
Precise Asymptotic Generalization for Multiclass Classification with Overparameterized Linear Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

We study the asymptotic generalization of an overparameterized linear model for multiclass classification under the Gaussian covariates bi-level model introduced in [Subramanian et al. \(2022\)](#), where the number of data points, features, and classes all grow together. We fully resolve the conjecture posed in [Subramanian et al. \(2022\)](#), matching the predicted regimes for generalization. Furthermore, our new lower bounds are akin to an information-theoretic strong converse: they establish that the misclassification rate goes to 0 or 1 asymptotically. One surprising consequence of our tight results is that the min-norm interpolating classifier can be asymptotically suboptimal relative to noninterpolating classifiers in the regime where the min-norm interpolating regressor is known to be optimal.

The key to our tight analysis is a new variant of the Hanson-Wright inequality which is broadly useful for multiclass problems with sparse labels. As an application, we show that the same type of analysis can be used to analyze the related multilabel classification problem under the same bi-level ensemble.

1 Introduction

In this paper, we directly follow up on a specific line of work initiated by [Subramanian et al. \(2022\)](#); [Anonymous \(2023\)](#). For the sake of self-containedness, we briefly reiterate the context, directing the reader to [Subramanian et al. \(2022\)](#) and the references cited therein for more. A broader story can be found in [Bartlett et al. \(2021\)](#); [Belkin \(2021\)](#); [Dar et al. \(2021\)](#); [Oneto et al. \(2023\)](#).

Classical statistical learning theory intuition predicts that highly expressive models, which can interpolate random labels ([Zhang et al., 2016; 2021](#)), ought not to generalize well. However, deep learning practice has seen such models performing well when trained with good labels. Resolving this apparent contradiction has recently been the focus of a multitude of works, and this paper builds on one particular thread of investigation that can be rooted in [Bartlett et al. \(2020\)](#); [Muthukumar et al. \(2020\)](#) where the concept of benign/harmless interpolation was crystallized in the context of overparameterized linear regression problems and conditions given for when this can happen. In [Muthukumar et al. \(2021\)](#), a specific toy "bi-level model" with Gaussian features was introduced to study overparameterized binary classification and show that successful generalization could happen even beyond the conditions for benign interpolation for regression. Following the introduction of the corresponding multi-class problem in [Wang et al. \(2021\)](#) with a constant number of classes, an asymptotic setting where the number of classes can grow with the number of training examples was introduced in [Subramanian et al. \(2022\)](#) where a conjecture was presented for when minimum-norm interpolating classifiers will generalize. We are now in a position to state our main contributions; afterwards, we expand on the related works.

35 Our contributions

36 Our main contribution is crisply identifying the asymptotic regimes where an overparameterized
37 linear model which performs minimum-norm interpolation does and does not generalize for multiclass
38 classification under a Gaussian features assumption, thus resolving the main conjecture posed by
39 (Subramanian et al., 2022). We improve on the analysis of Subramanian et al. (2022); Anonymous
40 (2023), covering all regimes with the asymptotically optimal misclassification rate. When the model
41 generalizes, it does so with a misclassification rate $o(1)$, and we show a matching "strong converse"
42 establishing when it misclassifies, it does so with rate $1 - o(1)$, where the explicit rate is nearly
43 identical to that of random guessing. The critical component of our analysis is a new variant of the
44 Hanson-Wright inequality, which applies to bilinear forms between a vector with subgaussian entries
45 and a vector that is bounded and has *soft sparsity*, a notion we will define in Section 4.2. We show
46 how this tool can be used to analyze other multiclass problems, such as multilabel classification.

47 1.1 Brief treatment of related work

48 Our thread begins with a recent line of work that analyzes the generalization behavior of overparam-
49 eterized linear models for regression (Hastie et al., 2022; Mei and Montanari, 2022; Bartlett et al., 2020;
50 Belkin et al., 2020; Muthukumar et al., 2020). These simple models demonstrate how the capacity
51 to interpolate noise can actually aid in generalization: training noise can be harmlessly absorbed
52 by the overparameterized model without contaminating predictions on test points. In effect, extra
53 features can be regularizing (in the context of descent algorithms' implicit regularization (Soudry
54 et al., 2018; Ji and Telgarsky, 2019; Engl et al., 1996; Gunasekar et al., 2018)), but an excessive
55 amount of such regularization causes regression to fail because even the true signal will not survive
56 the training process. Although works in this thread focus on very shallow networks, Chatterji and
57 Long (2023) established that deeper networks can behave similarly. Note that recently, Mallinar et al.
58 (2022) called-out an alternative regime (behaving like 1-nearest-neighbor learning) called "tempered"
59 overfitting in which training noise is not completely absorbed but the true signal does survive training.

60 The thread continues in a line of work that studies binary classification (Muthukumar et al., 2021;
61 Chatterji and Long, 2021; Wang and Thrampoulidis, 2021) in similar overparameterized linear models.
62 While confirming that the basic story is similar to regression, these works identify a further surprise:
63 binary classification can work in some regimes where the corresponding regression problem would
64 not work¹ due to the regularizing effect of overparameterization being too strong. Just as in the
65 regression case, the results here are sharp in toy models: we can exactly characterize where binary
66 classification using an interpolating classifier asymptotically generalizes.

67 With binary classification better understood, the thread continues to multiclass classification. After all,
68 the current wave of deep learning enthusiasm originated in breakthrough performance in multiclass
69 classification, and we have seen a decade of ever larger networks trained on ever larger datasets
70 with ever more classes Kaplan et al. (2020). Using similar toy models (Muthukumar et al., 2020;
71 2021; Wang et al., 2022), the constant number of classes case was studied in Wang et al. (2021) to
72 recover results similar to binary classification. Subramanian et al. (2022) further introduced a model
73 where the number of classes grows with the number of training points and proved an achievability
74 result on how fast the number of classes can grow while still allowing the interpolating classifier to
75 asymptotically generalize. While Subramanian et al. (2022) gave a conjecture for what the full region
76 should be, there was no converse proof, and they could not show generalization in entire conjectured
77 region. Anonymous (2023) proved a partial weak converse; they showed that the misclassification
78 rate is bounded away from 0 — rather than tending to 1 — in some of the predicted regimes.

79 2 Problem setup

80 We consider the multiclass classification problem with k classes. The following exposition is lifted
81 from Subramanian et al. (2022), but we include it for the sake of being self-contained. The training

¹Regression failing in the overparameterized regime is linked to the empirical covariance of the limited data not revealing the spiked reality of the underlying covariance (Wang and Fan, 2017). See Appendix J of Subramanian et al. (2022). When regression doesn't generalize, we also get "support-vector proliferation" in classification problems (Muthukumar et al., 2021; Hsu et al., 2021) which is also intimately related to the phenomenon of "neural collapse" (Papayan et al., 2020) as discussed, for example, in Xu et al. (2023).

82 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 standard Gaussian vectors². We assume
 83 that the labels $\ell_i \in [k]$ are generated as follows.

84 **Assumption 1** (1-sparse noiseless model). *The class labels ℓ_i are generated based on which of the*
 85 *first k dimensions of a point \mathbf{x}_i has the largest value,*

$$\ell_i = \arg \max_{m \in [k]} \mathbf{x}_i[m]. \quad (1)$$

86 For a vector \mathbf{x} , we index its j th entry with $\mathbf{x}[j]$. Hence, under Assumption 1, $\mathbf{x}_i[m]$ can be interpreted
 87 as how representative of class m the i th training point is.

88 For clarity of exposition, we make explicit a feature weighting that transforms the training points:

$$\mathbf{x}_i^w[j] = \sqrt{\lambda_j} \mathbf{x}_i[j] \quad \forall j \in [d]. \quad (2)$$

89 Here $\boldsymbol{\lambda} \in \mathbb{R}^d$ contains the squared feature weights. The feature weighting serves the role of favoring
 90 the true pattern, something that is essential for good generalization.³

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

$$\mathbf{X}^w = [\mathbf{x}_1^w \quad \cdots \quad \mathbf{x}_n^w]^\top = [\sqrt{\lambda_1} \mathbf{z}_1 \quad \cdots \quad \sqrt{\lambda_d} \mathbf{z}_d] \quad (3)$$

92 where we introduce the notation $\mathbf{z}_j \in \mathbb{R}^n$ to contain the j^{th} feature from the n training points. Note
 93 that $\mathbf{z}_j \sim N(0, \mathbf{I}_n)$ are i.i.d Gaussians. We use a one-hot encoding for representing the labels as the
 94 matrix $\mathbf{Y}^{\text{oh}} \in \mathbb{R}^{n \times k}$

$$\mathbf{Y}^{\text{oh}} = [\mathbf{y}_1^{\text{oh}} \quad \cdots \quad \mathbf{y}_k^{\text{oh}}], \quad \text{where} \quad \mathbf{y}_m^{\text{oh}}[i] = \begin{cases} 1, & \text{if } \ell_i = m \\ 0, & \text{otherwise} \end{cases}. \quad (4)$$

95 Since we consider linear models, we center the one-hot encodings and define

$$\mathbf{y}_m \triangleq \mathbf{y}_m^{\text{oh}} - \frac{1}{k} \mathbf{1}. \quad (5)$$

96 Our classifier consists of k coefficient vectors $\hat{\mathbf{f}}_m$ for $m \in [k]$ that are learned by minimum-norm
 97 interpolation (MNI) of the zero-mean one-hot variants using the weighted features:⁴

$$\hat{\mathbf{f}}_m = \arg \min_{\mathbf{f}} \|\mathbf{f}\|_2 \quad (6)$$

$$\text{s.t. } \mathbf{X}^w \mathbf{f} = \mathbf{y}_m. \quad (7)$$

98 We can express these coefficients in closed form as

$$\hat{\mathbf{f}}_m = (\mathbf{X}^w)^\top (\mathbf{X}^w (\mathbf{X}^w)^\top)^{-1} \mathbf{y}_m. \quad (8)$$

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

$$\hat{\ell} = \arg \max_{1 \leq m \leq k} \hat{\mathbf{f}}_m^\top \mathbf{x}_{\text{test}}^w. \quad (9)$$

²Following previous work, we are staying within a Gaussian features framework. However, recent developments have confirmed that these models are actually predictive when the features arise from nonlinearities in a lifting, as long as there is enough randomness underneath (Hu and Lu, 2022; Lu and Yau, 2022; Goldt et al., 2022; Misiakiewicz, 2022; McRae et al., 2022; Pesce et al., 2023; Kaushik et al., 2023).

³Our weighted feature model is equivalent to other works (e.g. Muthukumar et al. (2021)) that assume that the covariates come from a d -dimensional anisotropic Gaussian with a covariance matrix Σ that favors the truly important directions (Wei et al., 2022). These directions do not have to be axis-aligned — we make that assumption only for notational convenience. In reality, the optimizer will never know these directions *a priori*.

⁴The classifier learned via this method is equivalent to those obtained by other natural training methods (SVMs or gradient-descent with exponential tailed losses like cross-entropy) under sufficient overparameterization (Wang et al., 2021; Kaushik et al., 2023). Recently, Lai and Muthukumar (2023) showed via an extension of Ji and Telgarsky (2021) that a much broader category of losses also asymptotically result in convergence to the same MNI solution for sufficiently overparameterized classification problems.

By assumption, a misclassification event \mathcal{E}_{err} occurs whenever

$$\arg \max_{1 \leq m \leq k} \mathbf{x}_{\text{test}}[m] \neq \arg \max_{1 \leq m \leq k} \hat{\mathbf{f}}_m^\top \mathbf{x}_{\text{test}}^w. \quad (10)$$

We study where the MNI generalizes in an asymptotic regime where the number of training points, features, classes, and feature weights all scale according to the bi-level ensemble model⁵:

Definition 1 (Bi-level ensemble). *The bi-level ensemble is parameterized by p, q, r and t where $p > 1$, $0 \leq r < 1$, $0 < q < (p - r)$ and $0 \leq t < r$. 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), and number of classes (k) all scale with the number of training points (n) as follows:*

$$d = \lfloor n^p \rfloor, s = \lfloor n^r \rfloor, a = n^{-q}, k = c_k \lfloor n^t \rfloor, \quad (11)$$

where c_k is a positive integer. Define the feature weights by

$$\sqrt{\lambda_j} = \begin{cases} \sqrt{\frac{ad}{s}}, & 1 \leq j \leq s \\ \sqrt{\frac{(1-a)d}{d-s}}, & \text{otherwise} \end{cases}. \quad (12)$$

We introduce the notation $\lambda_F \triangleq \frac{ad}{s}$ and $\lambda_U \triangleq \frac{(1-a)d}{d-s}$ to distinguish between the (squared) favored and unfavored weights, respectively.

We visualize the bi-level model in Fig. 1, reproduced from Subramanian et al. (2022).

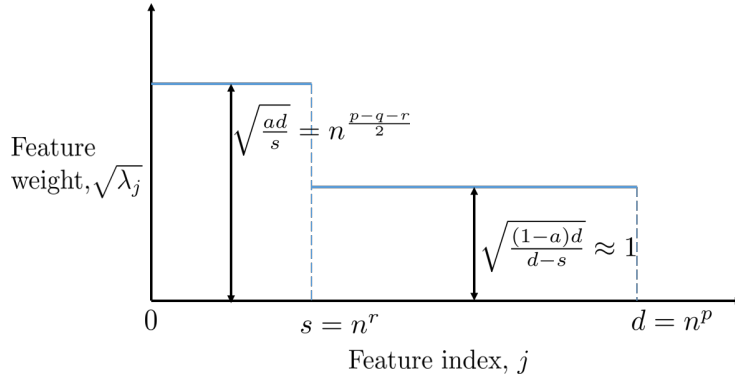


Figure 1: Bi-level feature weighting model. The first s features have a higher weight and are favored during minimum-norm interpolation. These can be thought of as the square-roots of the eigenvalues of the feature covariance matrix Σ in a Gaussian model for the covariates as in Bartlett et al. (2020).

3 Main results

In this section we state our main results and compare them to what was known and conjectured previously. Subramanian et al. (2022) use heuristic calculations to conjecture necessary and sufficient conditions for the bi-level model to generalize; we restate the conjecture here for reference.

Conjecture 3.1 (Conjectured bi-level regions). *Under the bi-level ensemble model (Definition 1), when the true data generating process is 1-sparse (Assumption 1), as $n \rightarrow \infty$, the probability of misclassification $\Pr[\mathcal{E}_{\text{err}}]$ for MNI as described in Eq. (6) satisfies*

$$\Pr[\mathcal{E}_{\text{err}}] \rightarrow \begin{cases} 0, & \text{if } t < \min \{1 - r, p + 1 - 2 \max \{1, q + r\}\} \\ 1, & \text{if } t > \min \{1 - r, p + 1 - 2 \max \{1, q + r\}\} \end{cases}. \quad (13)$$

Our main theorem establishes that Conjecture 3.1 indeed captures the correct generalization behavior of the overparameterized linear model.

⁵Such models are widely used to study learning even beyond this particular thread of work. For example, Tan et al. (2023) uses this to understand the privacy/generalization tradeoff of overparameterized learning.

Theorem 3.2 (Generalization for bi-level-model). *Under the bi-level ensemble model (Definition 1), when the true data generating process is 1-sparse (Assumption 1), Conjecture 3.1 holds.*

For comparison, we quote the best known previous positive and negative results for the bi-level model, which only hold in the restricted regime where regression fails ($q + r > 1$).

Theorem 3.3 (Generalization for bi-level model (Subramanian et al., 2022)). *In the same setting as Conjecture 3.1, in the regime where regression fails ($q + r > 1$), as $n \rightarrow \infty$ we have $\Pr[\mathcal{E}_{\text{err}}] \rightarrow 0$ if*

$$t < \min \{1 - r, p + 1 - 2(q + r), p - 2, 2q + r - 2\}. \quad (14)$$

Theorem 3.4 (Misclassification in bi-level model (Anonymous, 2023)). *In the same setting as Conjecture 3.1, in the regime where regression fails ($q + r > 1$), as $n \rightarrow \infty$ we have $\Pr[\mathcal{E}_{\text{err}}] \geq \frac{1}{2}$ if*

$$t > \min \{1 - r, p + 1 - 2(q + r)\}. \quad (15)$$

For ease of comparison between our main result and Theorems 3.3 and 3.4, we visualize the regimes in Fig. 2, as in Subramanian et al. (2022); Anonymous (2023). In particular, the blue starred and dashed regions in Fig. 2 indicate how Theorem 3.3 only applies where regression fails. In contrast, our new result holds regardless of whether regression fails or not, as in the green diamond region and light blue triangle regions. The regions are also completely tight; the looseness between the prior Theorem 3.3 and our result can be seen in the light blue square region.

The weak converse in the prior Theorem 3.4 captures some of the correct conditions for misclassification, but again only when $q + r > 1$. As depicted in the maroon X region for $r < 0.25$ in Fig. 2b, our main theorem gives a strong converse, whereas Theorem 3.4 has nothing to say because $q + r < 1$. Theorem 3.4 also only proves that the misclassification rate is asymptotically at least $\frac{1}{2}$. In the red circle and maroon X regions, we illustrate how our result pushes the misclassification rate to $1 - o(1)$, which requires a more refined analysis. We elaborate on this further in Section 4.

