mu}(\mathbb{E}[\mathbf{z}_\alpha^\top \mathbf{y}_\alpha] - \mathbb{E}[\mathbf{z}_\alpha^\top \mathbf{y}_\beta]) + c_9(\bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(nk)} + \Delta_\mu \sqrt{n/k}) \quad (97)$$

$$\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta\mathbf{y} \geq \bar{\mu}(\mathbb{E}[\mathbf{z}_\alpha^\top \mathbf{y}_\alpha] - \mathbb{E}[\mathbf{z}_\alpha^\top \mathbf{y}_\beta]) - c_9(\bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(nk)} + \Delta_\mu \sqrt{n/k}), \quad (98)$$

for some constant c_9 .

The following corollary of the above is used to lower-bound the relative survival $\hat{h}_{\alpha,\beta}[\alpha]$, which in turn bounds the SU/CN ratio and the denominator of the survival variation term:

Corollary D.2.1. *Let $\Delta\mathbf{y} = \mathbf{y}_\alpha - \mathbf{y}_\beta$. With probability at least $(1 - 5/(nk))$, we have,*

$$\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta\mathbf{y} \geq c_{10} \bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(k)} - c_9(\bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(nk)} + \Delta_\mu \sqrt{n/k}), \quad (99)$$

for some constant c_{10} .

D.1 Proof of Lemma D.1

We will write $\mathbf{A}^{-1} = \bar{\mu} \mathbf{I}_n + \Delta A_{inv}$, and split up the expression $\mathbf{z}_j^\top \mathbf{A}^{-1} \Delta\mathbf{y}$ into components involving $\bar{\mu} \mathbf{I}_n$, and components involving ΔA_{inv} . To bound the first term, we will use Hanson-Wright, and to bound the second we will use Cauchy-Schwartz. Throughout the proof, we rely on the concentration of the eigenvalues of \mathbf{A}^{-1} .

Next, we bound the first term (we set aside the constant $\bar{\mu}$ for now and deal with it later).

D.1.1 Bounds on $\mathbf{z}_j^\top (\mathbf{y}_\alpha - \mathbf{y}_\beta)$

Throughout this section, let j be a feature index distinct from α and β . Define the diagonal matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$ with diagonal entries given by:

$$M_{ii} = \begin{cases} 1, & \text{if } \Delta y[i] \neq 0 \\ 0, & \text{otherwise} \end{cases}. \quad (100)$$

In other words, M_{ii} is 1 only if training point i belongs to class α or β and is 0 otherwise. Thus for each $i \in [n]$, $M_{ii} \sim \text{Bernoulli}(2/k)$ and are independent of each other. We introduce this matrix \mathbf{M} to ensure that our bound reflects the fact that most of the entries of $\Delta\mathbf{y}$ are 0. In particular $\Delta y[i] \neq 0$ only if point i belongs to class α or β and only contains roughly $2n/k$ non-zero entries.¹⁶ Note that we have by definition,

$$\mathbf{z}_j^\top \Delta\mathbf{y} = \mathbf{z}_j^\top \mathbf{M} \Delta\mathbf{y}. \quad (101)$$

Our strategy is to bound $\mathbf{z}_j^\top \mathbf{M} \Delta\mathbf{y}$ for every typical realization \mathcal{M} of the random variable \mathbf{M} using the Hanson-Wright inequality. Subsequently, we will apply these bounds with high probability over typical realizations of \mathbf{M} that satisfy the Proposition below, which merely asserts that with high probability, the number of 1s in $\Delta\mathbf{y}$ is close to its expected value.

Proposition D.1. *For $\delta \in (0, 1)$, with probability at least $(1 - 2e^{-\frac{2n\delta^2}{3k}})$, the trace of \mathbf{M} is bounded as:*

$$(1 - \delta) \frac{2n}{k} \leq \text{Tr}(\mathbf{M}) \leq (1 + \delta) \frac{2n}{k}. \quad (102)$$

¹⁶An alternative bounding technique that first converted $\mathbf{z}_j^\top \Delta\mathbf{y}$ to a quadratic form and applied Hanson-Wright would be looser by a factor of $\sqrt{n/k}$ if we did not introduce \mathbf{M} .

Proof. Note that $\text{Tr}(\mathbf{M})$ is the sum of n i.i.d Bernoulli random variables with mean $2/k$. The result follows by application of the Chernoff bound. \square

Note that once we fix the realization \mathcal{M} , the distributions of \mathbf{z}_j and Δy will now have to be conditioned on this realization and we need to deal with the modified distributions while applying the Hanson-Wright inequality. In particular, once we know that a feature was not the winning feature, it is no longer zero-mean.

Now,

$$\mathbf{z}_j^T \mathcal{M} \Delta y = \sum_i z_j[i] \mathcal{M}_{ii} \Delta y[i] \quad (103)$$

$$= \sum_{i: \mathcal{M}_{ii}=1} z_j[i] \Delta y[i] \quad (104)$$

$$= \sum_{i: \mathcal{M}_{ii}=1} (z_j[i] - \mathbb{E}[z_j[i] | \mathcal{M}_{ii} = 1]) \Delta y[i] + \sum_{i: \mathcal{M}_{ii}=1} \mathbb{E}[z_j[i] | \mathcal{M}_{ii} = 1] \Delta y[i] \quad (105)$$

$$= \sum_{i: \mathcal{M}_{ii}=1} \tilde{z}_{j, \mathcal{M}}[i] \Delta y[i] + \sum_{i: \mathcal{M}_{ii}=1} \mathbb{E}[z_j[i] | \mathcal{M}_{ii} = 1] \Delta y[i], \quad (106)$$

where $\tilde{z}_{j, \mathcal{M}}[i]$ is now a zero-mean random variable conditioned on the realization \mathcal{M} .

First, we bound the term $\sum_{i: \mathcal{M}_{ii}=1} \tilde{z}_{j, \mathcal{M}}[i] \Delta y[i]$. We collect the elements corresponding to indices where $\mathcal{M}_{ii} = 1$ into the vectors $\mathbf{z}_{j, \mathcal{M}}^0$ and $\Delta y_{\mathcal{M}}^0$, which are both length $\text{Tr}(\mathcal{M})$ (Figure 3 shows an example of collecting elements).

$$\underbrace{\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}}_{\tilde{\mathbf{z}}_{j, \mathcal{M}}} , \underbrace{\begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}}_{\Delta y} \quad ! \quad \underbrace{\begin{bmatrix} 1 \\ 3 \\ 4 \end{bmatrix}}_{\mathbf{z}_{j, \mathcal{M}}^0} , \underbrace{\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}}_{\Delta y_{\mathcal{M}}^0}$$

Figure 3: An example of collecting elements at indices where $\mathcal{M}_{ii} = 1$ into smaller vectors of length $\text{Tr}(\mathcal{M})$. Recall that $\Delta y[i] \neq 0$ iff $\mathcal{M}_{ii} = 1$.

We can then express

$$\sum_{i: \mathcal{M}_{ii}=1} \tilde{z}_{j, \mathcal{M}}[i] \Delta y[i] \quad (107)$$

$$= (\mathbf{z}_{j, \mathcal{M}}^0)^T \Delta y_{\mathcal{M}}^0 \quad (108)$$

$$= \frac{1}{4} ((\mathbf{z}_{j, \mathcal{M}}^0 + \Delta y_{\mathcal{M}}^0)^T \mathbf{I}_{\text{Tr}(\mathcal{M})} (\mathbf{z}_{j, \mathcal{M}}^0 + \Delta y_{\mathcal{M}}^0) - (\mathbf{z}_{j, \mathcal{M}}^0 - \Delta y_{\mathcal{M}}^0)^T \mathbf{I}_{\text{Tr}(\mathcal{M})} (\mathbf{z}_{j, \mathcal{M}}^0 - \Delta y_{\mathcal{M}}^0)), \quad (109)$$

where we added and subtracted terms in the last equality.

We prove via the subsequent propositions that conditioned on the realization \mathcal{M} , the entries of $\mathbf{z}_{j, \mathcal{M}}^0$ and $\Delta y_{\mathcal{M}}^0$ are i.i.d. and sub-Gaussian with bounded norm. Thus, they satisfy the requirements to apply the Hanson-Wright inequality from Rudelson and Vershynin [65] to bound the two quadratic forms in the above expression (109).

Proposition D.2. *Conditioned on the realization \mathcal{M} , $z_{j, \mathcal{M}}^0[i^0]$ has sub-Gaussian norm at most 6.*

Proof. Let i be the original index from which $z_{j, \mathcal{M}}^0[i^0]$ was sampled.

If $j > k$, then $z_{j, \mathcal{M}}^0[i^0] = \tilde{z}_{j, \mathcal{M}}[i] = z_j[i]$ irrespective of the realization \mathcal{M} because feature j is not used in the comparison to determine the class label and is independent to y_α and y_β (and thus independent to \mathbf{M}). Further, $z_j[i]$ is simply a Gaussian (and therefore sub-Gaussian with sub-Gaussian

norm $kz_j[i]k_{\psi_2}$ 2. Here we use the definition of sub-Gaussian norm from (74) reproduced here for convenience:

The sub-Gaussian norm of a random variable ξ is given by,

$$k\xi k_{\psi_2} = \inf_{K>0} K \quad (110)$$

$$\text{s.t. } \mathbb{E} \exp(\xi^2/K^2) \leq 2. \quad (111)$$

Otherwise, if j is one of the k features that define classes, since

$$z_{j,\mathcal{M}}^0[i] = \tilde{z}_{j,\mathcal{M}}[i] \quad (112)$$

$$= z_j[i] \quad \mathbb{E}[z_j[i] j M_{ii} = 1], \quad (113)$$

the triangle inequality states that

$$k\tilde{z}_{j,\mathcal{M}}[i]k_{\psi_2} \leq kz_j[i]k_{\psi_2} + k\mathbb{E}[z_j[i] j M_{ii} = 1]k_{\psi_2}. \quad (114)$$

Note that the distribution of $z_j[i]$ conditioned on realization \mathcal{M} is equivalent to the distribution obtained by conditioning on the event $M_{ii} = 1$. So it is sufficient to compute these sub-Gaussian norms conditioned on the event $M_{ii} = 1$.

We will first bound $kz_j[i]k_{\psi_2}$. Let E_j be the event that $z_j[i]$ is the maximum out of the first k features, and let E_j^c be the complementary event.

First, without conditioning on E_j , we know by well-known results for the standard Gaussian that

$$\mathbb{E} \exp(\mathbf{z}_j[i]^2/5) = \sqrt{\frac{5}{3}} \frac{4}{3}. \quad (115)$$

Using the law of iterated expectation we can relate this to the expectation conditioned on the events E_j and E_j^c , noting that $P(E_j) = 1/k$:

$$\frac{4}{3} \mathbb{E} \exp(\mathbf{z}_j[i]^2/5) \quad (116)$$

$$= P(E_j) \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j) + P(E_j^c) \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c) \quad (117)$$

$$= \frac{1}{k} \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j) + \frac{k-1}{k} \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c). \quad (118)$$

Rearranging terms, we obtain,

$$\frac{k-1}{k} \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c) \leq \frac{4}{3} - \frac{1}{k} \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j) \quad (119)$$

$$\Rightarrow \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c) \leq \frac{k-1}{k} \left(\frac{4}{3} - \frac{1}{k} \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j) \right) \quad (120)$$

$$\frac{k-1}{k} \frac{4}{3} \quad (121)$$

$$2, \quad (122)$$

where in the second to last inequality we used the non-negativity of $\mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j)$ and in the last equality we assumed $k \geq 3$. We then have

$$\mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c) = \sum_{m \neq j} \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c \cap E_m) P(E_m | E_j^c) \quad (123)$$

$$= \frac{1}{k-1} \sum_{m \neq j} \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c \cap E_m) \quad (124)$$

where the last equality follows by symmetry. Further by symmetry, all the terms in the above summation that we are averaging are equal, so we can express it as an average of just the terms corresponding to $m = \alpha$ and $m = \beta$, as follows:

$$(124) = \frac{1}{2} \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c \cap E_\alpha) + \frac{1}{2} \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c \cap E_\beta) \quad (125)$$

$$= P(E_\alpha | E_j^c \cap (E_\alpha \cup E_\beta)) \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c \cap E_\alpha) \\ + P(E_\beta | E_j^c \cap (E_\alpha \cup E_\beta)) \mathbb{E} \exp(\mathbf{z}_j[i]^2/5|E_j^c \cap E_\beta), \quad (126)$$

again by symmetry. Since exactly one of E_α and E_β are true when conditioned on $E_j^c \setminus (E_\alpha \cup E_\beta)$, we can rewrite the above as our desired expectation

$$(126) = \mathbb{E} \exp(\mathbf{z}_j[i]^2 / 5jE_j^c \setminus (E_\alpha \cup E_\beta)) \quad (127)$$

$$= \mathbb{E} \exp(\mathbf{z}_j[i]^2 / 5jE_\alpha \cup E_\beta) \quad (128)$$

$$= \mathbb{E} \exp(\mathbf{z}_j[i]^2 / 5jM_{ii} = 1), \quad (129)$$

since $M_{ii} = 1$ is equivalent to the event $E_\alpha \cup E_\beta$. Thus, conditioned on the event $M_{ii} = 1$, $kz_j[i]k_{\psi_2} \leq \frac{\rho}{5}$.