We remark that it is simpler to analyze the case where regression fails, as the random matrices that arise in the analysis are *flat*, i.e. approximately equal to a scaled identity matrix. However, in the regime where regression works, the same matrices have a *spiked* spectrum, which complicates the analysis. To smoothly handle both cases, we leverage a new variant of the Hanson-Wright inequality to show concentration of certain sparse bilinear forms; see Section 4.1 for more details.

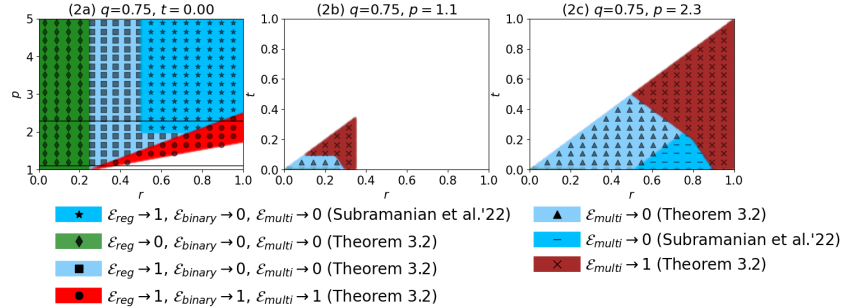


Figure 2: Example of regimes for multiclass/binary classification and regression. The white regions correspond to invalid regimes under the bi-level model. The entirety of 2b and all the light blue regions are new to this paper, as is showing that the error tends to 1 in the maroon regions.

4 Technical overview

We now sketch out the proof for our main theorem. As in Subramanian et al. (2022); Anonymous (2023), the starting point is writing out the necessary and sufficient conditions for misclassification.

Assume without loss of generality that the test point $\mathbf{x}_{\text{test}} \sim N(0, \mathbf{I}_d)$ has true label α for some $\alpha \in [k]$. Let $\mathbf{x}_{\text{test}}^w$ be the weighted version of this test point. From (10), an equivalent condition for misclassification is that for some $\beta \neq \alpha, \beta \in [k]$, we have $\hat{\mathbf{f}}_{\alpha}^{\top} \mathbf{x}_{\text{test}}^w < \hat{\mathbf{f}}_{\beta}^{\top} \mathbf{x}_{\text{test}}^w$, i.e. the score for β outcompetes the score for α . Define the Gram matrix $\mathbf{A} \triangleq \mathbf{X}^w (\mathbf{X}^w)^{\top}$, the relative label vector

155 $\Delta y \triangleq \mathbf{y}_\alpha - \mathbf{y}_\beta \in \{-1, 0, 1\}^n$, and the relative survival vector $\hat{\mathbf{h}}_{\alpha,\beta} \in \mathbb{R}^d$ which compares the signal
 156 from α and β :

$$\hat{\mathbf{h}}_{\alpha,\beta}[j] \triangleq \lambda_j^{-1/2}(\hat{\mathbf{f}}_\alpha[j] - \hat{\mathbf{f}}_\beta[j]) \quad (16)$$

$$= \mathbf{z}_j^\top \mathbf{A}^{-1} \Delta y, \quad (17)$$

157 where to obtain the last line we have used (8). By converting the misclassification condition into the
 158 unweighted feature space we see that we will have errors when

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

159 Define the contamination term $\text{CN}_{\alpha,\beta}$:

$$\text{CN}_{\alpha,\beta} \triangleq \sqrt{\sum_{j \notin \{\alpha,\beta\}} \lambda_j^2 (\hat{\mathbf{h}}_{\beta,\alpha}[j])^2}. \quad (19)$$

160 Note that $\text{CN}_{\alpha,\beta}$ normalizes the RHS of (18) into a standard Gaussian. Indeed, define

$$Z^{(\beta)} \triangleq \frac{1}{\text{CN}_{\alpha,\beta}} \sum_{j \notin \{\alpha,\beta\}} \lambda_j \hat{\mathbf{h}}_{\beta,\alpha}[j] \mathbf{x}_{\text{test}}[j] \sim N(0, 1). \quad (20)$$

161 Since $\alpha, \beta \in [k]$ are favored, we have $\lambda_\alpha = \lambda_\beta = \lambda_F$. Hence an equivalent condition for misclassifi-
 162 cation is that there exists some $\beta \neq \alpha, \beta \in [k]$ such that

$$\frac{\lambda_F}{\text{CN}_{\alpha,\beta}} (\hat{\mathbf{h}}_{\alpha,\beta}[\alpha] \mathbf{x}_{\text{test}}[\alpha] - \hat{\mathbf{h}}_{\beta,\alpha}[\beta] \mathbf{x}_{\text{test}}[\beta]) < Z^{(\beta)}. \quad (21)$$

163 We now translate the above criterion into *sufficient* conditions for correct classification and misclassi-
 164 fication and analyze these two cases separately.

165 **Correct classification:** For correct classification, it suffices for the maximum value of the LHS of
 166 Eq. (21) to outcompete the maximum value of the RHS, where the max is taken over $\beta \in [k], \beta \neq \alpha$.
 167 Some algebra, as in Subramanian et al. (2022), shows that we correctly classify if

$$\underbrace{\frac{\min_\beta \lambda_F \hat{\mathbf{h}}_{\alpha,\beta}[\alpha]}{\max_\beta \text{CN}_{\alpha,\beta}}}_{\text{SU/CN ratio}} \left(\underbrace{\min_\beta (\mathbf{x}_{\text{test}}[\alpha] - \mathbf{x}_{\text{test}}[\beta])}_{\text{closest feature margin}} - \underbrace{\max_\beta |\mathbf{x}_{\text{test}}[\beta]|}_{\text{largest competing feature}} \cdot \underbrace{\max_\beta \left| \frac{\hat{\mathbf{h}}_{\alpha,\beta}[\alpha] - \hat{\mathbf{h}}_{\beta,\alpha}[\beta]}{\hat{\mathbf{h}}_{\alpha,\beta}[\alpha]} \right|}_{\text{survival variation}} \right) > \underbrace{\max_\beta Z^{(\beta)}}_{\text{normalized contamination}}. \quad (22)$$

168 We will show that under the conditions specified in Conjecture 3.1, with high probability, the relevant
 169 survival to contamination ratio SU/CN grows at a polynomial rate n^v for some $v > 0$, whereas the
 170 term in the parentheses shrinks at a subpolynomial rate $\omega(n^{-\delta})$ for any $\delta > 0$. Further, by standard
 171 subgaussian maximal inequalities, the magnitudes of the *normalized contamination* is no more than
 172 $O(\sqrt{\log(nk)})$ with high probability. Thus with high probability the LHS outcompetes the RHS,
 173 leading to correct classification. See Section 4.1 for more discussion on how we prove tight bounds
 174 on the survival-to-contamination ratios.

175 **Misclassification:** On the other hand, for misclassification it suffices for the maximum *abso-*
 176 *lute* value of the LHS of Eq. (21) to be outcompeted by the maximum value of the RHS. Some
 177 manipulations yield the following sufficient condition for misclassification:

$$\underbrace{\frac{\max_\beta \lambda_F \left(\left| \hat{\mathbf{h}}_{\alpha,\beta}[\alpha] \right| + \left| \hat{\mathbf{h}}_{\beta,\alpha}[\beta] \right| \right)}{\min_\beta \text{CN}_{\alpha,\beta}}}_{\text{SU/CN ratio}} \cdot \underbrace{\max_{\gamma \in [k]} |\mathbf{x}_{\text{test}}[\gamma]|}_{\text{largest label-defining feature}} < \underbrace{\max_\beta Z^{(\beta)}}_{\text{normalized contamination}}. \quad (23)$$

178 We show that within the misclassification regimes in Conjecture 3.1, the survival-to-contamination
179 ratio SU/CN *shrinks* at a polynomial rate n^{-w} for some $w > 0$. By standard subgaussian maximal
180 inequalities, the largest label-defining feature is $O(\sqrt{\log(nk)})$ with high probability. Gaussian
181 anticoncentration implies that for some $\beta \neq \alpha, \beta \in [k]$, $Z^{(\beta)}$ outcompetes the LHS with probability at
182 least $\frac{1}{2} - o(1)$. Hence, we conclude that the model will misclassify with rate at least $\frac{1}{2}$ asymptotically.
183 Let us now describe how to boost the misclassification rate to $1 - o(1)$. Notice that the above
184 argument only considered the competition between the LHS of Eq. (23) and one of the $Z^{(\beta)}$'s on
185 the RHS instead of the maximum $Z^{(\beta)}$. It's not hard to see from the definition of $Z^{(\beta)}$ in Eq. (20)
186 that the $Z^{(\beta)}$ are jointly Gaussian. For intuition's sake, assuming the $Z^{(\beta)}$ were *independent*, then
187 $\max_{\beta} Z^{(\beta)}$ would outcompete with probability $(\frac{1}{2} - o(1))^{k-1}$.
188 In reality, the $Z^{(\beta)}$ are correlated, but we are able to show that the maximum correlation between the
189 $Z^{(\beta)}$ is $\frac{1}{2} + o(1)$ with high probability. An application of Slepian's lemma (Slepian (1962)) and some
190 explicit bounds on orthant probabilities (Pinasco et al. (2021)) implies that $\max_{\beta} Z^{(\beta)} > 0$ with
191 probability at least $1 - \frac{1}{k^{1+o(1)}}$. Another application of anticoncentration implies that $\max_{\beta} Z^{(\beta)} >$
192 n^{-w} with probability $1 - o(1)$, which finishes off the proof.

193 4.1 Bounding the survival-to-contamination ratio

194 Note that the critical *survival-to-contamination* ratio appears in both Eqs. (22) and (23). The most
195 involved part of the proof is nailing down the correct order of growth of the survival to contamination
196 ratio; a similar analysis tightly bounds the survival variation and the correlation structure of the $Z^{(\beta)}$.

197 To understand the relative survival and contamination, we must analyze the bilinear forms $\hat{h}_{\alpha,\beta}[j] =$
198 $\mathbf{z}_j^{\top} \mathbf{A}^{-1} \Delta \mathbf{y}$. Similarly, to control the correlation of the $Z^{(\beta)}$, we must understand the correlation
199 between the $\hat{h}_{\alpha,\beta}$ vectors, which reduces to understanding the bilinear forms $\mathbf{z}_j^{\top} \mathbf{A}^{-1} \mathbf{y}_{\alpha}$ for $j \in$
200 $[d], \alpha \in [k]$. The main source of inspiration for bounding these bilinear forms is the heuristic style of
201 calculation carried out in Appendix K of Subramanian et al. (2022) that leads to Conjecture 3.1.

202 To simplify the discussion, we temporarily restrict to the regime where regression fails ($q + r > 1$).
203 However, our main technical tool seamlessly generalizes to the regime where regression works
204 ($q + r < 1$). In the regime where regression fails, \mathbf{A}^{-1} turns out to have a *flat* spectrum: $\mathbf{A}^{-1} \approx \alpha \mathbf{I}$
205 for some constant $\alpha > 0$. Assume for now that \mathbf{A}^{-1} is *exactly* equal to a scaled identity matrix.
206 Then the survival is proportional to $\mathbf{z}_{\alpha}^{\top} \Delta \mathbf{y}$, which is a random inner product. Similarly, to bound the
207 contamination terms we must control the random inner product $\mathbf{z}_j^{\top} \Delta \mathbf{y}$ for $j \notin \{\alpha, \beta\}$.

208 Since $\Delta \mathbf{y}$ is a sparse vector — it only has $\frac{2n}{k}$ nonzero entries in expectation — a quick computation
209 reveals that $\mathbb{E}[\mathbf{z}_{\alpha}^{\top} \Delta \mathbf{y}] = \tilde{O}(\frac{n}{k})$ and $\mathbb{E}[\mathbf{z}_j^{\top} \Delta \mathbf{y}] = 0$. The deciding factor, then, is how tightly these
210 quantities concentrate around their means. A naïve application of Hoeffding implies a concentration
211 radius of order $\tilde{O}(\sqrt{n})$, which would lead to looseness in the overall result. The hope is to exploit
212 sparsity to get a concentration radius of order $\tilde{O}(\sqrt{n/k})$. This is where our new technical tool
213 Theorem 4.1 comes in, which may be of independent interest; we present it in the following section.

214 4.2 A new variant of the Hanson-Wright inequality

215 In reality, even in the regime where regression fails, \mathbf{A}^{-1} is not actually perfectly flat. Even worse,
216 in the regime where regression works, \mathbf{A}^{-1} is actually spiked. Thus, we cannot simply reduce the
217 bilinear form $\mathbf{z}_j^{\top} \mathbf{A}^{-1} \Delta \mathbf{y}$ to an inner product. Instead, we turn to the well-known Hanson-Wright
218 inequality (Rudelson and Vershynin, 2013), which tells us that quadratic forms of random vectors with
219 independent, mean zero, subgaussian entries concentrate around their mean. It was used extensively
220 to study binary classification (Muthukumar et al., 2021), and multiclass classification (Subramanian
221 et al., 2022; Anonymous, 2023).

222 However, just as Hoeffding is loose, so too is the standard form of Hanson-Wright, because it also
223 does not exploit sparsity. This motivates a new variant of Hanson-Wright which fully leverages the
224 (soft) sparsity inherent to multiclass problems with an increasing number of classes. We now formally
225 define the notions of soft and hard sparsity.

Variant	Assumptions on \mathbf{y}	Concentration radius
Classic quadratic ^a : $\mathbf{x}^\top \mathbf{M} \mathbf{x}$	same as \mathbf{x}	$\tilde{O}(\ \mathbf{M}\ _F)$
Sparse bilinear ^b : $\mathbf{x}^\top \mathbf{M}(\gamma \circ \mathbf{y})$	$\gamma_i \sim \text{Ber}(\pi)$, indep. of X_i but not Y_i	$\tilde{O}(\sqrt{\pi}\ \mathbf{M}\ _F)$
Sparse bilinear ^c : $\mathbf{x}^\top \mathbf{M}(\gamma \circ \mathbf{y})$	$\gamma_i \sim \text{Ber}(\pi)$, indep. of Y_i but not X_i	$\tilde{O}(\sqrt{\pi}\ \mathbf{M}\ _F)$
Theorem 4.1: $\mathbf{x}^\top \mathbf{M} \mathbf{y}$	$ Y_i \leq 1$ a.s., $\mathbb{E}Y_i^2 \leq \pi$	$\tilde{O}(\sqrt{\pi}\ \mathbf{M}\ _F)$

Table 1: Comparison of different variants of the Hanson-Wright inequality. In all variants, we assume that $(\mathbf{x}, \mathbf{y}) = (X_i, Y_i)_{i=1}^n$ are subgaussian, centered, and the pairs (X_i, Y_i) are independent across i . We use \circ to denote elementwise multiplication, which allows us to express hard sparsity with the sparsity mask $\gamma \in \{0, 1\}^n$. The concentration radius corresponds to the size of typical fluctuations guaranteed by the concentration inequality, i.e. the ϵ needed for high probability guarantees.

^a (Rudelson and Vershynin, 2013, Theorem 1.1); ^b (Park et al., 2022, Theorem 1); ^c (Anonymous, 2023, Theorem 4)

Definition 2 (Soft and hard sparsity). For $\pi \leq 1$, we say that random vector $\mathbf{y} = (Y_i)_{i=1}^n$ has soft sparsity at level π if $|Y_i| \leq 1$ almost surely and $\text{Var}(Y_i) \leq \pi$ for all i . On the other hand, we say that \mathbf{y} has hard sparsity at level π if at most a π fraction of the Y_i are nonzero.

In particular, our variant Theorem 4.1 below requires that one of the vectors in the bilinear form has soft sparsity at level π . Throughout, one should think of $\pi = o(1)$, and for us indeed $\pi = O(\frac{1}{k})$. One can check that a bounded random vector \mathbf{y} with hard sparsity level π must also have soft sparsity at level $O(\pi)$, so soft sparsity is more general for bounded random vectors. In Table 1 we compare our variant with several variants of Hanson-Wright which have appeared in the literature, some of which involve hard sparsity.