Next we consider $k\mathbb{E}[z_j[i] \mid M_{ii} = 1]k_{\psi_2}$. By a similar argument to above, we have that $\mathbb{E}[z_j[i] \mid M_{ii} = 1] = \mathbb{E}[z_j[i] \mid E_j^c]$, so we will focus on the second quantity instead. Bounds on the max of Gaussians (Lemma C.3) state that:

$$0 < \mathbb{E}[z_j[i] \mid E_j] \leq \sqrt{2 \log(k)} \quad (130)$$

$$\Rightarrow 0 > \mathbb{E}[z_j[i] \mid E_j^c] \geq \frac{1}{k-1} \sqrt{2 \log(k)} - 2 \quad (131)$$

$$\Rightarrow \exp\left(\frac{\mathbb{E}[z_j[i] \mid E_j^c]^2}{3^2}\right) < 2. \quad (132)$$

In the second last inequality we use the fact that the function $f(k) = \left| \frac{\sqrt{2 \log k}}{k-1} \right|$ is monotonically decreasing in k and assumed $k \geq 3$.

Thus, the (constant) random variable $\mathbb{E}[z_j[i] \mid M_{ii} = 1]$ is sub-Gaussian with parameter 3. So, by the triangle inequality, conditioned on $M_{ii} = 1$

$$k\tilde{z}_{j,m}[i]k_{\psi_2} \leq kz_j[i]k_{\psi_2} + k\mathbb{E}[z_j[i] \mid M_{ii} = 1]k_{\psi_2} \quad (133)$$

$$\leq \frac{\rho}{5} + 3 \quad (134)$$

$$\leq 6. \quad (135)$$

This completes the proof that conditioned on the realization \mathcal{M} , $z_{j,\mathcal{M}}^0[i^0]$ is sub-Gaussian with norm at most 6. \square

We can now prove our target result:

Proposition D.3. *With probability at least $(1 - 6/(ndk))$,*

$$\|\mathbf{z}_j^> \Delta \mathbf{y}\| \leq c_6 \sqrt{\frac{n}{k}} \sqrt{\log(ndk)}. \quad (136)$$

for universal constant c_6 .

Proof. Our strategy will be to bound $\mathbf{z}_j^> \Delta \mathbf{y} = \mathbf{z}_j^> \mathbf{M} \Delta \mathbf{y}$ for every typical realization \mathcal{M} of \mathbf{M} that satisfies Proposition D.1. Recall that for a given realization \mathcal{M} we have,

$$\mathbf{z}_j^T \mathcal{M} \Delta \mathbf{y} = \sum_{i: M_{ii}=1} \tilde{z}_{j,\mathcal{M}}[i] \Delta y[i] + \sum_{i: M_{ii}=1} \mathbb{E}[z_j[i] \mid M_{ii} = 1] \Delta y[i]. \quad (137)$$

We will use Hanson-Wright to bound the first term, which we previously expressed in (109) as:

$$\sum_{i: M_{ii}=1} \tilde{z}_{j,\mathcal{M}}[i] \Delta y[i] \quad (138)$$

$$= \frac{1}{4} \left((\mathbf{z}_{j,\mathcal{M}}^0 + \Delta \mathbf{y}_{\mathcal{M}}^0)^T \mathbf{I}_{\text{Tr}(\mathcal{M})} (\mathbf{z}_{j,\mathcal{M}}^0 + \Delta \mathbf{y}_{\mathcal{M}}^0) - (\mathbf{z}_{j,\mathcal{M}}^0 - \Delta \mathbf{y}_{\mathcal{M}}^0)^T \mathbf{I}_{\text{Tr}(\mathcal{M})} (\mathbf{z}_{j,\mathcal{M}}^0 - \Delta \mathbf{y}_{\mathcal{M}}^0) \right). \quad (139)$$

By Proposition D.2, the sub-Gaussian conditions for the entries of $\mathbf{z}_{j,m}^0$ are satisfied. Further, $\Delta \mathbf{y}_{\mathcal{M}}^0$ is bounded in $[-1, 1]$, so $k\Delta \mathbf{y}_{\mathcal{M}}^0k_{\psi_2} \leq 2$. Thus, by the triangle inequality, the sub-Gaussian norm of the entries of $\mathbf{z}_{j,\mathcal{M}}^0 - \Delta \mathbf{y}_{\mathcal{M}}^0$ is bounded by $K = 6 + 2 = 8$. Also note that conditioned on the

realization \mathcal{M} , $\mathbf{z}_{j,\mathcal{M}}^\theta$ is zero-mean by construction and $\Delta y_{\mathcal{M}}^\theta$ is zero-mean by symmetry between α and β , so we can now apply the Hanson-Wright inequality to both terms.

We choose parameter

$$t = \frac{K^2}{c} \sqrt{\text{Tr}(\mathcal{M})} \sqrt{\log(ndk)}. \quad (140)$$

where c is the constant from the Hanson-Wright result.

So

$$\frac{t^2}{K^4 k \mathbf{1}_{\text{Tr}(\mathcal{M})} k_{\text{F}}^2} = \frac{1}{c} \log(ndk) \quad (141)$$

$$\frac{t}{K^2 k \mathbf{1}_{\text{Tr}(\mathcal{M})} k_{\text{op}}} = \frac{1}{c} \sqrt{\text{Tr}(\mathcal{M})} \sqrt{\log(ndk)} > \frac{1}{c} \log(ndk). \quad (142)$$

The last inequality follows since with high probability $\text{Tr}(\mathcal{M}) = \Theta(\sqrt{n/k})$, by Proposition D.1, $\sqrt{\text{Tr}(\mathcal{M})} \sqrt{\log(ndk)} = \Theta(\sqrt{n \log(ndk)/k})$ grows faster than $\log(ndk)$.

Finally, note that:

$$\mathbb{E}[(\mathbf{z}_{j,\mathcal{M}}^\theta)^T \Delta y_{\mathcal{M}}^\theta \mathbf{1}_{\mathbf{M} = \mathcal{M}}] = \sum_{i:\mathcal{M}_{ii}=1} \mathbb{E}[\tilde{z}_{j,\mathcal{M}}[i] \Delta y[i] \mathbf{1}_{\mathbf{M} = \mathcal{M}}] \quad (143)$$

$$= \sum_{i:\mathcal{M}_{ii}=1} \mathbb{E}[\tilde{z}_{j,\mathcal{M}}[i] \Delta y[i] \mathbf{1}_{M_{ii} = 1}] \quad (144)$$

$$= \sum_{i:\mathcal{M}_{ii}=1} \frac{1}{2} \mathbb{E}[\tilde{z}_{j,\mathcal{M}}[i] \mathbf{1}_{\Delta y[i] = 1}] - \frac{1}{2} \mathbb{E}[\tilde{z}_{j,\mathcal{M}}[i] \mathbf{1}_{\Delta y[i] = -1}] \quad (145)$$

$$= 0, \quad (146)$$

where the last equation follows by symmetry. Knowing which of $\mathbf{z}_\alpha[i]$ or $\mathbf{z}_\beta[i]$ was the maximum does not change the conditional expectation of $\tilde{z}_{j,\mathcal{M}}[i]$.

So, applying Hanson-Wright, with probability at least $(1 - 4/(ndk))$ we have

$$\frac{K^2}{2} c_5 \sqrt{\text{Tr}(\mathcal{M})} \sqrt{\log(ndk)} = \frac{t}{2} \|\tilde{\mathbf{z}}_{j,m}^T \Delta \mathbf{y}\| \leq \frac{t}{2} = \frac{K^2}{2} c_5 \sqrt{\text{Tr}(\mathcal{M})} \sqrt{\log(ndk)}, \quad (147)$$

where $c_5 = \frac{1}{c}$.

We next consider the second term $\sum_{i:\mathcal{M}_{ii}=1} \mathbb{E}[z_j[i] \mathbf{1}_{M_{ii} = 1}] \Delta y[i]$ from (106) conditioned on the realization \mathcal{M} . By an identical symmetry argument as for the previous term we have, $0 = \mathbb{E}[z_j[i] \mathbf{1}_{E_j^c}] = \mathbb{E}[z_j[i] \mathbf{1}_{M_{ii} = 1}]$. Then as a consequence of Lemma C.3 and using the fact that $M_{ii} = 1$ implies $z_j[i]$ is not the maximum of k Gaussians we have, $\mathbb{E}[z_j[i] \mathbf{1}_{E_j^c}] \leq 2\sqrt{\log(k)}/(k-1)$. So we can bound

$$\left| \sum_{i:\mathcal{M}_{ii}=1} \mathbb{E}[z_j[i] \mathbf{1}_{M_{ii} = 1}] \Delta y[i] \right| \leq \frac{2\sqrt{\log(k)}}{k-1} \left| \sum_{i:\mathcal{M}_{ii}=1} \Delta y_i \right| \leq \frac{2\delta^\theta \sqrt{\log(k)}}{k-1}, \quad (148)$$

with probability $1 - 2e^{-\delta^2/(6 \text{Tr}(\mathcal{M}))}$, by application of the Chernoff bound and using the fact that conditioned on $M_{ii} = 1$, $\Delta y[i]$ takes value ± 1 with probability half by symmetry among features α and β .

Next, we apply the high probability bounds above on typical realizations \mathcal{M} . In particular, we substitute bounds on $\text{Tr}(\mathbf{M})$ from (102) from Proposition D.1 with $\delta = 1/2$ into (148), and set $\delta^\theta = \sqrt{6(1+\delta)(n/k) \log(ndk)}$. Then $e^{-\delta^2/(6 \text{Tr}(\mathbf{M}))} \leq 1/(ndk)$ and $e^{-\frac{2n\delta^2}{3k}} < 1/(ndk)$, so

using the union bound we have with probability at least $(1 - 4/(ndk) - 1/(ndk) - 1/(ndk))$,

$$|\mathbf{z}_j^T \Delta \mathbf{y}| \leq \left| \sum_{i: \mathcal{M}_{ii}=1} \tilde{z}_{j, \mathcal{M}}[i] \Delta y[i] \right| + \left| \sum_{i: \mathcal{M}_{ii}=1} \mathbb{E}[z_j[i] | j, M_{ii}=1] \Delta y[i] \right| \quad (149)$$

$$\frac{K^2}{2} c_5 \frac{\rho}{1+\delta} \sqrt{\frac{2n}{k}} \sqrt{\log(ndk)} + \frac{2\sqrt{(1+\delta)(n/k) \log(ndk)} \sqrt{\log(k)}}{k-1} \quad (150)$$

$$\frac{K^2}{2} c_5 \frac{\rho}{1+\delta} \sqrt{\frac{2n}{k}} \sqrt{\log(ndk)} + \frac{2\sqrt{(1+\delta)} \sqrt{\log(k)}}{k-1} \sqrt{\frac{n}{k}} \sqrt{\log(ndk)} \quad (151)$$

$$c_6 \sqrt{\frac{n}{k}} \sqrt{\log(ndk)}, \quad (152)$$

for a suitable choice of c_6 . \square

D.1.2 Bounds on $\mathbf{z}_j^> \mathbf{A}^{-1}(\mathbf{y}_\alpha - \mathbf{y}_\beta)$

We can now prove bounds on our target quantity. We restate the lemma that we are trying to prove below for convenience.

Lemma D.1. *Let $\Delta \mathbf{y} = \mathbf{y}_\alpha - \mathbf{y}_\beta$. Let α, β , and j be distinct. Then, with probability at least $(1 - 7/(ndk))$, we have,*

$$|\mathbf{z}_j^> \mathbf{A}^{-1} \Delta \mathbf{y}| \leq c_7 (\bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(ndk)} + \Delta_\mu n / \sqrt{k}), \quad (96)$$

for some constant c_7 .

Proof. We can rewrite

$$\mathbf{z}_j^> \mathbf{A}^{-1} \Delta \mathbf{y} = \mathbf{z}_j^> (\bar{\mu} \mathbf{I}_n + \Delta A_{inv}) \Delta \mathbf{y} \quad (153)$$

$$= \bar{\mu} \mathbf{z}_j^> \Delta \mathbf{y} + \mathbf{z}_j^> \Delta A_{inv} \Delta \mathbf{y}. \quad (154)$$

Next we can bound $|\mathbf{z}_j^> \Delta A_{inv} \Delta \mathbf{y}|$ simply as

$$|\mathbf{z}_j^> \Delta A_{inv} \Delta \mathbf{y}| \leq \|\mathbf{z}_j^>\|_2 \|\Delta A_{inv}\|_2 \|\Delta \mathbf{y}\|_2 \quad (155)$$

$$= k \Delta A_{inv} k_{op} k_{\mathbf{z}_j} k_{\Delta \mathbf{y}} \quad (156)$$

$$\Delta_\mu k_{\mathbf{z}_j} k_{\Delta \mathbf{y}}, \quad (157)$$

where we use the fact that ΔA_{inv} is a symmetric matrix and its 2-norm is its maximum absolute eigenvalue. We obtain the eigenvalue bounds for ΔA_{inv} from Lemma B.1, holding with probability at least $1 - 2e^{-n}$.

So, by the triangle inequality, we have with probability at least $(1 - 6/(ndk) - 2e^{-n} - 2e^{-\frac{2n\delta^2}{3k}} - 2e^{-n\delta^2})$

$$|\mathbf{z}_j^> \mathbf{A}^{-1} \Delta \mathbf{y}| \leq c_6 \bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(ndk)} + \Delta_\mu \sqrt{(1+\delta)n} \sqrt{(1+\delta) \frac{2n}{k}}. \quad (158)$$

The first term follows from Proposition D.3, and the second from our bound on $\text{Tr}(\mathbf{M}) = k \Delta \mathbf{y} k_2^2$ from Proposition D.1, as well as an analogous application of the chi-squared bound (Lemma C.2) on $k_{\mathbf{z}_j} k_2$.