Define the subgaussian norm $\|\xi\|_{\psi_2}$ (Vershynin, 2018) as

$$\|\xi\|_{\psi_2} = \inf_{K > 0} \{K : \mathbb{E} \exp(\xi^2/K^2) \leq 2\}, \quad (24)$$

Theorem 4.1 (Hanson-Wright for bilinear forms with soft sparsity). Let $\mathbf{x} = (X_1, \dots, X_n) \in \mathbb{R}^n$ and $\mathbf{y} = (Y_1, \dots, Y_n) \in \mathbb{R}^n$ be random vectors such that (X_i, Y_i) are independent pairs of (possibly correlated) centered random variables such that $\|X_i\|_{\psi_2} \leq K$ and Y_i has soft sparsity at level π , i.e. $|Y_i| \leq 1$ almost surely, and $\mathbb{E}[Y_i^2] \leq \pi$. Assume that conditioned on Y_j , $\|X_j\|_{\psi_2} \leq K$. Then there exists an absolute constant $c > 0$ such that for all $\mathbf{M} \in \mathbb{R}^{n \times n}$ and $\epsilon \geq 0$ we have

$$\Pr[|\mathbf{x}^\top \mathbf{M} \mathbf{y} - \mathbb{E}[\mathbf{x}^\top \mathbf{M} \mathbf{y}]| > \epsilon] \leq 2 \exp\left(-c \min\left\{\frac{\epsilon^2}{K^2 \pi \|\mathbf{M}\|_F^2}, \frac{\epsilon}{K \|\mathbf{M}\|_F}\right\}\right). \quad (25)$$

The full proof of Theorem 4.1 is deferred to Appendix G. The main proof techniques are heavily inspired by those of Rudelson and Vershynin (2013); Zhou (2019); Park et al. (2022). However, the proof of Theorem 4.1 is actually simpler than in Park et al. (2022); Anonymous (2023), as bounded with soft sparsity turns out to be easier to work with than subgaussian with hard sparsity. We refer readers to Anonymous (2023) for a more in-depth discussion of how these new “sparse” variants overcome the limitations of previous proof techniques used to study classification problems.

We briefly illustrate how Theorem 4.1 can be used to get tighter results throughout our analysis. A quick calculation reveals that the label vectors $\Delta \mathbf{y}$ and \mathbf{y}_α both have soft sparsity at level $\pi = O(1/k)$. However, \mathbf{y}_α does not have hard sparsity as required by the variants in Park et al. (2022); Anonymous (2023). Since $\|\mathbf{M}\|_F^2 \leq n \|\mathbf{M}\|_2^2$, we obtain a concentration radius ϵ which scales like $\sqrt{n/k}$ rather than \sqrt{n} (obtained via vanilla Hanson-Wright) or n/\sqrt{k} (obtained via Cauchy-Schwarz). This gain is crucial to tightly analyzing the survival, contamination, and correlation structure.

4.3 Completing the proof sketch

Theorem 4.1 and the above insights about sparsity and independence allow us to prove the following bounds on the relative survival and contamination terms which are tight up to log factors; see the Appendix for more details. For brevity’s sake, we introduce the notation $\mu \triangleq n^{q+r-1}$.

257 **Proposition 4.2** (Bounds on relative survival). *Suppose we are in the bi-level model. With probability*
 258 *at least $1 - O(1/n)$,*

$$\lambda_F \hat{\mathbf{h}}_{\alpha, \beta}[\alpha] = \min \{ \mu^{-1}, 1 \} \Theta(n^{-\min \{t, \frac{1}{2}\}}) \sqrt{\log k}.$$

259 Next, we state our bounds on contamination.

260 **Proposition 4.3** (Lower bound on contamination). *Suppose we are in the bi-level model. Then with*
 261 *probability at least $1 - O(1/n)$, the contamination satisfies*

$$\text{CN}_{\alpha, \beta} = \underbrace{\min \{ \mu^{-1}, 1 \} \Theta(n^{\frac{r-t-1}{2}})}_{\text{favored features}} + \underbrace{\Theta(n^{\frac{1-t-p}{2}})}_{\text{unfavored features}}. \quad (26)$$

262 Translating the parameters in Propositions 4.2 and 4.3 we see that (i) the relative survival is diminished
 263 by a factor $1/k$ as long as $k = o(\sqrt{n})$, and a factor $1/\sqrt{n}$ for $k = \Omega(\sqrt{n})$ (this looseness ends up
 264 being negligible for the final result) and (ii) the contamination is diminished by a factor of $1/\sqrt{k}$.
 265 This essentially matches the expected behavior from the heuristic calculation in Subramanian et al.
 266 (2022). Together with some straightforward algebra, Propositions 4.2 and 4.3 allow us to compute
 267 the regimes where the survival-to-contamination ratio SU/CN grows or decays polynomially. This
 268 yields the stated regimes in Conjecture 3.1; see the Appendix for more details.

269 For technical reasons, the analogous bounds in Subramanian et al. (2022) are loose, giving rise to
 270 unnecessary conditions for good generalization such as $t < p - 2$ and $t < 2q + r - 2$. Moreover, we
 271 are able to give both upper and lower bounds on the survival and contamination terms, whereas they
 272 only give one sided inequalities for each quantity.

273 5 Discussion

274 In this paper we resolve the main conjecture of Subramanian et al. (2022), identifying the exact
 275 regimes where an overparameterized linear model succeeds at multiclass classification. Our tech-
 276 niques also lay the foundation for investigating related generalization for other multiclass tasks and
 277 nonlinear algorithms. We hope that by bringing the rigorous proofs closer to the heuristic style of
 278 calculation, we open the path for analyzing more complicated and realistic models.

279 As an example application, we sketch out how our proof techniques imply precise conditions for
 280 a variant of the learning task called multilabel classification. In a simple model for multilabel
 281 classification, each datapoint can have several of k possible labels — corresponding to the positive
 282 valued features — but in the training set only one such correct label is provided at random for each
 283 datapoint. We deem that the model generalizes if for any queried label it successfully labels test
 284 inputs as positive or negative. We can use the MNI approach here to learn classifiers.

285 Some thought reveals that the main difference between multilabel classification and multiclass
 286 classification from a survival and contamination perspective is that positive features no longer need
 287 to outcompete other features. Thus, the main object of study would be the bilinear forms $\mathbf{z}_j^\top \mathbf{A}^{-1} \mathbf{y}_\alpha$,
 288 which is possible thanks to Theorem 4.1. The survival and contamination terms are only affected by
 289 the expected values of these bilinear forms, but the expected values match the multiclass behavior up
 290 to log factors, which do not affect the regimes where SU/CN will grow or shrink polynomially. A
 291 similar analysis thus reveals that MNI will generalize in exactly the same regimes as in Conjecture 3.1.
 292 Here, the model generalizes in the sense that with high probability over the labels the model will
 293 correctly classify, and failure to generalize means that the model will do no better than a coin toss.

294 Perhaps surprisingly, resolving Conjecture 3.1 also implies that MNI is asymptotically *suboptimal*
 295 compared to a natural *non-interpolative* approach: simply make $\hat{\mathbf{f}}_m$ equal to the average⁶ of all
 296 positive training examples of class m . A straightforward analysis, detailed in the supplementary
 297 material, reveals this scheme fails to generalize exactly when $t < \min \{1 - r, p + 1 - 2(q + r)\}$,
 298 even in the regime where regression succeeds ($q + r < 1$). This is particularly interesting because
 299 we have shown that in the regime where regression succeeds, MNI generalizes only when $t <$
 300 $\min \{1 - r, p - 1\}$, which is a smaller region. In light of this gap, it would be interesting to identify
 301 the *information-theoretic* barrier for multiclass classification, especially within the broader context of
 302 statistical-computation gaps (see e.g. (Wu and Xu, 2021; Brennan and Bresler, 2020)).

⁶Note that Frei et al. (2023) point out that even leaky ReLU networks trained with a gradient flow can behave like averages of training examples.

References

- Robert J Adler, Jonathan E Taylor, et al. *Random fields and geometry*, volume 80. Springer, 2007.
- Anonymous Anonymous. Lower bounds for multiclass classification with overparameterized linear models. In *International Symposium on Information Theory*, 2023.
- 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.
- Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint. *Acta numerica*, 30:87–201, 2021.
- Mikhail Belkin. Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation. *Acta Numerica*, 30:203–248, 2021.
- 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.
- Matthew Brennan and Guy Bresler. Reducibility and statistical-computational gaps from secret leakage. In *Conference on Learning Theory*, pages 648–847. PMLR, 2020.
- 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.
- Niladri S Chatterji and Philip M Long. Deep linear networks can benignly overfit when shallow ones do. *Journal of Machine Learning Research*, 24(117):1–39, 2023.
- Victor Chernozhukov, Denis Chetverikov, and Kengo Kato. Comparison and anti-concentration bounds for maxima of gaussian random vectors. *Probability Theory and Related Fields*, 162(1): 47–70, 2015.
- 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.
- Heinz Werner Engl, Martin Hanke, and Andreas Neubauer. *Regularization of inverse problems*, volume 375. Springer Science & Business Media, 1996.
- Spencer Frei, Gal Vardi, Peter L Bartlett, and Nathan Srebro. Benign overfitting in linear classifiers and leaky relu networks from kkt conditions for margin maximization. *arXiv preprint arXiv:2303.01462*, 2023.
- Sebastian Goldt, Bruno Loureiro, Galen Reeves, Florent Krzakala, Marc Mézard, and Lenka Zdeborová. The Gaussian equivalence of generative models for learning with shallow neural networks. In *Mathematical and Scientific Machine Learning*, pages 426–471. PMLR, 2022.
- 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.
- Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *The Annals of Statistics*, 50(2):949–986, 2022.
- Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- 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.
- Hong Hu and Yue M Lu. Universality laws for high-dimensional learning with random features. *IEEE Transactions on Information Theory*, 2022.
- Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning Theory*, pages 1772–1798, 2019.

347 Ziwei Ji and Matus Telgarsky. Characterizing the implicit bias via a primal-dual analysis. In
348 *Algorithmic Learning Theory*, pages 772–804. PMLR, 2021.

349 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott
350 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.
351 *arXiv preprint arXiv:2001.08361*, 2020.

352 Chiraag Kaushik, Andrew D McRae, Mark A Davenport, and Vidya Muthukumar. New equivalences
353 between interpolation and svms: Kernels and structured features. *arXiv preprint arXiv:2305.02304*,
354 2023.

355 Kuo-Wei Lai and Vidya Muthukumar. General loss functions lead to (approximate) interpolation in
356 high dimensions. *arXiv preprint arXiv:2303.07475*, 2023.

357 Yue M Lu and Horng-Tzer Yau. An equivalence principle for the spectrum of random inner-product
358 kernel matrices. *arXiv preprint arXiv:2205.06308*, 2022.

359 Neil Mallinar, James Simon, Amirhesam Abedsoltan, Parthe Pandit, Misha Belkin, and Preetum
360 Nakkiran. Benign, tempered, or catastrophic: Toward a refined taxonomy of overfitting. *Advances*
361 *in Neural Information Processing Systems*, 35:1182–1195, 2022.

362 Andrew D McRae, Santhosh Karnik, Mark Davenport, and Vidya K Muthukumar. Harmless inter-
363 polation in regression and classification with structured features. In *International Conference on*
364 *Artificial Intelligence and Statistics*, pages 5853–5875. PMLR, 2022.

365 Song Mei and Andrea Montanari. The generalization error of random features regression: Precise
366 asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*, 75
367 (4):667–766, 2022.

368 Theodor Misiakiewicz. Spectrum of inner-product kernel matrices in the polynomial regime and
369 multiple descent phenomenon in kernel ridge regression. *arXiv preprint arXiv:2204.10425*, 2022.

370 Vidya Muthukumar, Kailas Vodrahalli, Vignesh Subramanian, and Anant Sahai. Harmless interpo-
371 lation of noisy data in regression. *IEEE Journal on Selected Areas in Information Theory*, 1(1):
372 67–83, 2020.

373 Vidya Muthukumar, Adhyayan Narang, Vignesh Subramanian, Mikhail Belkin, Daniel J. Hsu, and
374 Anant Sahai. Classification vs regression in overparameterized regimes: Does the loss function
375 matter? *Journal of Machine Learning Research*, 22:222:1–222:69, 2021.

376 Luca Oneto, Sandro Ridella, and Davide Anguita. Do we really need a new theory to understand
377 over-parameterization? *Neurocomputing*, page 126227, 2023.

378 Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal
379 phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):
380 24652–24663, 2020.

381 Seongoh Park, Xinlei Wang, and Johan Lim. Estimating high-dimensional covariance and precision
382 matrices under general missing dependence. *Electronic Journal of Statistics*, 15(2):4868–4915,
383 2021.

384 Seongoh Park, Xinlei Wang, and Johan Lim. Sparse Hanson-Wright inequality for a Bilinear Form of
385 Sub-Gaussian variables. *arXiv preprint arXiv:2209.05685*, 2022.

386 Luca Pesce, Florent Krzakala, Bruno Loureiro, and Ludovic Stephan. Are gaussian data all you need?
387 extents and limits of universality in high-dimensional generalized linear estimation. *arXiv preprint*
388 *arXiv:2302.08923*, 2023.

389 Damián Pinasco, Ezequiel Smucler, and Ignacio Zalduendo. Orthant probabilities and the attainment
390 of maxima on a vertex of a simplex. *Linear Algebra and its Applications*, 610:785–803, 2021.

391 Mark Rudelson and Roman Vershynin. Non-asymptotic theory of random matrices: extreme singular
392 values. In *Proceedings of the International Congress of Mathematicians 2010 (ICM 2010) (In 4*
393 *Volumes) Vol. I: Plenary Lectures and Ceremonies Vols. II–IV: Invited Lectures*, pages 1576–1602.
394 World Scientific, 2010.

395 Mark Rudelson and Roman Vershynin. Hanson-Wright inequality and sub-Gaussian concentration.
396 *Electronic Communications in Probability*, 18:1–9, 2013.

397 David S. Slepian. The one-sided barrier problem for gaussian noise. *Bell System Technical Journal*,
398 41:463–501, 1962.

399 Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The
400 implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(1):
401 2822–2878, 2018.

402 Vignesh Subramanian, Rahul Arya, and Anant Sahai. Generalization for multiclass classification
403 with overparameterized linear models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and
404 Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL
405 <https://openreview.net/forum?id=ikWvMRVQBWW>.

406 Jasper Tan, Daniel LeJeune, Blake Mason, Hamid Javadi, and Richard G Baraniuk. A blessing of
407 dimensionality in membership inference through regularization. In *International Conference on*
408 *Artificial Intelligence and Statistics*, pages 10968–10993. PMLR, 2023.

409 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,
410 volume 47. Cambridge university press, 2018.

411 Guillaume Wang, Konstantin Donhauser, and Fanny Yang. Tight bounds for minimum ℓ_1 -norm
412 interpolation of noisy data. In *International Conference on Artificial Intelligence and Statistics*,
413 pages 10572–10602. PMLR, 2022.

414 Ke Wang and Christos Thrampoulidis. Benign overfitting in binary classification of Gaussian mixtures.
415 In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing*
416 *(ICASSP)*, pages 4030–4034. IEEE, 2021.

417 Ke Wang, Vidya Muthukumar, and Christos Thrampoulidis. Benign Overfitting in Multiclass
418 Classification: All Roads Lead to Interpolation. *arXiv e-prints*, art. arXiv:2106.10865, June 2021.

419 Weichen Wang and Jianqing Fan. Asymptotics of empirical eigenstructure for high dimensional
420 spiked covariance. *Annals of statistics*, 45(3):1342, 2017.

421 Alexander Wei, Wei Hu, and Jacob Steinhardt. More than a toy: Random matrix models predict how
422 real-world neural representations generalize. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song,
423 Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International*
424 *Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,
425 pages 23549–23588. PMLR, 17–23 Jul 2022. URL [https://proceedings.mlr.press/v162/](https://proceedings.mlr.press/v162/wei22a.html)
426 [wei22a.html](https://proceedings.mlr.press/v162/wei22a.html).

427 Yihong Wu and Jiaming Xu. Statistical problems with planted structures: Information-theoretical
428 and computational limits. *Information-Theoretic Methods in Data Science*, 383:13, 2021.

429 Mengjia Xu, Akshay Rangamani, Qianli Liao, Tomer Galanti, and Tomaso Poggio. Dynamics in deep
430 classifiers trained with the square loss: Normalization, low rank, neural collapse, and generalization
431 bounds. *Research*, 6:0024, 2023.

432 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
433 deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.

434 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep
435 learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115,
436 2021.

437 Shuheng Zhou. Sparse Hanson–Wright inequalities for subgaussian quadratic forms. *Bernoulli*, 25
438 (3):1603–1639, 2019.