The proof follows by setting δ to any value in $(0, 1)$, choosing an appropriate constant c_7 , and noting that for large enough n , $1/(ndk) - d_1 e^{-d_2 n/k}$ for any positive constants d_1, d_2 . \square

D.2 Proof of Lemma D.2

Next we use a similar technique as in Appendix D.1 to bound $\mathbf{z}_\alpha^> \mathbf{A}^{-1} \Delta \mathbf{y}$. We will write $\mathbf{A}^{-1} = \bar{\mu} \mathbf{I}_n + \Delta A_{inv}$, and split up the expression $\mathbf{z}_\alpha^> \mathbf{A}^{-1} \Delta \mathbf{y}$ into components involving $\bar{\mu} \mathbf{I}_n$, and components involving ΔA_{inv} .

Proposition D.4. Consider two arbitrary length- n zero-mean vectors \mathbf{y} and \mathbf{z} whose components each has sub-Gaussian norm at most K . With probability at least $1 - 4/(nk)$ we have each of

$$\mathbf{z} \succ \mathbf{y} \quad \mathbb{E}[\mathbf{z} \succ \mathbf{y}] + 2c_8 \frac{\rho_-}{n} \sqrt{\ln(nk)} \quad (159)$$

$$\mathbf{z} \succ \mathbf{y} \quad \mathbb{E}[\mathbf{z} \succ \mathbf{y}] - 2c_8 \frac{\rho_-}{n} \sqrt{\ln(nk)}, \quad (160)$$

for some universal constant c_8 .

Proof. The upper-bound follows as

$$\mathbf{z} \succ \mathbf{y} = \frac{1}{4} ((\mathbf{z} + \mathbf{y}) \succ (\mathbf{z} + \mathbf{y}) - (\mathbf{z} - \mathbf{y}) \succ (\mathbf{z} - \mathbf{y})) \quad (161)$$

$$\mathbb{E}[\mathbf{z} \succ \mathbf{y}] + \frac{K^2}{2} \frac{\rho_-}{c} \frac{\rho_-}{n} \frac{\rho_-}{\ln nk}, \quad (162)$$

with probability at least $(1 - 4/(nk))$, where we apply the Hanson-Wright inequality to each of the quadratic terms with $t = \frac{K^2}{c} \frac{\rho_-}{n} \sqrt{\ln(nk)}$ and use the fact that, letting $\mathbf{M} = \mathbf{I}_n$, $k\mathbf{M}k_F^2 = n$, $k\mathbf{M}k_{op} = 1$. The lower-bound can be obtained analogously, and an appropriate choice of c_8 completes the proof. \square

From this, we can now prove Lemma D.2, restated below for convenience:

Lemma D.2. Let $\Delta \mathbf{y} = \mathbf{y}_\alpha - \mathbf{y}_\beta$. With probability at least $(1 - 5/(nk))$, we have each of

$$\mathbf{z}_\alpha \succ \mathbf{A}^{-1} \Delta \mathbf{y} \quad \bar{\mu}(\mathbb{E}[\mathbf{z}_\alpha \succ \mathbf{y}_\alpha] - \mathbb{E}[\mathbf{z}_\alpha \succ \mathbf{y}_\beta]) + c_9(\bar{\mu} \frac{\rho_-}{n} \sqrt{\ln(nk)} + \Delta_\mu \frac{n}{\sqrt{k}}) \quad (97)$$

$$\mathbf{z}_\alpha \succ \mathbf{A}^{-1} \Delta \mathbf{y} \quad \bar{\mu}(\mathbb{E}[\mathbf{z}_\alpha \succ \mathbf{y}_\alpha] - \mathbb{E}[\mathbf{z}_\alpha \succ \mathbf{y}_\beta]) - c_9(\bar{\mu} \frac{\rho_-}{n} \sqrt{\ln(nk)} + \Delta_\mu \frac{n}{\sqrt{k}}), \quad (98)$$

for some constant c_9 .

Proof. We have

$$\mathbf{z}_\alpha \succ \mathbf{A}^{-1}(\mathbf{y}_\alpha - \mathbf{y}_\beta) = \mathbf{z}_\alpha \succ (\bar{\mu} \mathbf{I}_n + \Delta \mathbf{A}_{inv})(\mathbf{y}_\alpha - \mathbf{y}_\beta) \quad (163)$$

$$= \bar{\mu} \mathbf{z}_\alpha \succ (\mathbf{y}_\alpha - \mathbf{y}_\beta) + \mathbf{z}_\alpha \succ \Delta \mathbf{A}_{inv}(\mathbf{y}_\alpha - \mathbf{y}_\beta) \quad (164)$$

$$= \bar{\mu} \mathbf{z}_\alpha \succ (\mathbf{y}_\alpha - \mathbf{y}_\beta) + \mathbf{z}_\alpha \succ \Delta \mathbf{A}_{inv}(\mathbf{y}_\alpha^{oh} - \mathbf{y}_\beta^{oh}). \quad (165)$$

We again simply bound

$$|\mathbf{z}_\alpha \succ \Delta \mathbf{A}_{inv} \Delta \mathbf{y}| \leq k \mathbf{z}_\alpha k_2 k \Delta \mathbf{A}_{inv} \Delta \mathbf{y} k_2 \quad (166)$$

$$= k \Delta \mathbf{A}_{inv} k_{op} k \mathbf{z}_\alpha k_2 k \Delta \mathbf{y} k_2 \quad (167)$$

$$= \Delta_\mu k \mathbf{z}_\alpha k_2 k \Delta \mathbf{y} k_2 \quad (168)$$

$$\Delta_\mu \sqrt{(1 + \delta)n} \sqrt{(1 + \delta) \frac{2n}{k}} \quad (169)$$

$$= \Delta_\mu (1 + \delta) \frac{\rho_-}{2} \frac{n}{k}, \quad (170)$$

with probability $(1 - 2e^{-n\delta^2} - 2e^{-\frac{2n\delta^2}{3k}})$, using chi-squared bounds for \mathbf{z}_α (Lemma C.2) and Chernoff bounds for $\Delta \mathbf{y}$ (Proposition D.1).

With probability $(1 - 2e^{-n\delta^2} - 2e^{-\frac{2n\delta^2}{3k}})$, we get each of

$$\mathbf{z}_\alpha^T \mathbf{A}^{-1}(\mathbf{y}_\alpha - \mathbf{y}_\beta) \leq \bar{\mu} \mathbf{z}_\alpha^T (\mathbf{y}_\alpha - \mathbf{y}_\beta) + \Delta_\mu (1 + \delta) \frac{\rho_-}{2} \frac{n}{k} \quad (171)$$

$$\mathbf{z}_\alpha^T \mathbf{A}^{-1}(\mathbf{y}_\alpha - \mathbf{y}_\beta) \geq \bar{\mu} \mathbf{z}_\alpha^T (\mathbf{y}_\alpha - \mathbf{y}_\beta) - \Delta_\mu (1 + \delta) \frac{\rho_-}{2} \frac{n}{k}. \quad (172)$$

By applying Proposition D.4 on the relevant terms, setting δ to be an arbitrary value in $(0, 1)$, and choosing an appropriate constant c_9 , we obtain with probability $(1 - 5/(nk))$ each of

$$\mathbf{z}_\alpha^T \mathbf{A}^{-1}(\mathbf{y}_\alpha - \mathbf{y}_\beta) \leq \bar{\mu}(\mathbb{E}[\mathbf{z}_\alpha^T \mathbf{y}_\alpha] - \mathbb{E}[\mathbf{z}_\alpha^T \mathbf{y}_\beta]) + c_9(\bar{\mu} \frac{\rho_-}{n} \sqrt{\ln(nk)} + \Delta_\mu \frac{n}{\sqrt{k}}) \quad (173)$$

$$\mathbf{z}_\alpha^T \mathbf{A}^{-1}(\mathbf{y}_\alpha - \mathbf{y}_\beta) \geq \bar{\mu}(\mathbb{E}[\mathbf{z}_\alpha^T \mathbf{y}_\alpha] - \mathbb{E}[\mathbf{z}_\alpha^T \mathbf{y}_\beta]) - c_9(\bar{\mu} \frac{\rho_-}{n} \sqrt{\ln(nk)} + \Delta_\mu \frac{n}{\sqrt{k}}). \quad (174)$$

The probability comes from the union bound $(1 - 2e^{-n\delta^2} - 2e^{-\frac{2n\delta^2}{3k}} - 4/(nk)) \geq 1 - 5/(nk)$ (for sufficiently large n). \square

D.3 Proof of Corollary D.2.1

We claim the following bound:

Proposition D.5. *Bounds on $\mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\alpha]$.*

$$\rho \frac{1}{\pi \ln 2} \frac{n}{k} \rho \frac{1}{\ln k} \mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\alpha] \leq \frac{\rho}{2} \frac{n}{k} \rho \frac{1}{\ln k} \quad (175)$$

Proof.

$$\mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\alpha] = \mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\alpha^{oh}] \quad \mathbb{E}[\mathbf{z}_\alpha^> \frac{1}{c} \mathbf{1}] \quad (176)$$

$$= \mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\alpha^{oh}] \quad (177)$$

$$= n (\mathbb{E}[z_{\alpha,i} y_{\alpha,i}^{oh} j y_{\alpha,i}^{oh} = 1] P(y_{\alpha,i}^{oh} = 1) + \mathbb{E}[z_{\alpha,i} y_{\alpha,i}^{oh} j y_{\alpha,i}^{oh} = 0] P(y_{\alpha,i}^{oh} = 0)) \quad (178)$$

$$= \frac{n}{k} \mathbb{E}[z_{\alpha,i} j y_{\alpha,i}^{oh} = 1]. \quad (179)$$

So the desired bound follows from the bounds in Lemma C.3. \square

We can obtain a similar bound for when $\beta \notin \alpha$:

Proposition D.6. *Bounds on $\mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\beta]$.*

$$\frac{\rho}{2} \frac{n}{k} \frac{1}{k-1} \rho \frac{1}{\ln k} \mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\beta] \leq \frac{\rho}{\pi \ln 2} \frac{n}{k} \frac{1}{k-1} \rho \frac{1}{\ln k} \quad (180)$$

Proof. Observe that,

$$\mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\beta] = \mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\beta^{oh}] \quad \mathbb{E}[\mathbf{z}_\alpha^> \frac{1}{k} \mathbf{1}] \quad (181)$$

$$= \mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\beta^{oh}] \quad \frac{1}{k} \mathbb{E}[\mathbf{z}_\alpha] > \mathbf{1} \quad (182)$$

$$= \mathbb{E}[\mathbf{z}_\alpha^> \mathbf{y}_\beta^{oh}] \quad (183)$$

$$= \sum_i \mathbb{E}[z_{\alpha,i} y_{\beta,i}^{oh}] \quad (184)$$

$$= n (\mathbb{E}[z_{\alpha,i} y_{\beta,i}^{oh} j y_{\beta,i}^{oh} = 1] P(y_{\beta,i}^{oh} = 1) + \mathbb{E}[z_{\alpha,i} y_{\beta,i}^{oh} j y_{\beta,i}^{oh} = 0] P(y_{\beta,i}^{oh} = 1)) \quad (185)$$

$$= \frac{n}{k} \mathbb{E}[z_{\alpha,i} j y_{\beta,i}^{oh} = 1] \quad (186)$$

Now, observe that

$$\mathbb{E}[z_{\alpha,i} j y_{\alpha,i}^{oh} = 0] = \sum_{\beta \notin \alpha} \mathbb{E}[z_{\alpha,i} j y_{\beta,i}^{oh} = 1] \Pr(y_{\beta,i} = 1 \mid y_{\alpha,i}^{oh} = 0) \quad (187)$$

$$= \frac{1}{k-1} \sum_{\beta \notin \alpha} \mathbb{E}[z_{\alpha,i} j y_{\beta,i}^{oh} = 1] \quad (188)$$

$$= \mathbb{E}[z_{\alpha,i} j y_{\beta,i}^{oh} = 1] \quad (189)$$

for a particular β , by symmetry over the possible β .