439	Contents	
440	1 Introduction	1
441	1.1 Brief treatment of related work	2
442	2 Problem setup	2
443	3 Main results	4
444	4 Technical overview	5
445	4.1 Bounding the survival-to-contamination ratio	7
446	4.2 A new variant of the Hanson-Wright inequality	7
447	4.3 Completing the proof sketch	8
448	5 Discussion	9
449	A Preliminaries and notation	15
450	A.1 Proof of Theorem 3.2	16
451	B Main tools	20
452	B.1 Hanson-Wright Inequality	20
453	B.2 Gram matrices and the Woodbury formula	20
454	B.3 Concentration of spectrum	22
455	C Utility bounds: applying the tools	23
456	D Bounding the survival	26
457	E Bounding the contamination	27
458	E.1 Upper bounding the contamination from label-defining+favored features	28
459	E.2 Lower bounding the contamination from label-defining+favored features	29
460	E.3 Bounding the unfavored contamination	31
461	F Obtaining tight misclassification rate	32
462	F.1 Main results for tight misclassification rates	33
463	F.2 Lower bounding the denominator	34
464	F.3 Upper bounding the numerator: the unnormalized correlation	34
465	F.3.1 Bounding the favored correlation	35
466	F.3.2 Bounding the unfavored correlation	37
467	G A new variant of the Hanson-Wright inequality	38
468	G.1 Diagonal terms	38
469	G.2 Offdiagonal terms	39
470	H Proofs of main lemmas for concentration of spectrum	40

471	I	Miscellaneous lemmas	43
472	J	Comparison to the straightforward non-interpolative scheme	44

473 A Preliminaries and notation

474 For positive integers n , we use the shorthand $[n] \triangleq \{1, \dots, n\}$. For a vector $\mathbf{v} \in \mathbb{R}^n$, $\|\mathbf{v}\|_2$ always
 475 denotes the Euclidean norm. We index entries by using square brackets, so $\mathbf{v}[j]$ denotes the j th
 476 entry of \mathbf{v} . For any matrix $\mathbf{M} \in \mathbb{R}^{m \times n}$, we denote its ij th entry by m_{ij} , $\|\mathbf{M}\|_2$ denotes the spectral
 477 norm, and $\|\mathbf{M}\|_F = \sqrt{\text{Tr}(\mathbf{M}^\top \mathbf{M})}$ denotes the Frobenius norm. We use $\sigma_{\max}(\mathbf{M})$ and $\sigma_{\min}(\mathbf{M})$ to
 478 denote the maximum and minimum singular values of \mathbf{M} , respectively. If $\mathbf{M} \in \mathbb{R}^{n \times n}$ is symmetric,
 479 we write $\mu_1(\mathbf{M}) \geq \mu_2(\mathbf{M}) \geq \dots \geq \mu_n(\mathbf{M})$ to denote the ordered eigenvalues of \mathbf{M} . Given two
 480 vectors $\mathbf{v}, \mathbf{u} \in \mathbb{R}^n$, we write $\mathbf{v} \circ \mathbf{u} \in \mathbb{R}^n$ to denote the entrywise product of \mathbf{v} and \mathbf{u} .

481 We make extensive use of big- O notation. In this paragraph, c refers to a positive constant which
 482 does not depend on n , and all statements hold for sufficiently large n . If $f(n) = O(g(n))$, then
 483 $f(n) \leq cg(n)$ for some c . If $f(n) = \tilde{O}(g(n))$, then $f(n) \leq cg(n) \log(n)$ for some c . If $f(n) =$
 484 $o(g(n))$, then for all $c > 0$ we have $f(n) \leq cg(n)$. We write $f(n) = \Omega(g(n))$ if $f(n) \geq cg(n)$
 485 for some c . Finally, we write $f(n) = \Theta(g(n))$ if there exists positive constants c_1 and c_2 such that
 $c_1 g(n) \leq f(n) \leq c_2 g(n)$.

Table 2: Notation

Symbol	Definition	Dimension	Source
k	Number of classes	Scalar	Sec. 2
n	Number of training points	Scalar	Sec. 2
d	Dimension of each point — the total number of features	Scalar	Sec. 2
s	The number of favored features	Scalar	Def. 1
a	The constant controlling the favored weights	Scalar	Def. 1
p	Parameter controlling overparameterization ($d = n^p$)	Scalar	Def. 1
r	Parameter controlling the number of favored features ($s = n^r$)	Scalar	Def. 1
q	Parameter controlling the favored weights ($a = n^{-q}$)	Scalar	Def. 1
t	Parameter controlling the number of classes ($k = c_k n^t$)	Scalar	Def. 1
c_k	The number of classes when $t = 0$ ($k = c_k n^t$)	Scalar	Def. 1
λ_j	Squared weight of the j th feature	Scalar	Def. 1
\mathbf{x}_i	i th training point (unweighted)	Length- d vector	Sec. 2
ℓ_i	Class label of i th training point	Scalar	Eqn. 1
\mathbf{w}_i	i th training point (weighted)	Length- d vector	Eqn. 2
\mathbf{X}^w	Weighted feature matrix	$(n \times d)$ -matrix	Eqn. 3
\mathbf{z}_j	The collected j th features of all training points	Length- n vector	Eqn. 3
\mathbf{y}_m^{oh}	One-hot encoding of all the training points for label m	Length- n vector	Eqn. 4
\mathbf{Y}^{oh}	One-hot label matrix	$(n \times k)$ -matrix	Eqn. 4
\mathbf{y}_m	Zero-mean encoding of the training points for label m	Length- n vector	Eqn. 5
$\hat{\mathbf{f}}_m$	Learned coefficients for label m using min-norm interpolation	Length- d vector	Eqn. 8
\mathbf{x}_{test}	A single test point	Length- d vector	Sec. 2
$\mathbf{x}_{\text{test}}^w$	A single weighted test point	Length- d vector	Sec. 2
\mathbf{A}	Gram matrix $\mathbf{A} = \mathbf{X}^w (\mathbf{X}^w)^\top$	$(n \times n)$ -matrix	Sec. 4
$\mu_i(\mathbf{A})$	The i th eigenvalue of matrix \mathbf{A} , sorted in descending order	Scalar	App. A
λ_F	Squared favored feature weights: $\lambda_F = \frac{ad}{s}$	Scalar	Def. 1
λ_U	Squared unfavored feature weights: $\lambda_U = \frac{(1-a)d}{d-s}$	Scalar	Def. 1
$\hat{\mathbf{h}}_{\alpha,\beta}$	Relative survival $\hat{\mathbf{h}}_{\alpha,\beta}[j] = \lambda_j^{-1/2}(\hat{f}_\alpha[j] - \hat{f}_\beta[j])$	Length- d vector	Eqn. 16
$\text{CN}_{\alpha,\beta}$	Normalizing factor $\text{CN}_{\alpha,\beta} = \sqrt{\left(\sum_{j \notin \{\alpha,\beta\}} \lambda_j^2 (\hat{\mathbf{h}}_{\beta,\alpha}[j])^2\right)}$	Scalar	Eqn. 19
$\ \cdot\ _{\psi_2}$	The sub-Gaussian norm of a scalar random variable	Scalar	Eqn. 24
μ	Factor controlling whether regression works, $\mu \triangleq n^{q+r-1}$	Scalar	App. A.1

486

487 Let us now describe the organization of the appendix. In Appendix A.1, we give a more detailed
 488 proof sketch and introduce the main propositions that complete the proof of Theorem 3.2. In
 489 Appendix B, we introduce the main tools that allow us to prove that the critical bilinear forms
 490 $\mathbf{z}_j^\top \mathbf{A}^{-1} \Delta \mathbf{y}$ concentrate: our new variant of the Hanson-Wright inequality, the Woodbury inversion

formula, and Wishart concentration to bound the spectra of the relevant random matrices that appear. In Appendix C we apply these tools to bound some useful quantities that repeatedly appear in the rest of the proofs. After that, we proceed to bound the survival, contamination, and correlation structure in Appendices D to F.

A.1 Proof of Theorem 3.2

In this section, we fill in some of the details of the proof sketch of Theorem 3.2. After recalling the beginning of the proof, we will split up the proof into two subtheorems: one for the positive result where MNI generalizes (Theorem A.4), and another for the negative result where MNI misclassifies (Theorem A.6).

Assume without loss of generality that the test point $\mathbf{x}_{\text{test}} \sim N(0, \mathbf{I}_d)$ has true label α for some $\alpha \in [k]$. Let $\mathbf{x}_{\text{test}}^w$ be the weighted version of this test point. From (10), an equivalent condition for misclassification is that for some $\beta \neq \alpha, \beta \in [k]$, we have $\widehat{\mathbf{f}}_\alpha^\top \mathbf{x}_{\text{test}}^w < \widehat{\mathbf{f}}_\beta^\top \mathbf{x}_{\text{test}}^w$, i.e. the score for β outcompetes the score for α . Define the Gram matrix $\mathbf{A} \triangleq \mathbf{X}^w (\mathbf{X}^w)^\top$, the relative label vector $\Delta \mathbf{y} \triangleq \mathbf{y}_\alpha - \mathbf{y}_\beta \in \{-1, 0, 1\}^n$, and the relative survival vector $\widehat{\mathbf{h}}_{\alpha, \beta} \in \mathbb{R}^d$ which compares the signal from α and β :

$$\widehat{\mathbf{h}}_{\alpha, \beta}[j] \triangleq \lambda_j^{-1/2} (\widehat{\mathbf{f}}_\alpha[j] - \widehat{\mathbf{f}}_\beta[j]) \quad (27)$$

$$= \mathbf{z}_j^\top \mathbf{A}^{-1} \Delta \mathbf{y}, \quad (28)$$

where to obtain the last line we have used the explicit formula for the MNI classifiers (8). By converting the misclassification condition into the unweighted feature space we see that we will have errors when

$$\lambda_\alpha \widehat{\mathbf{h}}_{\alpha, \beta}[\alpha] \mathbf{x}_{\text{test}}[\alpha] - \lambda_\beta \widehat{\mathbf{h}}_{\beta, \alpha}[\beta] \mathbf{x}_{\text{test}}[\beta] < \sum_{j \notin \{\alpha, \beta\}} \lambda_j \widehat{\mathbf{h}}_{\beta, \alpha}[j] \mathbf{x}_{\text{test}}[j]. \quad (29)$$

Define the contamination term $\text{CN}_{\alpha, \beta}$:

$$\text{CN}_{\alpha, \beta} \triangleq \sqrt{\sum_{j \notin \{\alpha, \beta\}} \lambda_j^2 (\widehat{\mathbf{h}}_{\beta, \alpha}[j])^2}. \quad (30)$$

Note that $\text{CN}_{\alpha, \beta}$ normalizes the RHS of (29) into a standard Gaussian. Indeed, define

$$Z^{(\beta)} \triangleq \frac{1}{\text{CN}_{\alpha, \beta}} \sum_{j \notin \{\alpha, \beta\}} \lambda_j \widehat{\mathbf{h}}_{\beta, \alpha}[j] \mathbf{x}_{\text{test}}[j] \sim N(0, 1). \quad (31)$$

Since $\alpha, \beta \in [k]$ are favored, we have $\lambda_\alpha = \lambda_\beta = \lambda_F$. Hence an equivalent condition for misclassification is that there exists some $\beta \neq \alpha, \beta \in [k]$ such that

$$\frac{\lambda_F}{\text{CN}_{\alpha, \beta}} (\widehat{\mathbf{h}}_{\alpha, \beta}[\alpha] \mathbf{x}_{\text{test}}[\alpha] - \widehat{\mathbf{h}}_{\beta, \alpha}[\beta] \mathbf{x}_{\text{test}}[\beta]) < Z^{(\beta)}. \quad (32)$$

We will translate the above criterion into *sufficient* conditions for correct classification and misclassification and analyze these two cases separately.

First, let us recall our tight characterization of the survival and contamination terms, which will be useful for both sides of the theorem. Recall our definition of $\mu \triangleq n^{q+r-1}$; whether this quantity polynomially shrinks or decays directly determines if regression works or fails.

Proposition A.1 (Bounds on relative survival). *Under the bi-level ensemble model (Definition 1), when the true data generating process is 1-sparse (Assumption 1), if $t < \frac{1}{2}$, then with probability at least $1 - O(1/nk)$*

$$\lambda_F \widehat{\mathbf{h}}_{\alpha, \beta}[\alpha] = c_7 \min\{\mu^{-1}, 1\} n^{-t} (1 \pm O(n^{-\kappa_5})) \sqrt{\log k},$$

where c_7 and κ_5 are positive constants.

If $t \geq \frac{1}{2}$, then

$$\lambda_F \left| \widehat{\mathbf{h}}_{\alpha, \beta}[\alpha] \right| \leq c_9 \min\{\mu^{-1}, 1\} n^{-\frac{1}{2}} \sqrt{\log(nk)},$$

where c_9 is a positive constant.

524 **Proposition A.2** (Bounds on contamination). *Under the bi-level ensemble model (Definition 1), when*
 525 *the true data generating process is 1-sparse (Assumption 1), with probability at least $1 - O(1/nk)$,*

$$\text{CN}_{\alpha,\beta} \leq \underbrace{\min\{\mu^{-1}, 1\} O(n^{\frac{r-t-1}{2}}) \log(nsk)}_{\text{favored features}} + \underbrace{O(n^{\frac{1-t-p}{2}}) \sqrt{\log(nsk)}}_{\text{unfavored features}}.$$

526 Furthermore, if $t > 0$, then with probability at least $1 - O(1/nk)$,

$$\text{CN}_{\alpha,\beta} \geq \underbrace{\min\{\mu^{-1}, 1\} \Omega(n^{\frac{r-t-1}{2}})}_{\text{favored features}} + \underbrace{\Omega(n^{\frac{1-t-p}{2}})}_{\text{unfavored features}}.$$

527 We defer the proof of Proposition A.1 to Appendix D and the proof of Proposition A.2 to Appendix E.
 528 Combining Propositions A.1 and A.2 yields the following sufficient conditions for when the SU/CN
 529 ratio grows or shrinks polynomially.

530 **Proposition A.3** (Regimes for survival-to-contamination). *Under the bi-level ensemble model (Def-*
 531 *inition 1), when the true data generating process is 1-sparse (Assumption 1), as $n \rightarrow \infty$, with*
 532 *probability at least $1 - O(1/n)$, the survival-to-contamination ratio satisfies*

$$\frac{\min_{\beta} \lambda_F \hat{\mathbf{h}}_{\alpha,\beta}[\alpha]}{\max_{\beta} \text{CN}_{\alpha,\beta}} \geq n^v \text{ for some } v > 0 \text{ if } t < \min\{1 - r, p + 1 - 2 \max\{1, q + r\}\} \quad (33)$$

$$\frac{\max_{\beta} \lambda_F |\hat{\mathbf{h}}_{\alpha,\beta}[\alpha]|}{\min_{\beta} \text{CN}_{\alpha,\beta}} \leq n^{-w} \text{ for some } w > 0 \text{ if } t > \min\{1 - r, p + 1 - 2 \max\{1, q + r\}\} \quad (34)$$

533 Here, the max and min are being taken over $\beta \neq \alpha, \beta \in [k]$.

534 *Proof.* We do casework on whether we want to prove an upper bound or lower bound
 535 on SU/CN. First, suppose we want to prove the lower bound, so assume $t <$
 536 $\min\{1 - r, p + 1 - 2 \max\{1, q + r\}\}$. Since $t < r$ by the definition of the bi-level ensemble
 537 (Definition 1), we have that $t < \frac{1}{2}$. So by union bounding over β , Proposition A.1 implies that with
 538 probability $1 - O(1/n)$

$$\min_{\beta} \lambda_F \hat{\mathbf{h}}_{\alpha,\beta}[\alpha] \geq \min\{\mu^{-1}, 1\} \Omega(n^{-t}) \sqrt{\log k}. \quad (35)$$

539 Then from Proposition A.2, by union bounding over β we see that with probability $1 - O(1/n)$,

$$\max_{\beta} \text{CN}_{\alpha,\beta} \leq \underbrace{\min\{\mu^{-1}, 1\} \tilde{O}(n^{\frac{r-t-1}{2}})}_{\text{favored features}} + \underbrace{\tilde{O}(n^{\frac{1-t-p}{2}})}_{\text{unfavored features}}.$$

540 Let us combine these two bounds. If we compare the survival to the contamination coming from
 541 favored features, we obtain

$$\frac{\min\{\mu^{-1}, 1\} n^{-t} \sqrt{\log k}}{\min\{\mu^{-1}, 1\} \tilde{O}(n^{\frac{r-t-1}{2}})} \geq \frac{n^{-t - \frac{r-t-1}{2}}}{\text{poly log}(n)} \quad (36)$$

$$\geq \frac{n^{\frac{1-r-t}{2}}}{\text{poly log}(n)}, \quad (37)$$

542 so in particular if $t < 1 - r$, the numerator grows polynomially and dominates the denominator. Now
 543 let's compare the survival to the contamination coming from unfavored fetures. This yields

$$\frac{\min\{\mu^{-1}, 1\} n^{-t} \sqrt{\log k}}{\tilde{O}(n^{\frac{1-t-p}{2}})} \geq \frac{\min\{\mu^{-1}, 1\} n^{-t - \frac{1-t-p}{2}}}{\text{poly log } n} \quad (38)$$

$$\geq \frac{n^{-\max\{q+r-1, 0\}} \cdot n^{\frac{p-t-1}{2}}}{\text{poly log}(n)} \quad (39)$$

$$\geq \frac{n^{\frac{p+1-2 \max\{1, q+r\}-t}{2}}}{\text{poly log}(n)}. \quad (40)$$

Hence, by union bounding, we see that with probability $1 - O(1/n)$,

$$\min_{\beta} \frac{\lambda_F \hat{\mathbf{h}}_{\alpha, \beta}[\alpha]}{\text{CN}_{\alpha, \beta}} \geq n^v, \quad (41)$$

where $v \triangleq \frac{1}{4}(\min\{1 - r, p + 1 - 2 \max\{1, q + r\}\} - t) > 0$ by assumption.

For the upper bound, suppose $t > \min\{1 - r, p + 1 - 2 \max\{1, q + r\}\}$. Hence $t > 0$, and by union bounding we conclude that with probability at least $1 - O(1/n)$,

$$\max_{\beta} \lambda_F \left| \hat{\mathbf{h}}_{\alpha, \beta}[\alpha] \right| \leq \min\{\mu^{-1}, 1\} O(n^{-\frac{1}{2}}) \sqrt{\log k} \quad (42)$$

and

$$\min_{\beta} \text{CN}_{\alpha, \beta} \geq \min\{\mu^{-1}, 1\} \Omega(n^{\frac{r-t-1}{2}}) + \Omega(n^{\frac{1-t-p}{2}}). \quad (43)$$

Combining these and union bounding yields that with probability $1 - O(1/n)$,

$$\min_{\beta} \frac{\lambda_F \hat{\mathbf{h}}_{\alpha, \beta}[\alpha]}{\text{CN}_{\alpha, \beta}} \leq n^{-w}, \quad (44)$$

where $w \triangleq \frac{1}{4}(t - \min\{1 - r, p + 1 - 2 \max\{1, q + r\}\}) > 0$ by assumption. \square

We now sketch out a proof of both the positive and negative sides of Theorem 3.2. We point out that the regimes for generalization and misclassification exactly match the regimes above for where the SU/CN ratio grows or shrinks polynomially.

Theorem A.4 (Positive side of Theorem 3.2). *Under the bi-level ensemble model (Definition 1), when the true data generating process is 1-sparse (Assumption 1), as $n \rightarrow \infty$, the probability of misclassification for MNI satisfies $\Pr[\mathcal{E}_{\text{err}}] \rightarrow 0$ if*

$$t < \min\{1 - r, p + 1 - 2 \max\{1, q + r\}\}.$$

Proof sketch. For correct classification, it suffices for the maximum value of the LHS of Eq. (32) to outcompete the maximum value of the RHS, where the max is taken over $\beta \in [k], \beta \neq \alpha$. Some algebra, as in Subramanian et al. (2022), shows that we correctly classify if

$$\underbrace{\frac{\min_{\beta} \lambda_F \hat{\mathbf{h}}_{\alpha, \beta}[\alpha]}{\max_{\beta} \text{CN}_{\alpha, \beta}}}_{\text{SU/CN ratio}} \left(\underbrace{\min_{\beta} (\mathbf{x}_{\text{test}}[\alpha] - \mathbf{x}_{\text{test}}[\beta])}_{\text{closest feature margin}} - \underbrace{\max_{\beta} |\mathbf{x}_{\text{test}}[\beta]|}_{\text{largest competing feature}} \cdot \underbrace{\max_{\beta} \left| \frac{\hat{\mathbf{h}}_{\alpha, \beta}[\alpha] - \hat{\mathbf{h}}_{\beta, \alpha}[\beta]}{\hat{\mathbf{h}}_{\alpha, \beta}[\alpha]} \right|}_{\text{survival variation}} \right) > \underbrace{\max_{\beta} Z^{(\beta)}}_{\text{normalized contamination}}. \quad (45)$$

By our lower bound on the survival to contamination ratio (Proposition A.3), assuming $t < \min\{1 - r, p + 1 - 2(q + r)\}$, then with probability at least $1 - O(1/n)$ we have that $\frac{\lambda_F \hat{\mathbf{h}}_{\alpha, \beta}[\alpha]}{\text{CN}_{\alpha, \beta}} \geq n^u$ for some constant $u > 0$. By Lemmas B.2 and B.3 in Subramanian et al. (2022) for every $\epsilon > 0$, with probability at least $1 - \epsilon$, we have $\min_{\beta} \mathbf{x}_{\text{test}}[\alpha] - \mathbf{x}_{\text{test}}[\beta] \geq \Omega(\frac{1}{\sqrt{\log k}})$.

Next, by standard subgaussian maxima tail bounds we have that $|\mathbf{x}_{\text{test}}[\beta]| \leq 2\sqrt{\log(nk)}$ and $Z^{(\beta)} \leq 2\sqrt{\log(nk)}$ with probability at least $1 - O(1/nk)$. Finally, applying our upper bound on the relative survival variance (Proposition A.5, which we prove below), the survival variation is at most a polynomially decaying n^{-w} with probability at least $1 - O(1/nk)$.

By union bounding, we see that with probability at least $1 - O(1/n) - \epsilon$, the LHS outcompetes the RHS, implying that the model correctly classifies.

\square

In fact, given Proposition A.1, it is straightforward to bound the survival variation.

572 **Proposition A.5** (Upper bound on the survival variation). *Suppose that $t < 1 - r$. With probability*
 573 *at least $1 - 2/n$, we have*

$$\left| \frac{\hat{\mathbf{h}}_{\alpha,\beta}[\alpha] - \hat{\mathbf{h}}_{\beta,\alpha}[\beta]}{\hat{\mathbf{h}}_{\alpha,\beta}[\alpha]} \right| \leq c_1 n^{-w}, \quad (46)$$

574 where c_1 and w are both positive constants.

575 *Proof.* Since we have $\hat{\mathbf{h}}_{\alpha,\beta}[\alpha] = \mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y$, the survival variation is

$$\frac{\hat{\mathbf{h}}_{\alpha,\beta}[\alpha] - \hat{\mathbf{h}}_{\beta,\alpha}[\beta]}{\hat{\mathbf{h}}_{\alpha,\beta}[\alpha]} = \frac{\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y + \mathbf{z}_\beta^\top \mathbf{A}^{-1} \Delta y}{\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y}$$

576 Since $t < 1 - r$ and $t < r$ by definition, we know that $t < \frac{1}{2}$ and we can apply Proposition A.1 to
 577 see that with probability at least $1 - 2/n$ we have

$$\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y = \max \{ \mu^{-1}, 1 \} n^{-t} (1 \pm O(n^{-\kappa_5})) \sqrt{\log k} = -\mathbf{z}_\beta^\top \mathbf{A}^{-1} \Delta y$$

578 Hence we have

$$\left| \frac{\hat{\mathbf{h}}_{\alpha,\beta}[\alpha] - \hat{\mathbf{h}}_{\beta,\alpha}[\beta]}{\hat{\mathbf{h}}_{\alpha,\beta}[\alpha]} \right| \leq c_1 n^{-\kappa_5} \quad (47)$$

579 where c_1 is an appropriately defined positive constant. \square

580 **Theorem A.6** (Negative side of Theorem 3.2). *Under the bi-level ensemble model (Definition 1),*
 581 *when the true data generating process is 1-sparse (Assumption 1), as $n \rightarrow \infty$, the probability of*
 582 *misclassification for MNI satisfies $\Pr[\mathcal{E}_{\text{err}}] \rightarrow 1$ if*

$$t > \min \{1 - r, p + 1 - 2 \max \{1, q + r\}\}.$$

583 *Proof sketch.* On the other hand, for misclassification it suffices for the maximum *absolute* value of
 584 the LHS of Eq. (32) to be outcompeted by the maximum value of the RHS. Some manipulations yield
 585 the following sufficient condition for misclassification:

$$\underbrace{\frac{\max_\beta \lambda_F \left(\left| \hat{\mathbf{h}}_{\alpha,\beta}[\alpha] \right| + \left| \hat{\mathbf{h}}_{\beta,\alpha}[\beta] \right| \right)}{\min_\beta \text{CN}_{\alpha,\beta}}}_{\text{SU/CN ratio}} \cdot \underbrace{\max_{\gamma \in [k]} |\mathbf{x}_{\text{test}}[\gamma]|}_{\text{largest label-defining feature}} < \underbrace{\max_\beta Z^{(\beta)}}_{\text{normalized contamination}}. \quad (48)$$

586 Within the misclassification regimes in Conjecture 3.1, Proposition A.3 implies that the survival-
 587 to-contamination ratio SU/CN *shrinks* at a polynomial rate n^{-w} for some $w > 0$. By standard
 588 subgaussian maximal inequalities, the largest label-defining feature is $O(\sqrt{\log(nk)})$ with high
 589 probability. Gaussian anticoncentration (Proposition I.2) implies that for some $\beta \neq \alpha, \beta \in [k]$, $Z^{(\beta)}$
 590 outcompetes the LHS, which is bounded above by n^{-w} , with probability at least $\frac{1}{2} - o(1)$. Hence,
 591 we conclude that the model will misclassify with rate at least $\frac{1}{2}$ asymptotically.

592 Let us now describe how to boost the misclassification rate to $1 - o(1)$. Notice that the above
 593 argument only considered the competition between the LHS of Eq. (48) and one of the $Z^{(\beta)}$'s on
 594 the RHS instead of the maximum $Z^{(\beta)}$. It's not hard to see from the definition of $Z^{(\beta)}$ in Eq. (31)
 595 that the $Z^{(\beta)}$ are jointly Gaussian. For intuition's sake, assuming the $Z^{(\beta)}$ were *independent*, then
 596 $\max_\beta Z^{(\beta)}$ would outcompete with probability $(\frac{1}{2} - o(1))^{k-1}$.

597 In reality, the $Z^{(\beta)}$ are correlated, but we are able to show that the maximum correlation between
 598 the $Z^{(\beta)}$ is $\frac{1}{2} + o(1)$ with high probability. An application of Slepian's lemma (Slepian (1962)) and
 599 some explicit bounds on orthant probabilities (Pinasco et al. (2021)) implies that $\max_\beta Z^{(\beta)} > 0$
 600 with probability at least $1 - \frac{1}{k^{1+o(1)}}$. An application of anticoncentration for Gaussian maxima
 601 (Chernozhukov et al., 2015) implies that $\max_\beta Z^{(\beta)} > n^{-w}$ with probability $1 - o(1)$, which finishes
 602 off the proof. \square

603 To fill in the details of the above proof sketch, we will prove the following proposition in Appendix F.
 604

Proposition A.7 (Correlation bound). *Assume we are in the bi-level ensemble model (Definition 1), the true data generating process is 1-sparse (Assumption 1), and the number of classes scales with n (i.e. $t > 0$). Then for every $\epsilon > 0$, we have*

$$\Pr \left[\max_{\beta \in [k], \beta \neq \alpha} Z^{(\beta)} > n^{-u} \right] \geq 1 - \Theta \left(\frac{1}{k^{1+o(1)}} \right) - \epsilon \quad (49)$$

for sufficiently large n and any $u > 0$.

B Main tools

In this section we introduce our suite of technical tools that allow us to prove the desired rates of growth for survival, contamination, and correlation.

B.1 Hanson-Wright Inequality

As established in Section 4, we need to use the Hanson-Wright inequality to prove our tight characterization of generalization. For the sake of precision, we explicitly state our definitions of subgaussian and subexponential which we use throughout the rest of the paper.

The subgaussian norm $\|\xi\|_{\psi_2}$ of a random variable ξ is defined as in Rudelson and Vershynin (2013),

$$\|\xi\|_{\psi_2} = \inf_{K > 0} \{K : \mathbb{E} \exp(\xi^2/K^2) \leq 2\}. \quad (50)$$

The sub-exponential norm $\|\xi\|_{\psi_1}$ is defined as in Vershynin (2018, Definition 2.7.5):

$$\|\xi\|_{\psi_1} = \inf_{K > 0} \{K : \mathbb{E} \exp(|\xi|/K) \leq 2\}. \quad (51)$$

We will occasionally need to use the following variant of Hanson-Wright for nonsparse bilinear forms, first proved in Park et al. (2021).

Theorem B.1 (Hanson-Wright for bilinear forms without sparsity). *Let $\mathbf{x} = (X_1, \dots, X_n) \in \mathbb{R}^n$ and $\mathbf{y} \in (Y_1, \dots, Y_n) \in \mathbb{R}^n$ be random vectors such that the pairs (X_i, Y_i) are all independent of each other (however X_i and Y_i can be correlated). Assume also that $\mathbb{E}[X_i] = \mathbb{E}[Y_i] = 0$ and $\max \{ \|X_i\|_{\psi_2}, \|Y_i\|_{\psi_2} \} \leq K$. Then there exists an absolute constant $c > 0$ such that for all $\mathbf{M} \in \mathbb{R}^{n \times n}$ and $\epsilon \geq 0$ we have*

$$\Pr [|\mathbf{x}^\top \mathbf{M} \mathbf{y} - \mathbb{E}[\mathbf{x}^\top \mathbf{M} \mathbf{y}]| > \epsilon] \leq 2 \exp \left(-c \min \left\{ \frac{\epsilon^2}{K^4 \|\mathbf{M}\|_F^2}, \frac{\epsilon}{K^2 \|\mathbf{M}\|_2} \right\} \right). \quad (52)$$

Finally, we restate our new version of Hanson-Wright for bilinear forms with soft sparsity, which we prove in Appendix G.

Theorem 4.1 (Hanson-Wright for bilinear forms with soft sparsity). *Let $\mathbf{x} = (X_1, \dots, X_n) \in \mathbb{R}^n$ and $\mathbf{y} \in (Y_1, \dots, Y_n) \in \mathbb{R}^n$ be random vectors such that (X_i, Y_i) are independent pairs of (possibly correlated) centered random variables such that $\|X_i\|_{\psi_2} \leq K$ and Y_i has soft sparsity at level π , i.e. $|Y_i| \leq 1$ almost surely, and $\mathbb{E}[Y_i^2] \leq \pi$. Assume that conditioned on Y_j , $\|X_j\|_{\psi_2} \leq K$. Then there exists an absolute constant $c > 0$ such that for all $\mathbf{M} \in \mathbb{R}^{n \times n}$ and $\epsilon \geq 0$ we have*

$$\Pr [|\mathbf{x}^\top \mathbf{M} \mathbf{y} - \mathbb{E}[\mathbf{x}^\top \mathbf{M} \mathbf{y}]| > \epsilon] \leq 2 \exp \left(-c \min \left\{ \frac{\epsilon^2}{K^2 \pi \|\mathbf{M}\|_F^2}, \frac{\epsilon}{K \|\mathbf{M}\|_2} \right\} \right). \quad (25)$$

B.2 Gram matrices and the Woodbury formula

In order to apply Hanson-Wright to the bilinear form $\mathbf{x}^\top \mathbf{M} \mathbf{y}$, we need to have a deterministic matrix \mathbf{M} such that the hypotheses are satisfied. However, in our setting we study bilinear forms such as $\mathbf{z}_j^\top \mathbf{A}^{-1} \Delta \mathbf{y}$. Here, the inverse Gram matrix \mathbf{A}^{-1} is not independent of \mathbf{z}_j or $\Delta \mathbf{y}$, so we cannot simply condition on \mathbf{A}^{-1} . The way around this is to cleverly decompose \mathbf{A}^{-1} using the so-called Woodbury inversion formula (stated formally below), which generalizes the leave-one-out trick and

Sherman-Morrison used to study binary classification in [Muthukumar et al. \(2021\)](#). To that end, we will explicitly decompose the Gram matrix $\mathbf{A} \triangleq \sum_{j \in [d]} \lambda_j \mathbf{z}_j \mathbf{z}_j^\top$ based on whether the features \mathbf{z}_j are favored or not.

We now introduce some notation to keep track of which matrices contain or leave out which indices. In general, we use subscripts to denote which sets of features we preserve or leave out; we use a minus sign to signify leaving out. The k label-defining features are represented with a subscript k , whereas the $s - k$ favored but not label defining features are represented with a subscript F . The rest of the $d - s$ unfavored features are represented with a subscript U .

For notational convenience, we introduce some new notation for the weighted features, as the superscript w to denote weighted features is rather cumbersome. We denote the weighted label-defining feature matrix by $\mathbf{W}_k \triangleq [\mathbf{w}_1 \ \cdots \ \mathbf{w}_k] \in \mathbb{R}^{n \times k}$, where the vectors $\mathbf{w}_i \triangleq \sqrt{\lambda_i} \mathbf{z}_i \in \mathbb{R}^n$ denote the weighted observations for feature i . Define the unweighted label-defining feature matrix $\mathbf{Z}_k \triangleq [\mathbf{z}_1 \ \cdots \ \mathbf{z}_k] \in \mathbb{R}^{n \times k}$. Similarly, define $\mathbf{W}_F \triangleq [\mathbf{w}_{k+1} \ \cdots \ \mathbf{w}_s] \in \mathbb{R}^{n \times (s-k)}$, which contains the rest of the weighted favored features and the corresponding unweighted version \mathbf{Z}_F .

Let $\mathbf{A}_{-k} \triangleq \sum_{i \notin [k]} \mathbf{w}_i \mathbf{w}_i^\top$ denote the leave- k -out Gram matrix which removes the k label-defining features. Similarly let $\mathbf{A}_{-F} \triangleq \sum_{i \notin [s] \setminus [k]} \mathbf{w}_i \mathbf{w}_i^\top \in \mathbb{R}^{n \times n}$ to denote leave- $(s - k)$ -out Gram matrix which removes the favored but not label-defining features. Finally, let $\mathbf{A}_U \triangleq \sum_{i \notin [s]} \mathbf{w}_i \mathbf{w}_i^\top \in \mathbb{R}^{n \times n}$ denote the leave- s -out matrix which only retains the unfavored features. We will also sometimes write \mathbf{A}_{-s} instead of \mathbf{A}_U to emphasize that the s favored features have all been removed.