Next we bound $\mathbb{E}[z_{\alpha,i} j y_{\alpha,i}^{oh} = 0]$ as follows:

$$\begin{aligned} & \mathbb{E}[z_{\alpha,i} j y_{\alpha,i}^{oh} = 1] P(y_{\alpha,i}^{oh} = 1) \\ + & \mathbb{E}[z_{\alpha,i} j y_{\alpha,i}^{oh} = 0] P(y_{\alpha,i}^{oh} = 0) = \mathbb{E}[z_{\alpha,i}] = 0 \end{aligned} \quad (190)$$

$$\Rightarrow \mathbb{E}[z_{\alpha,i} j y_{\alpha,i}^{oh} = 0] \frac{k-1}{k} = \mathbb{E}[z_{\alpha,i} j y_{\alpha,i}^{oh} = 1] \frac{1}{k} \quad (191)$$

$$\Rightarrow \mathbb{E}[z_{\alpha,i} j y_{\alpha,i}^{oh} = 0] = \mathbb{E}[z_{\alpha,i} j y_{\alpha,i}^{oh} = 1] \frac{1}{k-1} \quad (192)$$

Thus, substituting in the results from Lemma C.3, and plugging back into (186), we obtain

$$\frac{\rho_{-}}{2} \frac{n}{k} \frac{1}{k-1} \frac{\rho_{-}}{\ln k} \mathbb{E}[\mathbf{z}_{\alpha}^{\top} \mathbf{y}_{\beta}] = \frac{1}{\pi \ln 2} \frac{n}{k} \frac{1}{k-1} \frac{\rho_{-}}{\ln k}, \quad (193)$$

the desired result. \square

We can now prove Corollary D.2.1, which we restate below for convenience:

Corollary D.2.1. *Let $\Delta \mathbf{y} = \mathbf{y}_{\alpha} - \mathbf{y}_{\beta}$. With probability at least $(1 - 5/(nk))$, we have,*

$$\mathbf{z}_{\alpha}^{\top} \mathbf{A}^{-1} \Delta \mathbf{y} \leq c_{10} \bar{\mu} \frac{n}{k} \sqrt{\ln(k)} + c_9 (\bar{\mu} \frac{\rho_{-}}{n} \sqrt{\ln(nk)} + \Delta_{\mu} n / \sqrt{k}), \quad (99)$$

for some constant c_{10} .

Proof. This follows by substituting the lower bound from (175) in Proposition D.5 and the upper bound from (180) in Proposition D.6 into (98) from Lemma D.2, making an appropriate choice for c_{10} . \square

E Misclassification Events: Proof of Lemmas used in Theorem 5.1

With the previous section's utility bounds that allow us to deal with multiclass training data in hand, we are in a position to establish all the lemmas that we need to analyze misclassification.

E.1 Proof of Lemma B.2: Lower bound on $\min_{\beta} (X_{\alpha} - X_{\beta})$

With these bounds in hand, we can look at each misclassification event in turn. The first event to consider is if the best competing feature is unusually close to the true (maximum) feature.

Lemma B.2. *(Lower bound on the closest feature margin as $k \rightarrow \infty$): For any constant $\varepsilon > 0$, there exists a constant θ such that, for sufficiently large k with probability at least $(1 - \varepsilon)$,*

$$\min_{\beta:1} \min_{\beta \neq \alpha} (x_{test}[\alpha] - x_{test}[\beta]) \geq \frac{\theta}{\sqrt{2 \ln(k)}}. \quad (45)$$

Here, α is fixed and corresponds to the index of the true class — i.e. α corresponds to the index of the maximum feature among the first k features.

Proof. The following result from [34] whose proof we reproduce here¹⁷, enables us to bound the closest feature margin as:

$$\Pr(\min_{\beta} (x_{test}[\alpha] - x_{test}[\beta]) > \theta / \sqrt{2 \ln(k)}) \leq c_{11} e^{-\theta}, \quad (194)$$

for some universal positive constant c_{11} , for sufficiently large k . Thus, by selecting a constant θ such that $c_{11} e^{-\theta} = 1 - \varepsilon$ and choosing a sufficiently large k , we have that with probability $(1 - \varepsilon)$:

$$\min_{\beta} (x_{test}[\alpha] - x_{test}[\beta]) \geq \theta / \sqrt{2 \ln k}. \quad (195)$$

The proof [34] is reproduced below, with slight adaptations to match our use-case: Let β be the index of the largest competing feature to $x_{test}[\alpha]$. Then, their joint PDF becomes

$$f(x_{test}[\beta], x_{test}[\alpha]) = k(k-1)F(x_{test}[\beta])^{k-2} f(x_{test}[\beta])f(x_{test}[\alpha])\mathbf{1}(x_{test}[\beta] < x_{test}[\alpha]). \quad (196)$$

where F and f are the CDF and PDF of the standard Gaussian. Let

$$x = \frac{\theta}{\sqrt{2 \ln(k)}}. \quad (197)$$

Thus,

$$\Pr(x_{test}[\alpha] - x_{test}[\beta] > x) = k(k-1)J, \quad (198)$$

¹⁷We do this for the convenience of the reviewers since the source we are citing is a URL online. We believe that this is in the spirit of fair use.

where J is defined as

$$J = \int_0^1 \int_x^1 F(w)^{k-2} f(w) f(v+w) dv dw \quad (199)$$

$$= \int_0^1 F(w)^{k-2} f(w) \left(\int_x^1 f(v+w) dv \right) dw \quad (200)$$

$$= \int_0^1 F(w)^{k-2} f(w) (1 - F(x+w)) dw. \quad (201)$$

Substituting $u = F(w)$, we have that

$$J = \int_0^1 u^{k-2} (1 - F(x + F^{-1}(u))) du \quad (202)$$

$$= \int_0^{\ln(k)^2/k} u^{k-2} (1 - F(x + F^{-1}(u))) du + \int_{1/\ln(k)^2/k}^1 u^{k-2} (1 - F(x + F^{-1}(u))) du + \int_{1/(k \ln(k))}^1 u^{k-2} (1 - F(x + F^{-1}(u))) du, \quad (203)$$

splitting $[0, 1]$ into three intervals, and integrating separately over each one. Let the three integrals be J_1 , J_2 , and J_3 .

We have that

$$J_1 = \int_0^{\ln(k)^2/k} u^{k-2} (1 - F(x + F^{-1}(u))) du \quad (204)$$

$$\int_0^{\ln(k)^2/k} u^{k-2} du \quad (205)$$

$$(1 - \ln(k)^2/k)^{k-2} \quad (206)$$

$$\exp\left(-\frac{k-2}{k} \ln(k)^2\right) \quad (207)$$

$$= o\left(\frac{1}{k^2}\right). \quad (208)$$

Similarly,

$$J_3 = \int_{1/(k \ln(k))}^1 u^{k-2} (1 - F(x + F^{-1}(u))) du \quad (209)$$

$$= \int_{1/(k \ln(k))}^1 u^{k-2} (1 - F(F^{-1}(u))) du \quad (210)$$

$$\int_{1/(k \ln(k))}^1 u^{k-2} (1 - u) du \quad (211)$$

$$\frac{1}{k \ln(k)} \int_{1/(k \ln(k))}^1 u^{k-2} du \quad (212)$$

$$= o\left(\frac{1}{k^2}\right). \quad (213)$$

Finally, in the intermediate interval $u \geq [1 - \ln(k)^2/k, 1 - 1/(k \ln(k))]$, as $k \rightarrow \infty$, we see that $u \rightarrow (1 - \ln(k)^2/k) \rightarrow 1$, $x \rightarrow 0$, and $F^{-1}(u) \rightarrow \infty$, so for sufficiently large k ,

$$1 - F(x + F^{-1}(u)) \sim \frac{f(x + F^{-1}(u))}{x + F^{-1}(u)} \quad (214)$$

$$\sim \frac{f(x + F^{-1}(u))}{F^{-1}(u)} \quad (215)$$

$$\sim \frac{f(F^{-1}(u))}{F^{-1}(u)} e^{-x F^{-1}(u)} \quad (216)$$

$$\sim (1 - F(F^{-1}(u))) e^{-x F^{-1}(u)} \quad (217)$$

$$\sim (1 - u) e^{-x F^{-1}(u)}, \quad (218)$$

applying the well-known approximation for the Gaussian CCDF $1 - F(w) \sim f(w)/w$ for large w (for example, see Eqn. 8.2.38 from Gallager [27]), and substituting in the Gaussian PDF.

Further, since $F^{-1}(u) \rightarrow \infty$, we have that

$$1 - F(F^{-1}(u)) \sim \frac{f(F^{-1}(u))}{F^{-1}(u)} \quad (219)$$

$$\Rightarrow 1 - u \sim \frac{e^{-(F^{-1}(u))^2/2}}{F^{-1}(u) \sqrt{2\pi}} \quad (220)$$

$$\sim e^{-(F^{-1}(u))^2/(2+o(1))} \quad (221)$$

$$\Rightarrow F^{-1}(u) \sim \sqrt{2 \ln(1 - u)} \quad (222)$$

$$\sim \sqrt{2 \ln(k)}, \quad (223)$$

where the last step follows from the bounds on u in the intermediate interval.

Substituting the bounds from (218) and (223) into the expression for J_2 , and applying the definition of x from (197), we have,

$$J_2 \sim \int_{1 - \ln(k)^2/k}^{1 - 1/(k \ln(k))} u^{k-2} (1 - F(x + F^{-1}(u))) du \quad (224)$$

$$\sim \int_{1 - \ln(k)^2/k}^{1 - 1/(k \ln(k))} u^{k-2} (1 - u) e^{-x F^{-1}(u)} du \quad (225)$$

$$\sim \int_{1 - \ln(k)^2/k}^{1 - 1/(k \ln(k))} u^{k-2} (1 - u) e^{-(\theta/\sqrt{2 \ln(k)}) \sqrt{2 \ln(k)}} du \quad (226)$$

$$\sim \int_{1 - \ln(k)^2/k}^{1 - 1/(k \ln(k))} u^{k-2} (1 - u) e^{-\theta} du \quad (227)$$

$$\sim e^{-\theta} \left[\frac{u^{k-1}}{k-1} - \frac{u^{k-2}}{k-2} \right]_{1 - \ln(k)^2/k}^{1 - 1/(k \ln(k))} \quad (228)$$

$$\sim \frac{e^{-\theta}}{(k-1)(k-2)}. \quad (229)$$

Combining the terms from (208), (213), and (229), and substituting back into (198), we see that

$$\Pr(x_{test}[\alpha] - x_{test}[\beta] > x) = k(k-1)J \quad (230)$$

$$= k(k-1)(J_1 + J_2 + J_3) \quad (231)$$

$$\sim k(k-1) \left(\frac{e^{-\theta}}{(k-1)(k-2)} + o\left(\frac{1}{k^2}\right) \right) \quad (232)$$

$$\sim e^{-\theta}. \quad (233)$$

Expressing this as a non-asymptotic lower-bound on the probability, holding for sufficiently large k , yields the cited result in (194). \square

Lemma B.3. (Lower bound on the closest feature margin when k is constant): *If $k = c_k$ for some fixed constant c_k , for any constant $\varepsilon > 0$, there exists a constant $\varepsilon^\theta > 0$ such that*

$$\Pr \left(\min_{\beta, \gamma: 1 \leq \beta \neq \gamma \leq c_k} |x_{test}[\beta] - x_{test}[\gamma]| \geq \varepsilon^\theta \right) \geq 1 - \varepsilon. \quad (46)$$

Thus, with probability at least $(1 - \varepsilon)$,

$$\min_{\beta: 1 \leq \beta \neq \alpha \leq k} (x_{test}[\alpha] - x_{test}[\beta]) \geq \varepsilon^\theta. \quad (47)$$

Here, α is fixed and corresponds to the index of the true class — i.e. α corresponds to the index of the maximum feature among the first k features.

Proof. Observe that,

$$\min_{\beta \neq \alpha} (x_{test}[\alpha] - x_{test}[\beta]) \geq \min_{\beta \neq \gamma} |x_{test}[\beta] - x_{test}[\gamma]|. \quad (234)$$

In other words, rather than bounding the margin between the largest and second-largest features, we will lower-bound the absolute difference between any pair of features.

Consider a particular (β, γ) tuple. Observe that $x_{test}[\beta] - x_{test}[\gamma] \sim N(0, 2)$, since each feature is drawn independently from a standard Gaussian. For any $\varepsilon^\theta > 0$, we can upper-bound

$$\Pr (|x_{test}[\beta] - x_{test}[\gamma]| \geq \varepsilon^\theta) \leq \frac{2\varepsilon^\theta}{\pi} \quad (235)$$

by taking the product of the maximum value of the Gaussian pdf and the width, $2\varepsilon^\theta$, of the region we are interested in. Taking the union bound across all (β, γ) tuples, we find that

$$\Pr \left(\min_{\beta \neq \gamma} |x_{test}[\beta] - x_{test}[\gamma]| \geq \varepsilon^\theta \right) \leq \frac{c_k^2 \varepsilon^\theta}{\pi}. \quad (236)$$

So for any given $\varepsilon > 0$, we can choose $\varepsilon^\theta = \varepsilon \frac{\pi}{c_k^2}$, and have that

$$\Pr \left(\min_{\beta \neq \gamma} |x_{test}[\beta] - x_{test}[\gamma]| \geq \varepsilon \right) \geq 1 - \varepsilon. \quad (237)$$

\square

E.2 Lower bound on $\frac{\lambda \hat{h}_{\alpha, \beta}[\alpha]}{\max_{\beta} \text{CN}_{\alpha, \beta}}$

Next, we will find a lower bound for survival-contamination ratio within the regime with low survival variance.

Lemma B.4. (Lower bound on relative survival of true feature): *For any fixed $\beta \geq [k]$, $\beta \neq \alpha$, with $\lambda_\alpha = \lambda_\beta = \lambda$ we have with probability at least $(1 - 5/(nk))$,*

$$\lambda \hat{h}_{\alpha, \beta}[\alpha] \geq \lambda \left(c_{10} \bar{\mu} \frac{n}{k} \sqrt{\ln(k)} - c_9 (\bar{\mu} \frac{D}{n} \sqrt{\ln(nk)} + \Delta_\mu \frac{D}{n/k}) \right), \quad (48)$$

for universal positive constants c_9 and c_{10} .