Define the so-called hat matrices by

$$\mathbf{H}_k \triangleq \mathbf{W}_k^\top \mathbf{A}_{-k}^{-1} \mathbf{W}_k \in \mathbb{R}^{k \times k} \quad (53)$$

$$\mathbf{H}_F \triangleq \mathbf{W}_F^\top \mathbf{A}_{-F}^{-1} \mathbf{W}_F \in \mathbb{R}^{(s-k) \times (s-k)}. \quad (54)$$

These hat matrices appear in the Woodbury inversion formula. For the sake of notational compactness, define

$$\mathbf{M}_k \triangleq \mathbf{W}_k (\mathbf{I}_k + \mathbf{H}_k)^{-1} \mathbf{W}_k^\top \in \mathbb{R}^{n \times n} \quad (55)$$

$$\mathbf{M}_F \triangleq \mathbf{W}_F (\mathbf{I}_{s-k} + \mathbf{H}_F)^{-1} \mathbf{W}_F^\top \in \mathbb{R}^{n \times n}. \quad (56)$$

The Woodbury inversion formula yields

$$\mathbf{A}^{-1} = (\mathbf{W}_k \mathbf{W}_k^\top + \mathbf{A}_{-k})^{-1} \quad (57)$$

$$= \mathbf{A}_{-k}^{-1} - \mathbf{A}_{-k}^{-1} \mathbf{W}_k (\mathbf{I}_k + \mathbf{H}_k)^{-1} \mathbf{W}_k^\top \mathbf{A}_{-k}^{-1} \quad (58)$$

$$= \mathbf{A}_{-k}^{-1} - \mathbf{A}_{-k}^{-1} \mathbf{M}_k \mathbf{A}_{-k}^{-1}. \quad (59)$$

Left multiplying (58) by \mathbf{W}_k^\top yields

$$\mathbf{W}_k^\top \mathbf{A}^{-1} = \mathbf{W}_k^\top \mathbf{A}_{-k}^{-1} - \mathbf{H}_k (\mathbf{I}_k + \mathbf{H}_k)^{-1} \mathbf{W}_k^\top \mathbf{A}_{-k}^{-1} \quad (60)$$

$$= (\mathbf{I}_k - \mathbf{H}_k (\mathbf{I}_k + \mathbf{H}_k)^{-1}) \mathbf{W}_k^\top \mathbf{A}_{-k}^{-1} \quad (61)$$

$$= (\mathbf{I}_k + \mathbf{H}_k)^{-1} \mathbf{W}_k^\top \mathbf{A}_{-k}^{-1}. \quad (62)$$

We can derive completely analogous identities using \mathbf{A}_{-F}^{-1} instead of \mathbf{A}_{-k}^{-1} . The above exposition is summarized by the following lemma.

Lemma B.2. *We have*

$$\mathbf{W}_k^\top \mathbf{A}^{-1} \Delta y = (\mathbf{I}_k + \mathbf{H}_k)^{-1} \mathbf{W}_k^\top \mathbf{A}_{-k}^{-1} \Delta y \quad (63)$$

$$\mathbf{W}_F^\top \mathbf{A}^{-1} \Delta y = (\mathbf{I}_{s-k} + \mathbf{H}_F)^{-1} \mathbf{W}_F^\top \mathbf{A}_{-F}^{-1} \Delta y. \quad (64)$$

Lemma B.2 is quite powerful. Indeed, consider the action of the linear operator $\mathbf{W}_k^\top \mathbf{A}^{-1} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ on Δy . The action is identical to that of the linear operator $\mathbf{W}_k^\top \mathbf{A}_{-k}^{-1} : \mathbb{R}^n \rightarrow \mathbb{R}^k$, up to some invertible transformation. This new linear operator is nice because \mathbf{A}_{-k}^{-1} is independent of \mathbf{W}_k and Δy , as it removes all of the label-defining features. Reclaiming independence sets the stage for using our variant of Hanson-Wright.

How does the invertible operator $(\mathbf{I}_k + \mathbf{H}_k)^{-1}$ act? Our general strategy is to show that \mathbf{H}_k is itself close to a scaled identity matrix, i.e. $\mathbf{H}_k \approx \nu \mathbf{I}_k$ for an appropriately defined ν . Then for any $i \in [k]$, we have that

$$\mathbf{w}_i^\top \mathbf{A}^{-1} \Delta y \approx (1 + \nu)^{-1} \mathbf{w}_i^\top \mathbf{A}_{-k}^{-1} \Delta y.$$

Of course, there will be some error in this approximation, as \mathbf{H}_k is not *exactly* equal to $\nu \mathbf{I}_k$. Nevertheless, we can bound away the error that arises from this approximation.

B.3 Concentration of spectrum

As foreshadowed in the previous section, we will leverage the fact that the hat matrices such as \mathbf{H}_k are close to a scaled identity. To formalize this, we appeal to random matrix theory and show that the spectra of various random matrices are very close to being flat (i.e. all eigenvalues are within $1 + o(1)$ of each other). To that end, we present the following standard characterization of the spectrum of a standard Wishart matrix, which is Equation 2.3 in [Rudelson and Vershynin \(2010\)](#).

Lemma B.3 (Concentration of spectrum for Wishart matrices). *Let $\mathbf{M} \in \mathbb{R}^{M \times m}$ with $M > m$ be a real matrix with iid $N(0, 1)$ entries. Then for any $\epsilon \geq 0$, we have with probability at least $1 - 2e^{-\epsilon^2/2}$ that*

$$\sqrt{M} - \sqrt{m} - \epsilon \leq \sigma_{\min}(\mathbf{M}) \leq \sigma_{\max}(\mathbf{M}) \leq \sqrt{M} + \sqrt{m} + \epsilon. \quad (65)$$

In other words, the singular values of \mathbf{M} satisfy subgaussian concentration.

Since $\mu_m(\mathbf{M}^\top \mathbf{M}) = \sigma_{\min}(\mathbf{M})^2$ and $\mu_1(\mathbf{M}^\top \mathbf{M}) = \sigma_{\max}(\mathbf{M})^2$, we can conclude that if $m = o(M)$, then for any $\epsilon > 0$ we have

$$M - 2\sqrt{Mm} - \epsilon + o(\sqrt{Mm}) \leq \mu_m(\mathbf{M}^\top \mathbf{M}) \leq \mu_1(\mathbf{M}^\top \mathbf{M}) \leq M + 2\sqrt{Mm} + \epsilon + o(\sqrt{Mm}), \quad (66)$$

with probability at least $1 - 2e^{-\epsilon^2/2}$.

On the other hand, consider $\mathbf{M}\mathbf{M}^\top \in \mathbb{R}^{M \times M}$. Its spectrum is just that of $\mathbf{M}^\top \mathbf{M} \in \mathbb{R}^{m \times m}$ with an additional $M - m$ zeros corresponding to the fact that $m < M$.

We can use Lemma B.3 to prove concentration of the spectrum of the various matrices introduced in Appendix B.2. Let us summarize some convenient forms of these results; their proofs are deferred to Appendix H.

Proposition B.4 (Gram matrices have a flat spectrum). *Recall that $\mathbf{A}_U = \mathbf{A}_{-s} = \sum_{j>s} \lambda_j \mathbf{z}_j \mathbf{z}_j^\top \in \mathbb{R}^{n \times n}$ is the unfavored Gram matrix and $\mathbf{A}_{-k} = \sum_{j>k} \lambda_j \mathbf{z}_j \mathbf{z}_j^\top \in \mathbb{R}^{n \times n}$ is the leave- k -out Gram matrix.*

Then the following hold with probability at least $1 - 2e^{-n} - 2e^{-\sqrt{n}}$,

(a) *For all $i \in [n]$, we have $\mu_i(\mathbf{A}_U) = n^p(1 \pm O(n^{-\kappa_7}))$.*

(b) *For all $i \in [s - k]$, we have*

$$\mu_i(\mathbf{A}_{-k}) = (1 + \mu^{-1})n^p(1 \pm O(n^{-\kappa_9})), \quad (67)$$

where κ_9 is a positive constant. Moreover, for all $i \in [n] \setminus [s - k]$, we have

$$\mu_i(\mathbf{A}_{-k}) = n^p(1 \pm O(n^{-\kappa_7})), \quad (68)$$

where κ_7 is a positive constant.

As a simple corollary, we can obtain the following cruder bounds on the trace and spectral norm of \mathbf{A}_{-k}^{-1} and \mathbf{A}_{-s}^{-1} .

Corollary B.5 (Trace and spectral norm of \mathbf{A}_{-k}^{-1}). *In the bi-level model, with probability at least $1 - 2e^{-n}$, we have*

$$\text{Tr}(\mathbf{A}_U^{-1}) = n^{1-p}(1 \pm O(n^{-\kappa_7}))\sqrt{\log k} \quad (69)$$

$$\text{Tr}(\mathbf{A}_{-k}^{-1}) = n^{1-p}(1 \pm O(n^{-\kappa_3}))\sqrt{\log k} \quad (70)$$

and

$$\max \left\{ \|\mathbf{A}_{-k}^{-1}\|_2, \|\mathbf{A}_U^{-1}\|_2 \right\} \leq c_2 n^{-p}, \quad (71)$$

where c_2 , κ_7 , and κ_3 are all positive constants.

707 *Proof.* We prove the claim for \mathbf{A}_{-k}^{-1} ; the proof for \mathbf{A}_U^{-1} is similar or easier because \mathbf{A}_U^{-1} has a flat
708 spectrum (Proposition B.4).

709 If $q + r < 1$, the upper bound for the spectral norm similarly follows. For the trace bounds, we can
710 apply Proposition B.4, we have

$$\text{Tr}(\mathbf{A}_{-k}^{-1}) = (n - n^r + n^t)n^{-p}(1 \pm O(n^{-\kappa_7})) + (n^r - n^t) \cdot (1 + \mu^{-1})n^{-p}(1 \pm O(n^{-\kappa_9})) \quad (72)$$

$$= n^{1-p}(1 \pm O(n^{-\kappa_1})) \quad (73)$$

711 where

$$\kappa_1 = \min \{r - 1, 2 - q - 2r\} > 0,$$

712 as $q + 2r < 2(q + r) < 2$ by assumption.

713 On the other hand, the claim is obviously true when $q + r > 1$, as the entire spectrum of \mathbf{A}_{-k}^{-1} is
714 $(1 \pm O(n^{-\kappa_2}))n^{-p}$ with an appropriately defined positive constant κ_2 . The spectral norm bound
715 follows by defining c_2 to be any positive constant greater than 1 which absorbs the $o(1)$ deviation
716 terms in the spectrum.

717 The proof concludes by setting $\kappa_3 = \min \{\kappa_1, \kappa_2\}$. \square

718 Finally, we have the following proposition which controls the spectrum of hat matrices such as
719 $\mathbf{H}_k \triangleq \mathbf{W}_k^\top \mathbf{A}_{-k}^{-1} \mathbf{W}_k \in \mathbb{R}^{k \times k}$. The intuition is that even though the spectrum of \mathbf{A}_{-k}^{-1} may be spiked,
720 the spectrum of $\mathbf{W}_k^\top \mathbf{A}_{-k}^{-1} \mathbf{W}_k$ is ultimately flat because we are taking an extremely low dimensional
721 projection which is unlikely to see significant contribution from the spiked portion of \mathbf{A}_{-k}^{-1} .

722 In fact, we can prove a more general statement, which will be useful for us in the proof. Let
723 $\emptyset \neq T \subseteq S \subseteq [s]$; here T and S index nonempty subsets of the s favored features. Then
724 we can define \mathbf{W}_T to be the matrix of weighted features in T and the leave- T -out Gram matrix
725 $\mathbf{A}_{-T} \triangleq \sum_{j \notin T} \lambda_j \mathbf{z}_j \mathbf{z}_j^\top$. Now define the (T, S) hat matrix as $\mathbf{H}_{T,S} \triangleq \mathbf{W}_T^\top \mathbf{A}_{-S}^{-1} \mathbf{W}_T$. Evidently we
726 have $\mathbf{H}_k = \mathbf{H}_{[k],[k]}$, so our notion is more general. The full proof is deferred to Appendix H.

727 **Proposition B.6** (Generalized hat matrices are flat). *Assume we are in the bi-level ensemble Defini-*
728 *tion 1. For any nonempty $T \subseteq S \subseteq [s]$, with probability at least $1 - 2e^{-\sqrt{n}} - 2e^{-n}$, we have all the*
729 *eigenvalues tightly controlled:*

$$\mu_i((\mathbf{I}_{|T|} + \mathbf{H}_{T,S})^{-1}) = \min \{\mu, 1\} (1 \pm c_{T,S} n^{-\kappa_{11}}). \quad (74)$$

730 where $c_{T,S}$ and κ_{11} are positive constants that depend on $|T|$ and $|S|$.

731 C Utility bounds: applying the tools

732 Wishart concentration allows us to tightly bound the hat matrix and pass to studying bilinear forms of
733 the form $\mathbf{w}_i^\top \mathbf{A}_{-k}^{-1} \Delta y$ rather than $\mathbf{w}_i^\top \mathbf{A}^{-1} \Delta y$. Since \mathbf{A}_{-k}^{-1} is independent of \mathbf{W}_k and Δy , we can
734 condition on \mathbf{A}_{-k}^{-1} and then apply Hanson-Wright (Theorem 4.1) to these bilinear forms for every
735 realization of \mathbf{A}_{-k}^{-1} . In this section, we will explicitly calculate the scaling of the typical value of
736 these bilinear forms using the bi-level ensemble scaling; these will prove to be useful throughout the
737 rest of the paper.

738 We first state the following proposition which bounds the correlation between the relevant label-
739 defining features and the label vectors; it is a combination of Propositions D.5 and D.6 in (Subrama-
740 nian et al., 2022).

741 **Proposition C.1.** *For any distinct $\alpha, \beta \in [k]$, we have*

$$\frac{1}{\sqrt{\pi \ln 2}} \cdot \frac{n}{k} \cdot \sqrt{\ln k} \leq \mathbb{E}[\mathbf{z}_\alpha^\top \mathbf{y}_\alpha] \leq \sqrt{2} \cdot \frac{n}{k} \cdot \sqrt{\ln k} \quad (75)$$

742 and

$$-\sqrt{2} \cdot \frac{n}{k} \cdot \frac{1}{k-1} \cdot \sqrt{\ln k} \leq \mathbb{E}[\mathbf{z}_\alpha^\top \mathbf{y}_\beta] \leq -\frac{1}{\sqrt{\pi \ln 2}} \cdot \frac{n}{k} \cdot \frac{1}{k-1} \cdot \sqrt{\ln k} \quad (76)$$

743 With the above proposition in hand, we can prove the following lemma which gives concentration of
744 the bilinear forms that we study.

745 **Lemma C.2.** *Let $i \in [d]$ and $\Delta y = \mathbf{y}_\alpha - \mathbf{y}_\beta$ where $\alpha, \beta \in [k]$ and $\beta \neq \alpha$. Let $\mathbf{M} \in \mathbb{R}^{n \times n}$ be a
746 (random) matrix which is independent of \mathbf{z}_i and Δy . Then conditioned on \mathbf{M} , with probability at
747 least $1 - 1/nk$,*

$$|\mathbf{z}_i^\top \mathbf{M} \Delta y - \mathbb{E}[\mathbf{z}_i^\top \mathbf{M} \Delta y | \mathbf{M}]| \leq c_3 \sqrt{\frac{n}{k}} \|\mathbf{M}\|_2 \sqrt{\log(nk)},$$

748 and the same holds with Δy replaced with \mathbf{y}_α . Here, c_3 is an appropriately chosen universal positive
749 constant.

750 Moreover, we have

751 (1) For any distinct $\alpha, \beta \in [k]$, we have

$$\mathbb{E}[\mathbf{z}_\alpha^\top \mathbf{M} \Delta y | \mathbf{M}] = c_7 \frac{\sqrt{\log k}}{k} \text{tr}(\mathbf{M}) = -\mathbb{E}[\mathbf{z}_\beta^\top \mathbf{M} \Delta y] \quad (77)$$

$$\mathbb{E}[\mathbf{z}_\alpha^\top \mathbf{M} \mathbf{y}_\alpha | \mathbf{M}] = c_4 \frac{\sqrt{\log k}}{k} \text{tr}(\mathbf{M}), \quad (78)$$

752 where c_7 and c_4 are positive constants.

753 (2) For $i \in [d] \setminus \{\alpha, \beta\}$, we have

$$\mathbb{E}[\mathbf{z}_i^\top \mathbf{M} \Delta y | \mathbf{M}] = 0. \quad (79)$$

754 (3) For $i \in [d] \setminus \{\alpha\}$, we have

$$\mathbb{E}[\mathbf{z}_i^\top \mathbf{M} \mathbf{y}_\alpha | \mathbf{M}] = -c_5 \frac{\sqrt{\log k}}{k(k-1)}, \quad (80)$$

755 where c_5 is a positive constant.