Proof. Using Corollary D.2.1, we lower bound $\hat{h}_{\alpha, \beta}[\alpha]$ with probability at least $(1 - 5/(nk))$ as

$$\hat{h}_{\alpha, \beta}[\alpha] = \lambda_\alpha^{1/2} (\hat{f}_\alpha[\alpha] - \hat{f}_\beta[\alpha]) \quad (238)$$

$$= \mathbf{z}_\alpha^\top \mathbf{A}^{-1} \mathbf{y}_\alpha - \mathbf{z}_\alpha^\top \mathbf{A}^{-1} \mathbf{y}_\beta \quad (239)$$

$$\geq c_{10} \bar{\mu} \frac{n}{k} \sqrt{\ln(k)} - c_9 (\bar{\mu} \frac{D}{n} \sqrt{\ln(nk)} + \Delta_\mu \frac{D}{n/k}). \quad (240)$$

Multiplying through by λ gives the desired result. \square

From the above result, under the scalings of our bi-level model we obtain:

Corollary B.4.1. *Under the bi-level ensemble model 4.2, for any fixed $\beta \geq [k]$, $\beta \notin \alpha$, $\lambda_\alpha = \lambda_\beta = \lambda$ if $t < 1/2$, $t < 2(q+r-1)$ and $1 < q+r < (p+1)/2$, with probability at least $(1-5/(nk))$,*

$$\lambda \widehat{h}_{\alpha,\beta}[\alpha] = c_{12} n^{1-q-r-t} \sqrt{\ln(k)}, \quad (49)$$

for universal positive constant c_{12} .

Proof. Substituting our asymptotic scalings into the results from Lemma B.4 and using the decay rate of $\bar{\mu} = n^{-p}$ from Corollary C.1.1 (which we can do since $1 < q+r < (p+1)/2$), we find that

$$\lambda \widehat{h}_{\alpha,\beta}[\alpha] = n^{p-q-r} \left(c_{10} \bar{\mu} \frac{n}{k} \sqrt{\ln(k)} - c_9 (\bar{\mu}^{-p} n^{-p} \sqrt{\ln(nk)} + \Delta_\mu n / \bar{k}) \right) \quad (241)$$

$$= c_{10} n^{1-q-r-t} \sqrt{\ln(k)} - c_9 n^{1/2-q-r} \sqrt{\ln(nk)} - c_9 n^{2-2q-2r-t/2} \quad (242)$$

$$= c_{12} n^{\max(1-q-r-t, 2-2q-2r-t/2)} \sqrt{\ln(k)} \quad (243)$$

$$= c_{12} n^{1-q-r-t+\max(0, 1-q-r-t/2)} \sqrt{\ln(k)} \quad (244)$$

$$= c_{12} n^{1-q-r-t} \sqrt{\ln(k)}, \quad (245)$$

for an appropriately chosen universal constant c_{12} and sufficiently large n . \square

Next we upper bound $\max_\beta \text{CN}_{\alpha,\beta}$.

Lemma B.5. (Upper bound on contamination): *For any fixed $\beta \geq [k]$, $\beta \notin \alpha$, with probability at least $(1-7/(nk))$,*

$$\text{CN}_{\alpha,\beta} \leq c_7 \left(\bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(ndk)} + \Delta_\mu n / \bar{k} \right) \sqrt{\sum \lambda_j^2}, \quad (50)$$

for universal positive constant c_7 .

Proof. For each β we have,

$$\text{CN}_{\alpha,\beta} = \sqrt{\left(\sum_{j \notin \alpha, \beta} \lambda_j^2 (\widehat{h}_{\beta,\alpha}[j])^2 \right)} \quad (246)$$

For $j \notin \alpha, \beta$, by Lemma D.1,

$$\left| \widehat{h}_{\beta,\alpha}[j] \right| = \left| \widehat{h}_{\alpha,\beta}[j] \right| \quad (247)$$

$$= \left| \widehat{f}_j - \widehat{g}_j \right| \quad (248)$$

$$= \left| \mathbf{z}_j^\top \mathbf{A}^{-1} \mathbf{y}_\alpha - \mathbf{z}_j^\top \mathbf{A}^{-1} \mathbf{y}_\beta \right| \quad (249)$$

$$= \left| \mathbf{z}_j^\top \mathbf{A}^{-1} (\mathbf{y}_\alpha - \mathbf{y}_\beta) \right| \quad (250)$$

$$\leq c_7 \left(\bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(ndk)} + \Delta_\mu n / \bar{k} \right), \quad (251)$$

with probability $1 - 7/(nk)$.

So taking the union bound over all $d \geq 2$ terms in the expression for the contamination, we can upper-bound it as

$$\text{CN}_{\alpha,\beta} \leq c_7 \left(\bar{\mu} \sqrt{\frac{n}{k}} \sqrt{\ln(ndk)} + \Delta_\mu n / \bar{k} \right) \sqrt{\sum \lambda_j^2}, \quad (252)$$

with probability $(1 - 7/(nk))$, the desired result. \square

Corollary B.5.1. *Under the bi-level model 4.2, in the regime $1 < q + r < (p + 1)/2$, with probability at least $(1 - 7/(nk))$,*

$$\text{CN}_{\alpha,\beta} \leq c_{13} n^{(1-t-p)/2 + \max(0, 3/2 - q - r) + \max(0, p/2 - q - r/2)} \sqrt{\ln(ndk)}, \quad (51)$$

for universal positive constant c_{13} .

Proof. Since $1 < q + r < (p + 1)/2$, we can apply Corollary C.1.1 to the result from Lemma B.5 and substitute in the known scalings of various terms, to obtain

$$\text{CN}_{\alpha,\beta} \leq c_7 (n^{1/2 - t/2 - p} \sqrt{\ln(ndk)} + c_4 n^{2 - p - q - r - t/2}) (n^{p - q - r/2} + n^{p/2}) \quad (253)$$

$$c_{13} n^{(1-t-p)/2 + \max(0, 3/2 - q - r) + \max(0, p/2 - q - r/2)} \sqrt{\ln(ndk)}, \quad (254)$$

for an appropriately chosen universal positive constant c_{13} . □

E.3 Proof of Lemma B.6: Bounds on Survival Variance

Finally, we look at the error event where a competing feature has unusually high survival relative to the true feature, so it is incorrectly selected.

Lemma B.6. (Upper bound on survival variance): *For any fixed competing feature $\beta \in [k]$, $\beta \neq \alpha$ with $\lambda_\alpha = \lambda_\beta$, we have with probability at least $(1 - 15/(nk))$,*

$$\frac{\widehat{h}_{\alpha,\beta}[\alpha] - \widehat{h}_{\beta,\alpha}[\beta]}{\widehat{h}_{\alpha,\beta}[\alpha]} \leq \frac{2c_9 (\bar{\mu}^{\rho_{\bar{n}}} \sqrt{\ln(nk)} + \Delta_\mu n / \bar{k})}{c_{10} \bar{\mu}^{\frac{\rho_{\bar{n}}}{k}} \sqrt{\ln(k)} - c_9 (\bar{\mu}^{\rho_{\bar{n}}} \sqrt{\ln(nk)} + \Delta_\mu n / \bar{k})}, \quad (52)$$

for universal positive constants c_9 and c_{10} .

Proof. We first consider the numerator of the LHS of (52). By Lemma D.2, with probability at least $(1 - 5/(nk))$,

$$\widehat{h}_{\alpha,\beta}[\alpha] = \lambda_\alpha^{1/2} (\widehat{f}_\alpha[\alpha] - \widehat{f}_\beta[\alpha]) \quad (255)$$

$$= \mathbf{z}_\alpha^\top \mathbf{A}^{-1} \mathbf{y}_\alpha - \mathbf{z}_\alpha^\top \mathbf{A}^{-1} \mathbf{y}_\beta \quad (256)$$

$$\bar{\mu} (E[\mathbf{z}_\alpha^\top \mathbf{y}_\alpha] - E[\mathbf{z}_\alpha^\top \mathbf{y}_\beta]) + c_9 (\bar{\mu}^{\rho_{\bar{n}}} \sqrt{\ln(nk)} + \Delta_\mu n / \bar{k}). \quad (257)$$

Similarly, with probability at least $(1 - 5/(nk))$,

$$\widehat{h}_{\beta,\alpha}[\beta] = \lambda_\beta^{1/2} (\widehat{f}_\beta[\beta] - \widehat{f}_\alpha[\alpha]) \quad (258)$$

$$= \mathbf{z}_\beta^\top \mathbf{A}^{-1} \mathbf{y}_\beta - \mathbf{z}_\beta^\top \mathbf{A}^{-1} \mathbf{y}_\alpha \quad (259)$$

$$\bar{\mu} (E[\mathbf{z}_\beta^\top \mathbf{y}_\beta] - E[\mathbf{z}_\beta^\top \mathbf{y}_\alpha]) - c_9 (\bar{\mu}^{\rho_{\bar{n}}} \sqrt{\ln(nk)} + \Delta_\mu n / \bar{k}). \quad (260)$$

By symmetry,

$$E[\mathbf{z}_\beta^\top \mathbf{y}_\beta] = E[\mathbf{z}_\alpha^\top \mathbf{y}_\alpha] \quad (261)$$

$$E[\mathbf{z}_\beta^\top \mathbf{y}_\alpha] = E[\mathbf{z}_\alpha^\top \mathbf{y}_\beta]. \quad (262)$$

Thus with probability at least $(1 - 10/(nk))$,

$$\widehat{h}_{\alpha,\beta}[\alpha] - \widehat{h}_{\beta,\alpha}[\beta] \leq 2c_9 (\bar{\mu}^{\rho_{\bar{n}}} \sqrt{\ln(nk)} + \Delta_\mu n / \bar{k}). \quad (263)$$

Using Corollary D.2.1 to lower-bound the denominator of the LHS of (52), we obtain with probability at least $(1 - 15/(nk))$

$$\frac{\widehat{h}_{\alpha,\beta}[\alpha] - \widehat{h}_{\beta,\alpha}[\beta]}{\widehat{h}_{\alpha,\beta}[\alpha]} \leq \frac{2c_9 (\bar{\mu}^{\rho_{\bar{n}}} \sqrt{\ln(nk)} + \Delta_\mu n / \bar{k})}{c_{10} \bar{\mu}^{\frac{\rho_{\bar{n}}}{k}} \sqrt{\ln(k)} - c_9 (\bar{\mu}^{\rho_{\bar{n}}} \sqrt{\ln(nk)} + \Delta_\mu n / \bar{k})}. \quad (264)$$

□

We can apply Corollary C.1.1 to simplify our results from Lemma B.6 in the asymptotic regime for the bi-level model.

Corollary B.6.1. *Under the bi-level ensemble model 4.2, for any fixed $\beta \geq 2/k$, $\beta \notin \alpha$, if $t < 1/2$, $t < 2(q+r-1)$, and $1 < q+r < (p+1)/2$, with probability at least $(1-15/(nk))$,*

$$\frac{\widehat{h}_{\alpha,\beta}[\alpha] - \widehat{h}_{\beta,\alpha}[\beta]}{\widehat{h}_{\alpha,\beta}[\alpha]} < n^{-u}, \quad (53)$$

for large enough n for some fixed $u > 0$.

Proof. Substituting, using Corollary C.1.1, in the regime where $1 < q+r < (p+1)/2$ and $t < 1/2$, we find that

$$\frac{\widehat{h}_{\alpha,\beta}[\alpha] - \widehat{h}_{\beta,\alpha}[\beta]}{\widehat{h}_{\alpha,\beta}[\alpha]} = \frac{2c_9(n^{1/2-p}\sqrt{\ln(nk)} + c_4n^{2-p-q-r-t/2})}{c_{10}n^{1-p-t}\sqrt{\ln(k)} - c_9(n^{1/2-p}\sqrt{\ln(nk)} + c_4n^{2-p-q-r-t/2})} \quad (265)$$

$$\frac{2c_9}{c_{10}} \frac{n^{1/2}\sqrt{\ln(nk)} + c_4n^{2-q-r-t/2}}{n^{1-t} - (c_9/c_{10})n^{1/2}\sqrt{\ln(n)} + (c_9 - c_4/c_{10})n^{2-q-r-t/2}} \quad (266)$$

$$c_{14} \frac{n^{1/2}\sqrt{\ln(nk)} + n^{2-q-r-t/2}}{n^{1-t}} \quad (267)$$

$$c_{14}n^{\max(t-1/2, t/2+1-q-r)}\sqrt{\ln(nk)}, \quad (268)$$

for sufficiently large n and an appropriate choice of positive constant c_{14} . Thus, if $\max(t-1/2, t/2+1-q-r) < 0$, our quantity of interest tends to zero at a polynomial rate as $n \rightarrow \infty$, completing the proof. \square

F Conjectured Looseness of Bound

In (155) in the proof of Lemma D.1, we upper bound $\mathbf{z}_j^\top \Delta A_{inv} \Delta y$ using the Cauchy-Schwarz inequality as

$$|\mathbf{z}_j^\top \Delta A_{inv} \Delta y| \leq \|\mathbf{z}_j\|_2 \|\Delta A_{inv} \Delta y\|_2 \quad (269)$$

$$\leq k \Delta A_{inv} k_{op} \|\mathbf{z}_j\|_2 \|\Delta y\|_2 \quad (270)$$

$$\leq \Delta_\mu \|\mathbf{z}_j\|_2 \|\Delta y\|_2. \quad (271)$$

This results in a high-probability bound of the order $\Delta_\mu n^{1/p} / \sqrt{k}$. Essentially this bound fears that ΔA_{inv} can, in worst case, align \mathbf{z}_j and Δy to be in the same direction. However, since there is only a weak dependence between ΔA_{inv} and \mathbf{z}_j and Δy this bound is likely overly cautious. We conjecture that this bound is loose by a factor \sqrt{n} . Why do we conjecture this? If we ignored the dependency of ΔA_{inv} on \mathbf{z}_j and Δy and blindly applied the Hanson-Wright inequality (with the \mathbf{M} matrix introduced as in Appendix D.1.1 to leverage the fact that Δy is mostly zeros) then we would obtain a high-probability upper bound of the form $\Delta_\mu \sqrt{n/k}$ (ignoring the logarithmic factors).

Assuming this tighter conjectured bound holds and similarly assuming an analogously tighter bound for $|\mathbf{z}_\alpha^\top \Delta A_{inv} \Delta y|$ in Appendix D.2 and following through with the rest of our analysis, we obtain the conjectured sufficient conditions for good generalization as in Equation (23) from Conjecture 6.1 for the regime $q+r > 1$.

It turns out that whenever the survival/contamination ratio grows at a polynomial rate n^v for $v > 0$ then the survival variation term also shrinks at a polynomial rate n^{-u} for $u > 0$. Thus ensuring the survival/contamination ratio is large enough (i.e. the number of classes is not too large relative to the level of favoring of potentially true features) is key to obtaining good generalization.

Although we focus on the regime $q+r > 1$ in our work, our proof technique is also applicable to the regime $q+r < 1$, i.e where regression works and by grinding through the math for this setting we should be able to get sufficient conditions for good generalization here as well. The survival in the multi-class setting in the regime $q+r < 1$ will scale roughly as $1/k$ due to the fewer positive training examples per class instead of behaving like the constant $\sqrt{2/\pi}$ as was the case for binary classification (Lemma 32, [56]). Moreover, Lemma 34 from Muthukumar et al. [56] shows that for

the binary classification setting the contamination scales as $n^{-\min(p-1, 1-r)/2}$ when $q+r=1$. In the multiclass setting the contamination will be lower by a factor of $\frac{1}{k}$ and substituting this in our error analysis we obtain Conjecture 6.1 for the regime $q+r < 1$.

Finally, we believe that we can adapt our analysis from the Proof of Theorem 5.1 in Appendix B to write a set of sufficient conditions for poor generalization. The primary condition for this would be for the relevant survival/contamination ratio to go to zero. We conjecture that computing conditions on p, q, r, t under which this occurs results in the converse result in the form of sufficient conditions for poor generalization present in Conjecture 6.1. Intuitively, if the survival/contamination ratio goes to zero, then the contamination can with significant probability flip the sign of a comparison involving the score that should be winning — this parallels the way that the converse is proved in Muthukumar et al. [56] for binary classification.

G Scaling parameters with the number of positive training examples per class

From our results in Figure 2 we observed that as the number of classes k increases (i.e. larger values of t), the region where multiclass classification generalizes well shrinks. A justification for this is when the number of classes k increases while the number of training points n stays constant, we have fewer positive training examples from each class, and this makes the task harder.

To see if the reduced number of positive training examples is indeed the dominant effect, we can explore what happens if we increase the number of total training points to compensate for this effect? Instead of scaling all parameters with the total number of training points, what happens if we scale them with the number of positive training examples per class?

Let $N = n^b$ be the new number of training points for some $b > 1$, while rest of the parameters in the bi-level model scale as before. We have,

$$N = n^b \quad (272)$$

$$d = n^p = N^{p/b} \quad (273)$$

$$s = n^r = N^{r/b} \quad (274)$$

$$a = n^{-q} = N^{-q/b} \quad (275)$$

$$k = c_k n^{-t} = c_k N^{-t/b}. \quad (276)$$

We can interpret this as our standard setup, albeit parameterized by N , rather than n . To keep the model well-defined we require the following:

- $b < p$, to ensure we are still overparameterized;
- $r < b$, to ensure the number of favored features does not exceed the total number of training points;
- $q < p - r$ to ensure we are actually favoring the first s features.

For this setup, Theorem 5.1 states that the probability of misclassification tends to zero if

$$\frac{t}{b} < \min\left(\frac{r}{b}, 1 - \frac{r}{b}, \frac{p}{b} + 1 - 2\left(\frac{q}{r} + \frac{r}{b}\right), \frac{p}{b} - 2, \frac{2q}{b} + \frac{r}{b} - 2\right) \quad (277)$$

$$\frac{q}{b} + \frac{r}{b} > 1. \quad (278)$$

Rearranging, we obtain the condition

$$t < \min(r, b - r, p + b - 2(q + r), p - 2b, 2q + r - 2b) \quad (279)$$

$$q + r > b. \quad (280)$$

To hold the number of training samples per class fixed we can set $b = t + 1$, so the ratio N/k becomes constant. Doing so, we obtain the following sufficient conditions for good generalization:

$$t < \min\left(r, \frac{p-2}{3}, \frac{2q+r-2}{3}\right) \quad (281)$$

$$0 < 1 - r \quad (282)$$

$$0 < p + 1 - 2(q + r) \quad (283)$$

$$t < q + r - 1. \quad (284)$$

Additionally for the model to be well defined we require $t < p - 1$. (The other conditions $r < t + 1$ and $q < p - r$ for model to be well defined are automatically satisfied if the above conditions for good generalization are satisfied).

If we assume Conjecture 6.1 then a set of sufficient conditions for good generalization is:

$$0 < p + 1 - 2(q + r) \quad (285)$$

$$r < 1 \quad (286)$$

$$t < r \quad (287)$$

$$t < p - 1. \quad (288)$$

The first two conditions must be satisfied for binary classification problem to generalize well and thus for multi-class classification to succeed in this setting we need to ensure binary classification succeeds. The condition $t < r$ arises because if we don't favor the features used in the comparison while assigning class labels then we have no hope of succeeding in overparameterized settings. The condition $t < p - 1$ ensures that the problem is overparameterized. If any of these conditions is not met then the probability of classification error will tend to 1.

Figure 4 visualizes the conjectured regimes for this alternative setup where the number of positive training examples per class is held fixed as we vary the number of classes for fixed values of p and q . In the white region, our model is not well defined. Note that in subfigure (a), the limiting factor to the model being well defined is the inequality $r < 1 + t$ (we must have more training examples than favored features) while in subfigure (b), the limiting factor for the model being well defined in the right-hand boundary is the inequality $r < p - q$ (we must put a larger weight on the features we favor as compared to those that we do not favor). In subfigure(b) we see that the top boundary for the model being well defined is the inequality $t < p - 1$ which is necessary for the problem to be overparameterized and support the existence of interpolating solutions. Further, the right-hand bound for good generalization in subfigure (a) corresponds to the inequality $r < 1$ while in subfigure (b) it corresponds to $p + 1 > 2(q + r)$. The left-hand boundary for good generalization in both figures is the inequality $t < r$, which reflects the fact that for MNI-based classification to succeed, all the features defining the classes must be favored.

It is interesting to note that when we add more training points so as to increase the number of positive examples, we are effectively decreasing the level of overparameterization in the problem. We know from Nakkiran [58] that adding training data in a way that reduces overparameterization can sometimes make performance worse instead of better. However, in the deeply overparameterized setting of the bi-level models explored here, this effect is counteracted by the survival benefits of having more positive examples — in effect, reducing the overall level of overparameterization reduces the shrinkage induced by the regularizing effect of overparameterization. This reduction in shrinkage compensates for the $\frac{1}{k}$ hit to survival induced by the larger number of classes.

H Additional related work

H.1 Comparisons to Wang et al. [77]

Recent work [77] provides an analysis of the generalization error of the minimum-norm interpolation of one-hot labels for multiclass classification with Gaussian features. Using the bi-level model, the authors present parameter regimes where multiclass classification error goes to zero asymptotically. While our work has many similarities in terms of model and problem setting, there are some key differences.

The first key difference is in how the training data is generated. In this paper, we assume the true label of a point is generated based on which of the first k dimensions is the largest, while Wang et al.

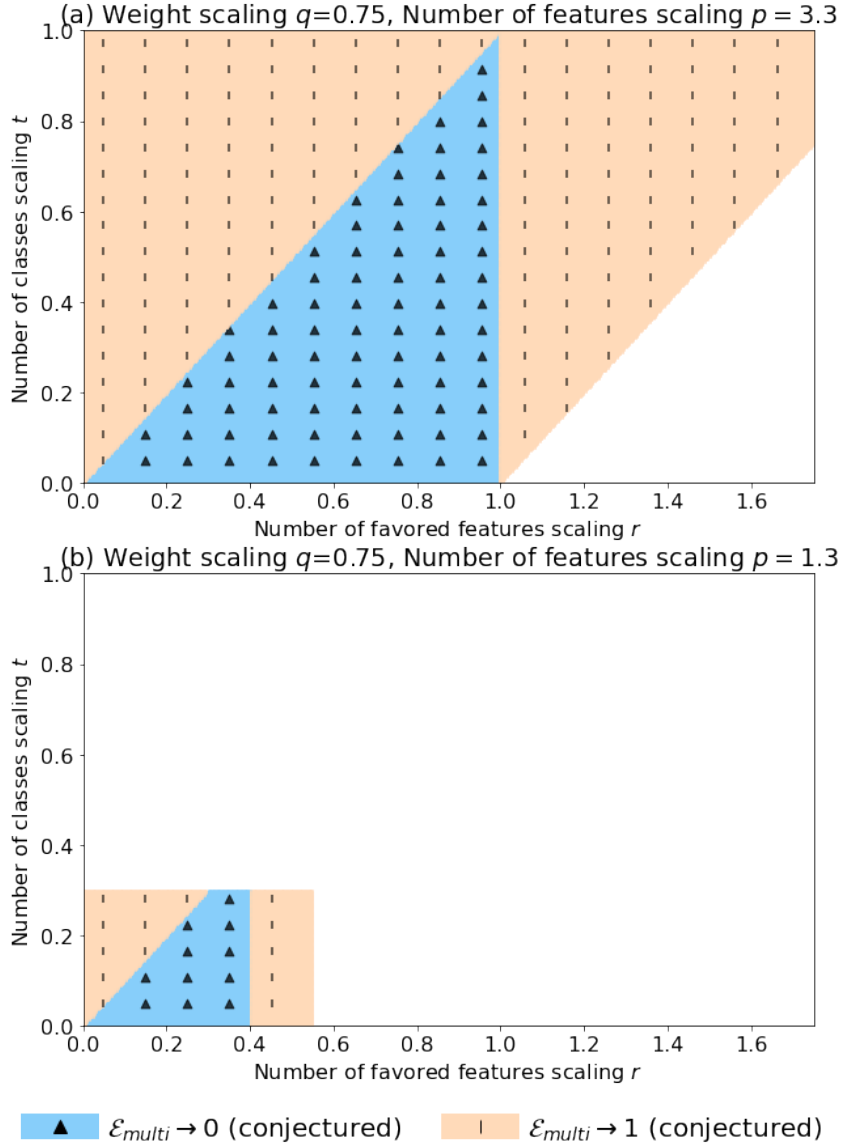


Figure 4: Visualization of the conjectured bi-level classification regimes when we scale everything with the number of positive training examples per class, instead of with the total number of training points.

[77] consider a Gaussian mixture model and a multinomial logistic model where the true labels have some randomness even conditioned on the first k dimensions. Like us, however, they also consider the case of orthogonal classes.

Second, we consider the asymptotic case where the number of classes, k scales with the number of training points as $k = cn^t$ for some positive integer c and non-negative real t . The work in Wang et al. [77] considers only the finite classes setting i.e. $t = 0$ in our model. The error analysis technique employed by us here in the form of a typicality-style argument featuring the feature margin (difference between the largest and second largest feature) is much tighter than the method employed in Wang et al. [77] and allows us to compute regimes where multiclass classification succeeds even when $t > 0$. A straight substitution into the analysis from Wang et al. [77] does not work since that analysis is too loose for this setting. Furthermore, in our expressions for survival and contamination

(Lemmas B.4 and B.5) we compute an exact dependence on k .¹⁸ The expressions from Wang et al. [77] don't compute this exact dependence because it is not required for their purposes. By using our novel analysis technique we are able to elucidate the challenges posed by fewer positive training examples per class in the multiclass setting and provide sufficient conditions for generalization when number of classes scales with the number of training points.