756 *Proof.* Let us check the conditions for our new variant of Hanson-Wright with soft sparsity (The-
757 orem 4.1). We want to apply it to the random vectors $(\mathbf{z}_i, \Delta y) = (\mathbf{z}_i[j], \Delta y[j])_{j=1}^n$. Some of the
758 hypotheses are immediate by definition. Evidently, $(\mathbf{z}_i[j], \Delta y[j])$ are independent across j , and are
759 mean zero. Since $\mathbf{z}_i[j] \sim N(0, 1)$, it is subgaussian with parameter at most $K = 2$. For the bounded
760 and soft sparsity assumption, we clearly have $|\Delta y[j]| \leq 1$ and $\mathbf{y}_\alpha[j] \leq 1$ almost surely. Also, since
761 $\Delta y[j]^2 \sim \text{Ber}(\frac{2}{k})$, we have $\mathbb{E}[\Delta y[j]^2] = \frac{2}{k}$. Similarly, $\mathbb{E}[\mathbf{y}_\alpha[j]^2] = \frac{1}{k}(1 - \frac{1}{k})^2 + (1 - \frac{1}{k})\frac{1}{k^2} \leq \frac{2}{k}$.

762 The more complicated condition is the subgaussianity of $\mathbf{z}_i[j]$ conditioned on the value of $\Delta y[j]$ or
763 $\mathbf{y}_\alpha[j]$. Regardless of whether we're conditioning on Δy or \mathbf{y}_α , it suffices to instead prove that $\mathbf{z}_i[j]$
764 is subgaussian conditioned on whether feature i won the competition for datapoint j . First, suppose
765 i won, i.e. $\mathbf{y}_i^{\text{oh}}[j] = 1$. Then the Borell-TIS inequality (Adler et al., 2007, Theorem 2.1.1) implies
766 that $\mathbf{z}_i[j]$ satisfies a subgaussian tail inequality. By the equivalent conditions for subgaussianity
767 Vershynin (2018, Proposition 2.5.2), it follows that $\mathbf{z}_i[j] - \mathbb{E}[\mathbf{z}_i[j] | \mathbf{y}_i^{\text{oh}}[j] = 1]$ has conditionally has
768 subgaussian norm bounded by some absolute constant K . If i doesn't win (or doesn't participate in the
769 competition), then Proposition D.2 in Subramanian et al. (2022) implies that $\mathbf{z}_i[j] - \mathbb{E}[\mathbf{z}_i[j] | \mathbf{y}_i^{\text{oh}}[j] =$
770 $0]$ conditionally has subgaussian norm bounded by 6.

771 Finally, since \mathbf{M} is independent of \mathbf{z}_i and Δy , we can condition on \mathbf{M} and apply Theorem 4.1 to
772 the bilinear form for every realization of \mathbf{M} .

773 Hence we conclude that with probability at least $1 - 1/nk$ we have

$$|\mathbf{z}_i^\top \mathbf{M} \Delta y - \mathbb{E}[\mathbf{z}_i^\top \mathbf{M} \Delta y | \mathbf{M}]| \leq c_3 \sqrt{\frac{n}{k}} \|\mathbf{M}\|_2 \sqrt{\log(nk)}, \quad (81)$$

774 where c_3 is an appropriately chosen absolute constant based on K and the constant c defined in
775 Theorem 4.1.

776 Now we can compute $\mathbb{E}[\mathbf{z}_i^\top \mathbf{M} \Delta y | \mathbf{M}]$ to prove the rest of the theorem. If $i = \alpha$, we have

$$\mathbb{E}[\mathbf{z}_\alpha^\top \mathbf{M} \Delta y | \mathbf{M}] = \text{tr}(\mathbf{M} \mathbb{E}[\Delta y \mathbf{z}_\alpha^\top]).$$

Let us now compute $\mathbb{E}[\Delta y z_\alpha^\top]$. From Eq. (75) in Proposition C.1, we have $\mathbb{E}[y_\alpha z_\alpha^\top] = c_4 \frac{\sqrt{\log k}}{k} \mathbf{I}_n$, where $\frac{1}{\sqrt{\pi \log 2}} \leq c_4 \leq \sqrt{2}$. Similarly we have $\mathbb{E}[y_\beta z_\alpha^\top] = -c_5 \frac{\sqrt{\log k}}{k(k-1)} \mathbf{I}_n$ where $\frac{1}{\sqrt{\pi \log 2}} \leq c_5 \leq \sqrt{2}$. It follows that $\mathbb{E}[\Delta y z_\alpha^\top] = \Theta\left(\frac{\sqrt{\log k}}{k}\right) \mathbf{I}_n$.

For $i \in [d] \setminus \{\alpha, \beta\}$, by symmetry we obtain $\mathbb{E}[y_\alpha z_i^\top] = \mathbb{E}[y_\beta z_i^\top]$. This implies $\mathbb{E}[\Delta y z_i^\top] = \mathbb{E}[y_\alpha z_i^\top] - \mathbb{E}[y_\beta z_i^\top] = 0$, so we obtain

$$\mathbb{E}[z_i^\top M \Delta y | M] = \text{tr}(M \mathbb{E}[\Delta y z_i^\top]) \quad (82)$$

$$= 0. \quad (83)$$

782

□

Plugging in the bi-level scaling, we obtain the following corollary.

Corollary C.3 (Asymptotic concentration of bilinear forms). *In the bi-level model, for any $i \in [k]$, we have with probability at least $1 - O(1/nk)$ that*

$$|z_i^\top A_{-k}^{-1} \Delta y - \mathbb{E}[z_i^\top A_{-k}^{-1} \Delta y]| \leq c_6 n^{\frac{1-t}{2}-p} \sqrt{\log(nk)}.$$

Moreover, we have

(1) For any distinct $\alpha, \beta \in [k]$,

$$\mathbb{E}[z_\alpha^\top A_{-k}^{-1} \Delta y] = c_7 n^{1-t-p} (1 \pm O(n^{-\kappa_3})) \sqrt{\log k} = -\mathbb{E}[z_\beta^\top A_{-k}^{-1} \Delta y] \quad (84)$$

The same statements hold (with different constants) if we replace A_{-k}^{-1} with A_{-s}^{-1} .

Proof. From Corollary B.5, we have $\|A_{-k}^{-1}\|_2 \leq c_6 n^{-p}$, where c_6 is an appropriately chosen universal positive constant based on c_3 . Recall that A_{-k} is obtained by removing the k label-defining features, so in particular A_{-k}^{-1} is independent of $(z_i, \Delta y)$ for $i \in [k]$. Hence, the conditions for Lemma C.2 are satisfied. Then applying the union bound for the spectral norm bound on A_{-k}^{-1} , we see that with probability at least $1 - O(1/nk)$, the deviation term from Hanson-Wright is at most $c_6 n^{\frac{1}{2}-p} \sqrt{\log(nk)}$.

We now turn to calculating the asymptotic scalings for the expectations. From Lemma C.2, we know that $\mathbb{E}[z_\alpha^\top A_{-k}^{-1} \Delta y | A_{-k}^{-1}] = c_7 \frac{\sqrt{\log k}}{k} \text{tr}(A_{-k}^{-1})$. Applying the high probability bound on $\text{tr}(A_{-k}^{-1})$ from Corollary B.5, we obtain that with probability at least $1 - O(1/nk)$ that

$$-\mathbb{E}[z_\beta^\top A_{-k}^{-1} \Delta y] = \mathbb{E}[z_\alpha^\top A_{-k}^{-1} \Delta y] = c_7 n^{-t} n^{1-p} (1 \pm O(n^{-\kappa_3})) \sqrt{\log k} \quad (85)$$

$$= c_7 n^{1-t-p} (1 \pm O(n^{-\kappa_3})) \sqrt{\log k} \quad (86)$$

where in the second line we have applied Corollary B.5 and c_7 is an appropriately chosen positive constant. This proves (84). □

With Corollary C.3 in hand, we are now in a position to do some straightforward calculations and bound some quantities which will pop up in the survival and contamination analysis.

Proposition C.4 (Worst-case bound based on Hanson-Wright). *Let $T \subseteq [s]$ be a subset of favored features such that $\{\alpha, \beta\} \subseteq T$. Assume that $|T| = n^\tau$ for some $\tau \leq r$. Then with probability at least $1 - O(1/nk)$, we have*

$$\|Z_T^\top A_{-T}^{-1} \Delta y\|_2 \leq c_8 (n^{1-t-p} + n^{\frac{1+\tau-t}{2}-p}) \sqrt{\log(nk|T|)}. \quad (87)$$

Proof. WLOG, suppose $\alpha = 1$ and $\beta = 2$. By Corollary C.3 we have with probability at least $1 - 1/n$ that

$$|Z_T^\top A_{-T}^{-1} \Delta y| \leq \begin{bmatrix} c_7 n^{1-t-p} \sqrt{\log k} \\ c_7 n^{1-t-p} \sqrt{\log k} \\ c_6 n^{\frac{1-t}{2}-p} \sqrt{\log(nk)} \\ \vdots \\ c_6 n^{\frac{1-t}{2}-p} \sqrt{\log(nk)} \end{bmatrix}. \quad (88)$$

807 Hence the norm of this vector is at most

$$\|Z_T^\top A^{-1} \Delta y\|_2 \leq 2c_7 n^{1-t-p} \sqrt{\log k} + n^{\frac{\tau}{2}} c_6 n^{\frac{1-t}{2}-p} \sqrt{\log(nk)} \quad (89)$$

$$\leq c_8 (n^{1-t-p} + n^{\frac{1+\tau-t}{2}-p}) \sqrt{\log(nk)}, \quad (90)$$

808 where c_8 is a positive constant.

809

□

810 D Bounding the survival

811 Recall that the relative survival was defined to be $\lambda_F \hat{\mathbf{h}}_{\alpha,\beta}[\alpha] = \lambda_F \mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y$. The strategy is to
 812 apply our variant of Hanson-Wright to $\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y$. Unfortunately, \mathbf{A}^{-1} is not independent of \mathbf{z}_α
 813 or Δy , so we need to use Woodbury to extract out the independent portions and bound away the
 814 dependent portion. As we'll see shortly, the error from the dependent portions can also be controlled
 815 using Hanson-Wright. Let us now recall Proposition A.1 for reference.

816 **Proposition A.1** (Bounds on relative survival). *Under the bi-level ensemble model (Definition 1),
 817 when the true data generating process is 1-sparse (Assumption 1), if $t < \frac{1}{2}$, then with probability at
 818 least $1 - O(1/nk)$*

$$\lambda_F \hat{\mathbf{h}}_{\alpha,\beta}[\alpha] = c_7 \min\{\mu^{-1}, 1\} n^{-t} (1 \pm O(n^{-\kappa_5})) \sqrt{\log k},$$

819 where c_7 and κ_5 are positive constants.

820 If $t \geq \frac{1}{2}$, then

$$\lambda_F |\hat{\mathbf{h}}_{\alpha,\beta}[\alpha]| \leq c_9 \min\{\mu^{-1}, 1\} n^{-\frac{1}{2}} \sqrt{\log(nk)},$$

821 where c_9 is a positive constant.

822 *Proof.* Recall that $\hat{\mathbf{h}}_{\alpha,\beta}[\alpha] = \mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y$. We first observe that for $i \in [k]$, the dependence between
 823 \mathbf{A}^{-1} and \mathbf{z}_i as well as \mathbf{A}^{-1} and Δy only comes through the k label defining features. Hence, we
 824 can use the Woodbury identity to extract out the independent portions of \mathbf{A}^{-1} .

825 Indeed, our “push through” lemma for Woodbury (Lemma B.2) and concentration of the hat matrix
 826 (Proposition B.6) implies that with extremely high probability

$$\mathbf{Z}_k^\top \mathbf{A}^{-1} \Delta y = (\mathbf{I}_k + \mathbf{H}_k)^{-1} \mathbf{Z}_k^\top \mathbf{A}_{-k}^{-1} \Delta y \quad (91)$$

$$= \min\{\mu, 1\} (\mathbf{I}_k + \mathbf{E}) \mathbf{Z}_k^\top \mathbf{A}_{-k}^{-1} \Delta y, \quad (92)$$

827 where $\|\mathbf{E}\|_2 = O(n^{-\kappa_{11}})$.

828 Let $\mathbf{u}_\alpha \in \mathbb{R}^k$ denote the α th row vector in \mathbf{E} , and let $\mathbf{u}_\alpha^- \in \mathbb{R}^{k-1}$ denote the subvector of \mathbf{u}_α without
 829 index α . By reading off the α th row of Eq. (92), we see that

$$\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y = \min\{\mu, 1\} (\mathbf{z}_\alpha^\top \mathbf{A}_{-k}^{-1} \Delta y + \langle \mathbf{u}_\alpha, \mathbf{Z}_k^\top \mathbf{A}_{-k}^{-1} \Delta y \rangle) \quad (93)$$

830 Since $\|\mathbf{u}_\alpha\|_2 \leq \|\mathbf{E}\|_2 = O(n^{-\kappa_{11}})$, it follows from Cauchy-Schwarz that

$$|\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y - \min\{\mu, 1\} \mathbf{z}_\alpha^\top \mathbf{A}_{-k}^{-1} \Delta y| \leq \min\{\mu, 1\} O(n^{-\kappa_{11}}) \|\mathbf{Z}_k^\top \mathbf{A}_{-k}^{-1} \Delta y\|_2. \quad (94)$$

831 Let us pause for a moment and interpret Eq. (94). The term $\min\{\mu, 1\}$ is merely capturing the
 832 difference in behavior when regression works and fails; if regression works ($q + r < 1$) then it
 833 becomes μ , and if regression fails ($q + r > 1$), then it becomes 1. This behavior should be expected:
 834 in the regression works case, we expect the effect of interpolation to be a *regularizing* one: the signals
 835 are attenuated by a factor of μ . The RHS of Eq. (94) is an error term, capturing how differently
 836 $\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y$ behaves from the expected behavior $\min\{\mu, 1\} \mathbf{z}_\alpha^\top \mathbf{A}_{-k}^{-1} \Delta y$.

837 Let us now bound the error term. From Proposition C.4 we have with probability at least $1 - O(1/nk)$
 838 that

$$\|\mathbf{Z}_k^\top \mathbf{A}_{-k}^{-1} \Delta y\|_2 \leq c_8 (n^{1-t-p} + n^{\frac{1}{2}-p}) \sqrt{\log(nk^2)}. \quad (95)$$

839 Let us do casework on t . For $t < \frac{1}{2}$, we have $\frac{1}{2} - p < 1 - t - p$, so we conclude that the error term
 840 is $\min\{\mu, 1\}O(n^{1-t-p} \cdot n^{-\kappa_4})\sqrt{\log(nk^2)}$, where κ_4 is a positive constant.

841 On the other hand, our Hanson-Wright calculations imply (Corollary C.3) that with probability at
 842 least $1 - O(1/nk)$ that

$$\left| \mathbf{z}_\alpha^\top \mathbf{A}_{-k}^{-1} \Delta y - c_7 n^{1-t-p} (1 \pm O(n^{-\kappa_3})) \sqrt{\log k} \right| \leq c_6 n^{\frac{1}{2}-p} \sqrt{\log(nk)}.$$

843 Again, since $t < \frac{1}{2}$, the deviation term is $o(n^{1-t-p})\sqrt{\log k}$.

844 Hence we conclude that with probability $1 - O(1/nk)$ we have

$$\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y = c_7 \min\{\mu, 1\} n^{1-t-p} (1 \pm O(n^{-\kappa_5})) \sqrt{\log k},$$

845 where κ_5 is a positive constant.