An equivalence between the solution obtained by minimum- ℓ_2 -norm interpolation on the adjusted zero-mean one-hot encoded labels that we perform in our approach (Equation (9)) and the solution obtained by other training methods has been established in [77]. In particular the minimum-norm interpolating solution is typically identical to the solution obtained via one-vs-all SVM and multi-class SVM (and thus gradient descent on cross-entropy loss due to its implicit bias [36, 68], under sufficient overparameterization. From Wang et al. [77], the sufficient conditions for the equivalence of solutions are,

$$\frac{\sum_{j=1}^n \lambda_j}{\lambda_1} > C_1 k^2 n \ln(kn), \quad (289)$$

$$\frac{(\sum_{j=1}^n \lambda_j)^2}{\sum_{j=1}^n \lambda_j^2} > C_2 (\ln(kn) + n), \quad (290)$$

where C_1, C_2 are positive constants. Under our bi-level model (Definition 4.2) these conditions translate to:

$$q + r > 2t + 1, \quad (291)$$

$$2p \max(2p - 2q - r, p) > 1. \quad (292)$$

This can be succinctly expressed as,

$$0 < t < \frac{q + r - 1}{2}. \quad (293)$$

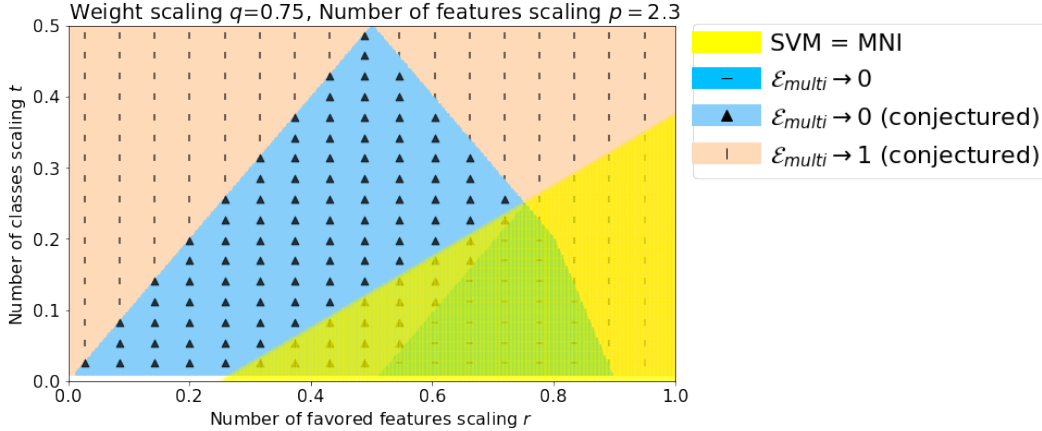


Figure 5: Visualization of regime where SVM solution is identical to MNI solution.

In Figure 5 we plot our provable as well as conjectured regimes alongside the regime where minimum-norm interpolation leads to same solution as other training methods. Notice the overlap. Thus our analysis is not limited only to the minimum-norm interpolator but holds for other training methods when the problem is sufficiently overparameterized. In this sense, the results in Wang et al. [77] and the present paper should be read together to tell a more full story of overparameterized multiclass classification.

¹⁸In particular, our analysis here brings out the fact that multiclass training data becomes less informative per training sample as the number of classes increases. This results in a $\frac{1}{k}$ scaling term in survival and a $\frac{1}{k^2}$ scaling in contamination. It is this effect that makes it possible in some regimes for the contamination from other favored features to dominate — whereas in the case of binary classification, it is always the contamination from unfavored features that dominates.

H.2 Comparisons to Muthukumar et al. [56]

The work in Muthukumar et al. [56] provides an analysis of the binary classification and regression problem with Gaussian features in the overparameterized regime and shows that binary classification is easier than regression by proving the existence of a regime in a bi-level model where binary classification generalizes well but regression does not. In this work we use a similar bi-level model and the signal-processing inspired concepts of survival and contamination in our proofs but the nature of the training data in the multiclass classification problem is the key challenge and complicates our analysis considerably. Since the true class labels are generated by comparing k features, we no longer have independence of the class label y with any of these features. This is relevant when we compute bounds on the term $z_\alpha^> \mathbf{A}^{-1}(y_\alpha - y_\beta)$ an integral part of our survival quantity (Equations (97),(98) from Lemma D.2), since the Hanson-Wright inequality is no longer applicable directly as was the case for the binary classification problem in prior work (Appendix D.3.1 of Muthukumar et al. [56]). Working through these challenges, we prove that the multiclass problem is fundamentally different from (and harder than) than the binary problem due to the effect of fewer informative samples (positive training examples) per class. In particular we show via dominant terms from Lemmas B.6 and B.5 that, as we increase the number of classes k , survival shrinks as $1/k$ while contamination shrinks only as $1/\sqrt{k}$. Thus the survival/contamination ratio which plays a key role in the expression for classification error decreases as $1/\sqrt{k}$ in the multiclass setting as we increase k . Thus, for good generalization we need to ensure number of classes is not too large in addition to having sufficient favoring of true features.

I Experimental results

Theorem 5.1 is proved rigorously and so we know that the asymptotic result is true. Underlying the result is the analysis of survival (how strongly is the true feature underlying this class represented in the learned score) and contamination (what is the standard deviation of the contamination in predictions that comes from learning nonzero coefficients to features that have nothing to do with this class). Multiclass classification asymptotically succeeds when the survival dominates the contamination.

In Fig. 6, we plot experimental results using the bi-level ensemble model (Definition 4.2) for a setting where regression does not work but multiclass classification is conjectured to work. We plot quantities from Equation (22) in our error analysis. From subfigures (a),(b) and (c) we observe that while both survival and contamination are decreasing as we increase n , the survival/contamination ratio increases. The survival/contamination ratio growing with n is important for correct classification. The trend is very clear from the experimental results and indicates that continuing to grow n (together with the number of classes k , the number of features d , the number of favored features s , and the level of favoring as per our bi-level idealized model) would result in ever improving performance. Furthermore, we see that the empirical slope of these quantities on a log-log plot (and thus the power-law scaling of these quantities with respect to n) agree with the theoretical slopes calculated based on our conjecture.

Subfigure (d) plots the binary classification error when only trying to distinguish between the true class and one other particular class. (In this experiment, the true class was determined by feature 1. We calculate the binary error as the probability of misclassifying a point from class 1 as belong to class 2 when we only compare the scores for class 1 and class 2.) We see that this error clearly decreases as we increase n . One way of thinking about successful multiclass classification is that the true class must win such pairwise competitions against all competing classes. Finally, subfigure (e) plots the total multiclass misclassification error overlaid with the number of classes. Here too we see a downward trend in classification error as n increases, even though we would have to go to significantly larger n than our compute could handle to see this error probability drop to very low values. Notice the integer effects arising from the number of classes k sometimes not growing with n that result in small upward spikes in the classification error.

J Comment on empirical eigenstructures of feature matrices

It is well known (for instance Remark 6 from Muthukumar et al. [56] that cites [78]) that for a spiked covariance model when the ratio of the top to the bottom eigenvalues grows as $\Omega(d/n)$, the top s eigenvalues can be estimated reliably from samples, even when the number of training samples n is less than the number of features d . The ratio of the top to the bottom eigenvalue in our bilevel model

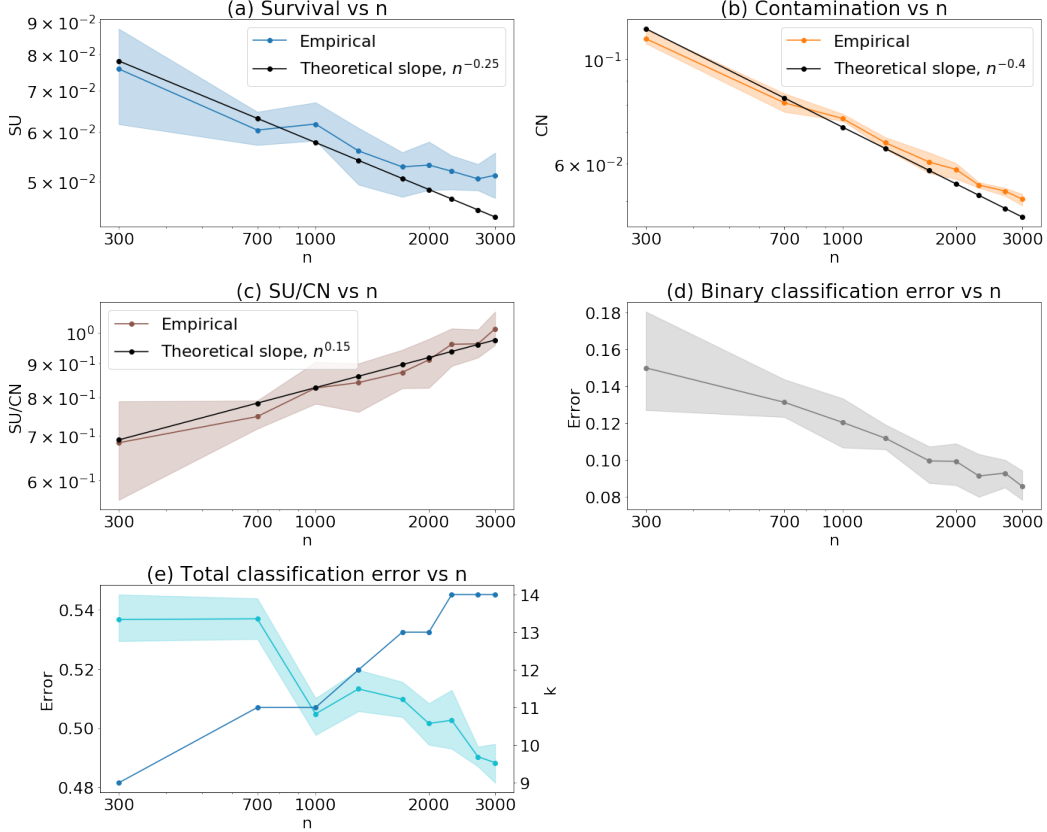


Figure 6: Experimental results using the bi-level ensemble model with $p = 1.5, q = 0.55, r = 0.5, t = 0.2$. Here, the number of training samples n varies from 300 to 3000 and the number of classes is computed as $k = b3n^t c$ and varies from 9 to 14. We calculated the classification errors over a batch size of 10000, and ran 10 trials. The plots show the mean plotted with error bars corresponding to the 10th, 90th percentile values. We also plot the theoretical slopes for survival, contamination and the survival/contamination ratio based on our conjecture and notice that it closely matches the empirical slope of the quantities when plotted on a log-log scale. Notice that jaggedness in the plots is often due to integer effects as k grows or does not grow with n .

scales as

$$\frac{\frac{ad}{s}}{\frac{(1-a)d}{d s}} = n^{p-q-r}, \quad (294)$$

and when $q+r < 1$, this ratio is larger than $d/n = n^{p-1}$. Fig. 7 shows empirical results of estimating the eigenvalues via the singular value decomposition of the training feature matrices. The visual distinction is quite striking. In the regime $q+r > 1$, the SVD of the training features matrix (and thus the empirical covariance matrix's eigenvalues) does not reveal that there are actually s favored features in the data. By contrast, in the regime $q+r < 1$, the SVD clearly shows an eigenvalue gap that reveals exactly what s is.

Theorem 5.1 shows that when $q+r > 1$ and there is not enough structure in the feature matrix to reliably estimate the top eigenvalues of the feature covariance matrix features, multiclass classification can still succeed for interpolating solutions as long as the number of classes does not increase too fast. It is because we are in the regime $q+r > 1$ that new techniques for analysis had to be developed in this paper, and the gap between the regime where we can prove the results and our conjectured results points interestingly to where there is a need for even better technique.

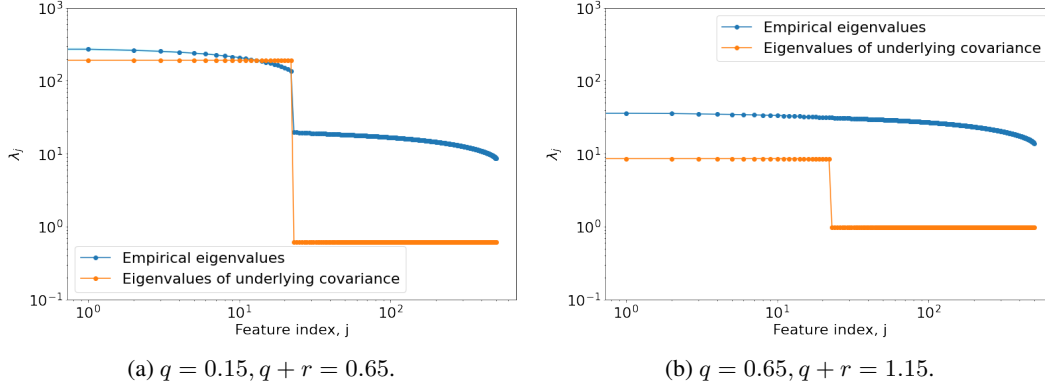


Figure 7: Estimating the eigenvalues of the covariance matrix of features empirically. Here $n = 400$ and the feature covariance structure follows the bi-level model with parameters $p = 1.5, r = 0.5$. Thus $d = 8001$ and $s = 21$ for both (a) and (b). The difference between (a) and (b) is in the level of favoring of favored features: with $q = 0.15$, (a) favors them more than (b) does with $q = 0.65$. In the regime where regression works, $q + r \leq 1$, we are able to accurately estimate the top s eigenvalues. In the regime where regression fails, $q + r > 1$, we are unable to estimate the top s eigenvalues accurately. The blue curve plots the estimated eigenvalues and the shaded region corresponds to the 10-90 percentile of the estimated values over 20 trials. Note, that there is only very small deviation across trials.

K A heuristic derivation of the main result

There is a significant amount technical work required in this paper to prove the main results, and Section F provides a conjecture for what the result should be, based on a more detailed look at our proof and where we suspect looseness. However, it is also possible to derive/conjecture the main result by another heuristic method in the style of Appendix B of Muthukumar et al. [55] and Appendix A of Muthukumar et al. [56]. At the heart of this heuristic method is an asymptotic linear-algebraic perspective driven by a few principles:

- When generalized linear models are sufficiently overparameterized, they tend to behave like the Fourier features model does on one-dimensional regularly spaced training data.
- If the favored features are approximately orthogonal on the training data, we can treat them as though they were exactly orthogonal on the training data as long as the number of favored features is sublinear in n , the number of training points.
- In the heavily overparameterized regime, for the bi-level model considered in this paper as well as in Muthukumar et al. [56] with s favored features, n training points, and d total features, min-norm-interpolation will behave as though there are $\frac{d}{n}$ exact aliases among the unfavored features for each of the favored features.
- Similarly, for this heavily overparameterized bi-level model, min-norm-interpolation will behave as though the unfavored features also provide another $n - s$ effectively orthogonal directions on the training data that also have $\frac{d}{n}$ exact aliases — all of these among the unfavored features.
- Any nonlinear function of a small group of training features can be viewed as having a linear part that is essentially captured by its projection onto that group of training features plus another component which is assumed to behave like "white noise" — equally spread across the n dimensions (assuming the size of the group is much smaller than n).
- On a test point, the different underlying features behave like uncorrelated random variables whose means are 0 and whose variances are proportional to the magnitude-squared of their coefficient times their inherent weighting squared.

These heuristics, together with other standard approximations like $1 + |x| \approx |x|$ when $|x| \gg 1$ and $1 + |x| \approx 1$ when $|x| \ll 1$ permit relatively simple calculations to be used together with the survival/contamination style of analysis to predict when regression and classification type problems

will succeed. This is done in Appendix B of Muthukumar et al. [56] for binary classification. Here, it is useful to recap the key insights from that paper in the context of the same bi-level model used here:

- When the underlying score function that we would like to learn is a linear combination of the favored features, the survival captures the extent to which that particular linear combination is present in the learned score. Without loss of generality, because all favored features are equally favored, it suffices to consider a desired score that is purely one of the favored features, normalized to have unit variance. In this case, the coefficient learned for that (normalized) feature represents the survival.
- The contamination captures everything else that is learned — all the learned coefficients of features other than the true feature — and is measured in terms of the standard deviation of the predictions due to those (falsely) learned coefficients.
- For learning in a regression problem to asymptotically generalize, the survival must tend to 1 and the contamination must tend to 0. This is because the score function itself must be learned in a mean-square-error sense.
- For learning in a binary classification problem to asymptotically generalize, all we need to get right is the sign of the learned score. To get the sign right, it suffices to have the ratio of survival to contamination go to infinity. This can happen even if both the survival and contamination tend to 0, as long as they do so at the right rates relative to each other.

K.1 Understanding what we learn from multiclass training data

Using the above heuristics and reusing calculations already done in Appendix A of Muthukumar et al. [56], we can see how survival and contamination asymptotically behave for multi-class training data. For class m , we consider the score function we learn by interpolating the zero-mean one-hot encoding for the n training points as represented by \mathbf{y}_m . Within \mathbf{y}_m , we are going to assume that there are exactly $\frac{n}{k}$ positive examples for the m -th class, and each of those is represented with a $1 - \frac{1}{k}$ in the appropriate position of \mathbf{y}_m . There are also $n - \frac{n}{k}$ negative examples for the m -th class, each represented with a $\frac{1}{k}$ in the appropriate position.

The total "energy" (norm-squared) in this vector \mathbf{y}_m is thus

$$\frac{n}{k} \left(1 - \frac{1}{k}\right)^2 + \left(n - \frac{n}{k}\right) \frac{1}{k^2} = \frac{n}{k} - \frac{n}{k^2}$$

where the final approximation is due to k asymptotically growing as n^t to infinity.

We can understand the process of min-norm interpolation as proceeding in two conceptual steps. First, this n -dimensional vector \mathbf{y}_m is decomposed into the s orthogonal directions represented by the s favored features and the $n - s$ unfavored synthetic directions according to the heuristic above. For this step, the level of favoring does not matter. After that, the level of favoring determines precisely how each of those directions is split across the representative favored feature (if any) and its $\frac{d}{n}$ unfavored aliases.

K.1.1 Survival

How much of \mathbf{y}_m ends up going into the direction representing the true favored feature m ? Because of orthogonality, we can simply look at the correlation between \mathbf{y}_m and a normalized vector in the direction of \mathbf{z}_m . We can approximate the standard Gaussian in \mathbf{z}_m which wins a competition with $k - 1$ other iid standard Gaussians as being a constant $\frac{1}{\sqrt{\ln k}}$, which as compared to the polynomial scalings relevant here, might as well just be the constant 1. With that, we can consider all the other entries of \mathbf{z}_m as being basically their mean, which by the same logic is essentially $\frac{1}{k}$. This gives a total correlation of

$$\frac{1}{n} \left(\frac{n}{k} \left(1 - \frac{1}{k}\right) + \left(n - \frac{n}{k}\right) \frac{1}{k} \right) = \frac{1}{n} \left(\frac{n}{k} - \frac{n}{k^2} + \frac{n}{k} - \frac{n}{k^2} \right) = \frac{1}{n} \left(\frac{2n}{k} - \frac{2n}{k^2} \right) \quad (295)$$

$$\frac{1}{n} \left(\frac{2n}{k} - \frac{2n}{k^2} \right) \quad (296)$$

with a normalized vector $\frac{\mathbf{z}_m}{k\mathbf{z}_m}$.

Notice that if we had done this correlation with \mathbf{z}_m itself instead of \mathbf{y}_m , we would have gotten simply $\frac{\rho_{\bar{n}}}{n}$. This shows how multiclass training data immediately reduces the survival by an asymptotic factor of k relative to noise-free regression training data.

K.1.2 Contamination

To understand contamination, we simply break it down into two sources: the other $s - 1 - s$ favored features and the unfavored features.

Here, we leverage the heuristic that the total $\frac{n}{k}$ energy of the training labels \mathbf{y}_m has to be split across the true feature m , the other k label-defining features that are not the true feature, the other $s - k - s$ favored features, and the unfavored features.

For positive training examples for class m , the other label-defining features have a mean value of $\frac{1}{k}$ since they are not the max. For negative training examples of class m , the mean value for other label-defining features is essentially zero. The consequence of this (by a calculation exactly analogous to (296)) is that the projected correlation for these is like $\frac{\rho_{\bar{n}}}{k^2}$ which might as well be zero.

How much energy is left in \mathbf{y}_m after we remove the components along the k label-defining directions that could have any linear relationship to it?

$$\frac{n}{k} - \left(\frac{\rho_{\bar{n}}}{k}\right)^2 = (k - 1) \left(\frac{\rho_{\bar{n}}}{k^2}\right)^2 = \frac{n}{k} - \frac{n}{k^2} - \frac{n}{k^3} \quad (297)$$

$$\frac{n}{k}. \quad (298)$$

Asymptotically, all the energy is still left. This can now be divided equally across the $n - k - n$ orthogonal directions by the heuristic.

This means that a fraction $\frac{s}{n}$ of this ends up as contamination in the favored features, which is a total energy of $\frac{s}{n} \frac{n}{k} = \frac{s}{k}$.

Meanwhile, there are $n - s$ other directions that are only represented by the unfavored features. This has total energy $\frac{n-s}{n} \frac{n}{k} = \frac{n-s}{k} \frac{n}{k}$.

K.1.3 When does survival dominate contamination?

There are qualitatively two kinds of contamination: coming from favored features and coming from unfavored features.

For favored features (like the true feature), there is a further split that happens as the min-norm-interpolation solution splits the coefficients themselves across the favored feature and its effectively $\frac{d}{n} = n^{p-1}$ unfavored aliases. Adapting the notation in Appendix A of Muthukumar et al. [56], let's say a fraction α goes onto the favored feature itself. Equation (23) in Muthukumar et al. [56] tells us that

$$\alpha = \begin{cases} 1 & \text{if } q < 1 - r \\ \frac{1}{n - (q - 1 - r)} & \text{if } q > 1 - r \end{cases} \quad (299)$$

Consequently, the actually survived signal scales as $\alpha \frac{\rho_{\bar{n}}}{k}$ while in those same units, the standard deviation of favored contamination scales as $\alpha \sqrt{\frac{s}{k}}$. From this we immediately see that the α cancel and the condition for classification working (as far as the favored features are concerned) is that

$$\frac{\rho_{\bar{n}}}{k} > \sqrt{\frac{s}{k}} \quad (300)$$

which implies

$$k > \frac{n}{s} \quad (301)$$

and since $k = n^t$ and $s = n^r$, this gives us a condition that $t < 1 - r$ for survival to dominate contamination from favored features.

This $t < 1 - r$ condition turns out not to care about whether we are in the regime where regression works (i.e. $q < 1 - r$) or not.

To understand the impact of contamination due to unfavored features, it is important to recall that all of that gets split across the $\frac{d}{n}$ aliases in our heuristic calculation and so the standard deviation of the contamination due to unfavored features scales as $\sqrt{\frac{n}{k} \frac{n}{d}} = \frac{\rho_n}{k d} = n^{1 - \frac{p}{2} - \frac{t}{2}}$. Meanwhile, survival behaves as $\alpha \frac{\rho_n}{k} = \alpha n^{\frac{1}{2} - t}$.

For survival to dominate contamination from unfavored features, we need:

$$\alpha n^{\frac{1}{2} - t} > n^{1 - \frac{p}{2} - \frac{t}{2}} \quad (302)$$

$$\alpha > n^{\frac{t - (p - 1)}{2}} \quad (303)$$

$$\alpha n^{\frac{p-1}{2}} > n^{\frac{t}{2}} \quad (304)$$

Substituting in for the case $q < 1 - r$ from (299), we immediately see the condition $t < p - 1$ for survival to dominate contamination from unfavored features.

For the case $q > 1 - r$, the substitution of $n^{-(q - (1 - r))}$ for α in (304) gives us:

$$\frac{t}{2} < \frac{p - 1}{2} - (q - (1 - r)) \quad (305)$$

$$t < (p - 1) - 2(q - (1 - r)). \quad (306)$$

Notice that at the boundary $q = (1 - r)$, the condition (306) matches the $t < p - 1$ condition for the region $q < 1 - r$.

Finally, note that our model requires $t < r$ so that the features that determine class labels are themselves favored. Putting all the terms together, we get the condition:

$$t < \min(r, 1 - r, (p - 1) - 2 \max(0, q - (1 - r))) \quad (307)$$

which agrees with the overall conjectured condition in (23).

K.2 Why is this sufficient for multiclass classification

Notice the heuristic derivation above is actually just saying that binary classification will succeed (i.e. we will typically have the m -th learned score agree in sign with the actual m -th feature for a test point) in the given regime, assuming we trained with multiclass training data.

But successful generalization for multiclass classification requires more — it requires the true class's score to win a competition against the competing classes.

Here, we can simply leverage a quick heuristic calculation involving order statistics. We know that the max of k i.i.d. random variables will concentrate to a neighborhood of where the Complementary Cumulative Distribution Function (CCDF) of the underlying random variable reaches $\frac{1}{k}$. Meanwhile, the second biggest of those variables will be around where CCDF hits $\frac{2}{k}$. For a standard Gaussian, these are basically at $\rho \frac{1}{2 \ln k}$ and then at $\sqrt{2(\ln k - \ln 2)}$. The gap between these is

$$\rho \frac{1}{2 \ln k} - \rho \frac{1}{2 \ln k - 2 \ln 2} = \frac{\rho}{2} \frac{1}{2 \ln k} \quad (308)$$

$$= \frac{\rho \ln 2}{2 \ln k} \quad (309)$$

and so what matters is that the contamination relative to the survival needs to go to zero faster than $\frac{\rho}{2 \ln k} = \frac{\rho}{2 \ln n}$. However, all the scalings here are polynomial in n and so if the contamination is asymptotically smaller than survival, it is smaller by a factor that is polynomial in n . This will beat the $\frac{\rho}{2 \ln n}$ scaling of this gap, resulting in successful multiclass classification.

Notice that we crucially used the Gaussian nature of the underlying features¹⁹ here. If the underlying features were uniformly distributed on $[-1, 1]$ instead, then the gap between the biggest and second biggest feature would instead decay polynomially like $\frac{1}{k}$. In that case, the relative contamination would have to go to zero fast enough to overcome this, making the conjectured multiclass region different.

¹⁹Any sufficiently thick tail would do. For example, an exponential or Laplacian feature would also work here.

K.3 Comments on the power of the heuristic calculation

The heuristic calculation being done here operated at a level of abstraction for which the fine details of the model do not matter. The total energy in the y_m just comes from the nature of multiclass training data where every point is labeled with only one class as its label and we have a balanced training set. The dominant contamination analysis follows simply from that.

All we need to verify is the survival analysis for alternative models. Here, it is also immediately clear that the asymptotic survival calculation would work just as well if one used a Gaussian mixture model with orthogonal class means or a logistic-softmax model with orthogonal classes — like the models in Wang et al. [77].

In fact, based on this heuristic calculation, we can conjecture that even if the training data were less informative (e.g. the training label did not actually correspond to the class whose feature had the true max but instead, just one of the positive feature classes was chosen uniformly at random as the training label), multiclass classification would still asymptotically succeed in the conjectured regimes in terms of reliably predicting the class whose feature is biggest. What matters is just that positive examples of a class show a significant amount of "classiness" on average in their training features.

Note however that the heuristic calculation here is just that — a heuristic calculation in what some might call the "physics style." Proving the conjectured region is in fact correct requires us to leverage the specific details of the models in our proofs.