846 Completely analogous logic handles the bounds for $\mathbf{z}_\beta^\top \mathbf{A}^{-1} \Delta y$. Let us now return back to the
 847 quantity of interest, $\lambda_F \hat{\mathbf{h}}_{\alpha,\beta}[\alpha]$. We can compute

$$\begin{aligned} \lambda_F \mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y &= n^{p-q-r} \cdot c_7 \min\{\mu, 1\} n^{1-t-p} (1 \pm O(n^{-\kappa_5})) \sqrt{\log k} \\ &= c_7 \mu^{-1} \min\{\mu, 1\} n^{-t} (1 \pm O(n^{-\kappa_5})) \sqrt{\log k} \\ &= c_7 \min\{1, \mu^{-1}\} n^{-t} (1 \pm O(n^{-\kappa_5})) \sqrt{\log k}. \end{aligned}$$

848 On the other hand, if $t \geq \frac{1}{2}$ the error terms all dominate, and we replace n^{1-t-p} with $n^{\frac{1}{2}-p}$
 849 everywhere. We conclude that with probability at least $1 - O(1/nk)$,

$$|\mathbf{z}_\alpha^\top \mathbf{A}^{-1} \Delta y| \leq c_9 \min\{\mu, 1\} n^{\frac{1}{2}-p} \sqrt{\log(nk)}, \quad (96)$$

850 where c_9 is a positive constant. Plugging in the scaling for λ_F yields the desired result. \square

851 E Bounding the contamination

852 In this section we give a tight analysis of the contamination term. First, we rewrite the squared
 853 contamination term and separate it out into the contamination from the $k - 2$ label-defining features
 854 which are not α or β , the rest of the $s - k$ favored features, and the remaining $d - s$ unfavored
 855 features. From Eq. (30), we have

$$\text{CN}_{\alpha,\beta}^2 = \sum_{j \in [d] \setminus \{\alpha,\beta\}} \lambda_j^2 (\mathbf{z}_j^\top \mathbf{A}^{-1} \Delta y)^2 \quad (97)$$

$$= \Delta y^\top \mathbf{A}^{-1} \left(\sum_{j \in [d] \setminus \{\alpha,\beta\}} \lambda_j^2 \mathbf{z}_j \mathbf{z}_j^\top \right) \mathbf{A}^{-1} \Delta y \quad (98)$$

$$\begin{aligned} &= \underbrace{\Delta y^\top \mathbf{A}^{-1} \left(\sum_{j \in [k] \setminus \{\alpha,\beta\}} \lambda_j^2 \mathbf{z}_j \mathbf{z}_j^\top \right) \mathbf{A}^{-1} \Delta y}_{\triangleq \text{CN}_{\alpha,\beta,L}^2} + \underbrace{\Delta y^\top \mathbf{A}^{-1} \left(\sum_{j \in [s] \setminus [k]} \lambda_j^2 \mathbf{z}_j \mathbf{z}_j^\top \right) \mathbf{A}^{-1} \Delta y}_{\triangleq \text{CN}_{\alpha,\beta,F}^2} \\ &\quad + \underbrace{\Delta y^\top \mathbf{A}^{-1} \left(\sum_{j > s} \lambda_j^2 \mathbf{z}_j \mathbf{z}_j^\top \right) \mathbf{A}^{-1} \Delta y}_{\triangleq \text{CN}_{\alpha,\beta,U}^2}. \end{aligned} \quad (99)$$

$$\quad (100)$$

856 Here, $\text{CN}_{\alpha,\beta,L}$ corresponds to contamination from label defining features, $\text{CN}_{\alpha,\beta,F}$ corresponds to
 857 contamination from favored features, and $\text{CN}_{\alpha,\beta,U}$ corresponds to contamination from unfavored
 858 features. The reason for separating out the contamination into these three subterms is that we will
 859 need slightly different arguments to bound each of them, although Hanson-Wright and Woodbury are

central to all of the arguments. In Appendix E.1 we prove the upper bound on $\text{CN}_{\alpha,\beta,L} + \text{CN}_{\alpha,\beta,F}$; in Appendix E.2 we prove the lower bound. Finally, in Appendix E.3 we bound $\text{CN}_{\alpha,\beta,U}$. After putting these bounds together, we will obtain the main bounds on the contamination, which we restate here for reference.

Proposition A.2 (Bounds on contamination). *Under the bi-level ensemble model (Definition 1), when the true data generating process is 1-sparse (Assumption 1), with probability at least $1 - O(1/nk)$,*

$$\text{CN}_{\alpha,\beta} \leq \underbrace{\min\{\mu^{-1}, 1\} O(n^{\frac{r-t-1}{2}}) \log(nsk)}_{\text{favored features}} + \underbrace{O(n^{\frac{1-t-p}{2}}) \sqrt{\log(nsk)}}_{\text{unfavored features}}.$$

Furthermore, if $t > 0$, then with probability at least $1 - O(1/nk)$,

$$\text{CN}_{\alpha,\beta} \geq \underbrace{\min\{\mu^{-1}, 1\} \Omega(n^{\frac{r-t-1}{2}})}_{\text{favored features}} + \underbrace{\Omega(n^{\frac{1-t-p}{2}})}_{\text{unfavored features}}.$$

E.1 Upper bounding the contamination from label-defining+favored features

In this section, we upper bound the contamination coming from the $s - 2$ favored features which are not α or β . This culminates in the following lemma.

Lemma E.1. *In the same setting as Proposition A.2, we have with probability $1 - O(1/nk)$ that*

$$\text{CN}_{\alpha,\beta,L}^2 + \text{CN}_{\alpha,\beta,F}^2 \leq c_{12}^2 \min\{1, \mu^{-2}\} n^{r-t-1} \log(nsk)^2,$$

where c_{12} is a positive constant.

Proof. Let $\mathbf{W}_R \in \mathbb{R}^{n \times (s-2)}$ as the weighted feature matrix which includes all of the $s - 2$ favored features aside from α, β . We can then define $\mathbf{A}_{-R} \triangleq \mathbf{A} - \mathbf{W}_R \mathbf{W}_R^\top$ and $\mathbf{H}_R = \mathbf{W}_R^\top \mathbf{A}_{-R}^{-1} \mathbf{W}_R$. Using Woodbury, an analogous computation to Lemma B.2 implies that

$$\mathbf{W}_R \mathbf{A}^{-1} \Delta y = (\mathbf{I}_{s-2} + \mathbf{H}_R)^{-1} \mathbf{W}_R \mathbf{A}_{-R}^{-1} \Delta y. \quad (101)$$

The contamination from all of the $s - 2$ favored features that are not α or β satisfies

$$\begin{aligned} \text{CN}_{\alpha,\beta,L}^2 + \text{CN}_{\alpha,\beta,F}^2 &= \lambda_F \Delta y^\top \mathbf{A}^{-1} \sum_{j \in [s] \setminus \{\alpha, \beta\}} \mathbf{w}_j \mathbf{w}_j^\top \mathbf{A}^{-1} \Delta y \\ &= \lambda_F \Delta y^\top \mathbf{A}^{-1} \mathbf{W}_R \mathbf{W}_R^\top \mathbf{A}^{-1} \Delta y \\ &= \lambda_F \Delta y^\top \mathbf{A}_{-R}^{-1} \mathbf{W}_R (\mathbf{I}_{s-2} + \mathbf{H}_R)^{-2} \mathbf{W}_R^\top \mathbf{A}_{-R}^{-1} \Delta y. \end{aligned}$$

Since Proposition B.6 implies that $\mu_1((\mathbf{I}_{s-2} + \mathbf{H}_R)^{-2}) \leq c_{10} \min\{\mu^2, 1\}$ with extremely high probability where c_{10} is a positive constant, we know the contamination is with extremely high probability upper bounded by the following quadratic form:

$$c_{10} \min\{\mu^2, 1\} \lambda_F \Delta y^\top \mathbf{A}_{-R}^{-1} \mathbf{W}_R \mathbf{W}_R^\top \mathbf{A}_{-R}^{-1} \Delta y \quad (102)$$

$$= c_{10} \min\{\mu^2, 1\} \lambda_F^2 \sum_{j \in [s] \setminus \{\alpha, \beta\}} \langle \mathbf{z}_j, \mathbf{A}_{-R}^{-1} \Delta y \rangle^2. \quad (103)$$

We still cannot apply Hanson-Wright, because \mathbf{A}_{-R}^{-1} is not independent of Δy . However, we can use Woodbury again to take out $\mathbf{z}_\alpha, \mathbf{z}_\beta$ from \mathbf{A}_{-R}^{-1} .

Define $\mathbf{W}_{\alpha,\beta} = [\mathbf{w}_\alpha \quad \mathbf{w}_\beta]$ and $\mathbf{H}_{\alpha,\beta}^{(s)} = \mathbf{W}_{\alpha,\beta}^\top \mathbf{A}_{-s}^{-1} \mathbf{W}_{\alpha,\beta}$. Then Woodbury implies that

$$\mathbf{A}_{-R}^{-1} = \mathbf{A}_{-s}^{-1} - \mathbf{A}_{-s}^{-1} \mathbf{W}_{\alpha,\beta} (\mathbf{I}_2 + \mathbf{H}_{\alpha,\beta}^{(s)})^{-1} \mathbf{W}_{\alpha,\beta}^\top \mathbf{A}_{-s}^{-1}. \quad (104)$$

Hence

$$\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \Delta y = \mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \Delta y - \mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{W}_{\alpha,\beta} (\mathbf{I}_2 + \mathbf{H}_{\alpha,\beta}^{(s)})^{-1} \mathbf{W}_{\alpha,\beta}^\top \mathbf{A}_{-s}^{-1} \Delta y. \quad (105)$$

We will use Hanson-Wright and Cauchy-Schwarz to argue that the second term in Eq. (105) above will be dominated by the first term. Indeed, Corollary C.3 implies that $\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \Delta y \leq$

885 $c_6 n^{\frac{1-t}{2}-p} \sqrt{\log(nsk)}$ with probability at least $1 - O(1/nsk)$, so it suffices to show that the other term
 886 is dominated by $n^{\frac{1-t}{2}-p}$. We will show that its contribution for each j is $\min\{1, \mu^{-1}\} \tilde{O}(n^{\frac{1}{2}-t-p})$.
 887 By Cauchy-Schwarz, the magnitude of the second term of Eq. (105) is at most

$$\|(\mathbf{I}_2 + \mathbf{H}_{\alpha,\beta})^{-1}\|_2 \|\mathbf{W}_{\alpha,\beta}^\top \mathbf{A}_{-s}^{-1} \mathbf{z}_j\|_2 \|\mathbf{W}_{\alpha,\beta}^\top \mathbf{A}_{-s}^{-1} \Delta y\|_2 \quad (106)$$

$$\leq c_{11} \lambda_F \min\{\mu, 1\} n^{\frac{1}{2}-p} \sqrt{\log(nsk)} \cdot n^{1-t-p} \sqrt{\log k} \quad (107)$$

$$\leq c_{11} \min\{1, \mu^{-1}\} n^{\frac{1}{2}-t-p} \log(nsk), \quad (108)$$

888 where c_{11} is a positive constant. In the second line, we have used Proposition B.6 to upper bound
 889 $\|(\mathbf{I}_2 + \mathbf{H}_{\alpha,\beta}^{(s)})^{-1}\|_2 \leq O(\min\{\mu, 1\})$ and we have used Theorem B.1 and Proposition B.4 to deduce
 890 that that $|\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{z}_a| \leq O(n^{\frac{1}{2}-p}) \sqrt{\log(nsk)}$ with probability at least $1 - O(1/nsk)$. Similarly, we
 891 used the scaling from Corollary C.3 to deduce that $|\mathbf{z}_\alpha^\top \mathbf{A}_{-s}^{-1} \Delta y| \leq O(n^{1-t-p}) \sqrt{\log k}$, and similarly
 892 for β .

893 Hence $\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \Delta y$ is $O(n^{\frac{1-t}{2}-p}) \log(nsk)$ with probability $1 - O(1/nsk)$. By union bounding
 894 over j and plugging our upper bound back into Eq. (103), we conclude that with probability at least
 895 $1 - O(1/nk)$

$$\text{CN}_{\alpha,\beta,L}^2 + \text{CN}_{\alpha,\beta,F}^2 \leq c_{10} \min\{\mu^2, 1\} \lambda_F^2 n^r \cdot O(n^{1-t-2p}) \log(nsk)^2 \quad (109)$$

$$= \mu^{-2} \min\{\mu^2, 1\} O(n^{r-t-1}) \log(nsk)^2 \quad (110)$$

$$\leq c_{12}^2 \min\{1, \mu^{-2}\} n^{r-t-1} \log(nsk)^2, \quad (111)$$

896 where c_{12} is a positive constant, concluding the proof. \square

897 E.2 Lower bounding the contamination from label-defining+ favored features

898 In this section, we upper bound the contamination coming from the $s - 2$ favored features which are
 899 not α or β . This culminates in the following lemma.

900 **Lemma E.2.** *In the same setting as Proposition A.2, if $t > 0$, with probability at least $1 - O(1/nk)$,
 901 we have*

$$\text{CN}_{\alpha,\beta,L}^2 + \text{CN}_{\alpha,\beta,F}^2 \geq c_{14}^2 \min\{1, \mu^{-2}\} n^{r-t-1},$$

902 where c_{14} is a positive constant.

903 *Proof.* Following the beginning of the proof of Lemma E.1 and what we know about the flatness of
 904 the spectra of hat matrices from Proposition B.6, we can deduce that there is some positive constant
 905 c_{13} such that with extremely high probability

$$\text{CN}_{\alpha,\beta,L}^2 + \text{CN}_{\alpha,\beta,F}^2 \geq c_{13} \min\{\mu^2, 1\} \lambda_F^2 \sum_{j \in [s] \setminus \{\alpha, \beta\}} \langle \mathbf{z}_j, \mathbf{A}_{-R}^{-1} \Delta y \rangle^2. \quad (112)$$

906 We will further lower bound this by throwing out all label-defining j . In other words, the goal now is
 907 to lower bound

$$\sum_{j \in [s] \setminus [k]} \langle \mathbf{z}_j, \mathbf{A}_{-R}^{-1} \Delta y \rangle^2 = \langle \mathbf{Z}_F^\top \mathbf{A}_{-R}^{-1} \Delta y, \mathbf{Z}_F^\top \mathbf{A}_{-R}^{-1} \Delta y \rangle. \quad (113)$$

908 The main idea is to use Bernstein's inequality, but unfortunately \mathbf{A}_{-R}^{-1} is not independent of Δy , so
 909 we will again resort to Woodbury to take out \mathbf{z}_α and \mathbf{z}_β . As in the proof for upper bounding the
 910 favored contamination, we have $\mathbf{W}_{\alpha,\beta} = [\mathbf{w}_\alpha \quad \mathbf{w}_\beta]$ and $\mathbf{H}_{\alpha,\beta}^{(s)} = \mathbf{W}_{\alpha,\beta}^\top \mathbf{A}_{-s}^{-1} \mathbf{W}_{\alpha,\beta}$. Then we can
 911 deduce from another application of Woodbury that

$$\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \Delta y = \mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \Delta y - \mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{W}_{\alpha,\beta} (\mathbf{I}_2 + \mathbf{H}_{\alpha,\beta})^{-1} \mathbf{W}_{\alpha,\beta}^\top \mathbf{A}_{-s}^{-1} \Delta y. \quad (114)$$

912 Again, we can argue that with probability $1 - O(1/nsk)$, the second term is upper bounded in
 913 magnitude by

$$\min\{1, \mu^{-1}\} n^{\frac{1}{2}-t-p} \log(nsk) = \min\{1, \mu^{-1}\} O(n^{\frac{1-t}{2}-p} \cdot n^{-\kappa_6}), \quad (115)$$

where κ_6 is a positive constant because $t > 0$. Since Hanson-Wright (Corollary C.3) implies that $\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \Delta y = \tilde{O}(n^{\frac{1-t}{2}-p})$, this implies that $\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \Delta y = \tilde{O}(n^{\frac{1-t}{2}-p})$, and similarly for β . Hence we have

$$(\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \Delta y)^2 = (\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \Delta y + \min\{1, \mu^{-1}\} O(n^{\frac{1-t}{2}-p-\kappa_6}))^2 \quad (116)$$

$$= (\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \Delta y)^2 + \min\{1, \mu^{-1}\} O(n^{\frac{1-t}{2}-p-\kappa_6}) \tilde{O}(n^{\frac{1-t}{2}-p}) \quad (117)$$

$$= (\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \Delta y)^2 + \min\{1, \mu^{-1}\} o(n^{1-t-2p}). \quad (118)$$

We are now in a position to analyze the contribution from the first term of Eq. (118) to Eq. (113): its contribution is $\langle \mathbf{Z}_F^\top \mathbf{A}_{-s}^{-1} \Delta y, \mathbf{Z}_F^\top \mathbf{A}_{-s}^{-1} \Delta y \rangle$. This does have all the independence required to apply Bernstein, because $(\mathbf{A}_{-s}^{-1}, \Delta y)$ are independent of \mathbf{Z}_F . Hence conditioned on \mathbf{A}_{-s}^{-1} and Δy , $\langle \mathbf{Z}_F^\top \mathbf{A}_{-s}^{-1} \Delta y, \mathbf{Z}_F^\top \mathbf{A}_{-s}^{-1} \Delta y \rangle$ is a sum of $s - k$ subexponential variables, and by Lemma 2.7.7 of Vershynin (2018) each of these random variables conditionally has subexponential norm at most $\|\mathbf{A}_{-s}^{-1} \Delta y\|_2^2$ and conditional mean $\langle \mathbf{A}_{-s}^{-1} \Delta y, \mathbf{A}_{-s}^{-1} \Delta y \rangle$.

We can use Hanson-Wright (Theorem 4.1) to bound both of these quantities. Indeed, it implies that with probability at least $1 - O(1/nk)$,

$$\|\mathbf{A}_{-s}^{-1} \Delta y\|_2^2 \leq O(n^{1-t-2p}). \quad (119)$$

Let us now compute the Hanson-Wright bound for $\langle \mathbf{A}_{-s}^{-1} \Delta y, \mathbf{A}_{-s}^{-1} \Delta y \rangle$. Note that \mathbf{A}_{-s}^{-1} is independent of Δy , so we can condition on \mathbf{A}_{-s}^{-1} and conclude that with probability at least $1 - O(1/nk)$

$$\langle \mathbf{A}_{-s}^{-1} \Delta y, \mathbf{A}_{-s}^{-1} \Delta y \rangle \geq \mathbb{E}[\langle \mathbf{A}_{-s}^{-1} \Delta y, \mathbf{A}_{-s}^{-1} \Delta y \rangle | \mathbf{A}_{-s}^{-1}] - O(n^{\frac{1-t}{2}}) \|\mathbf{A}_{-s}^{-2}\|_2 \sqrt{\log(nk)} \quad (120)$$

$$= \text{Tr}(\mathbf{A}_{-s}^{-2} \mathbb{E}[\Delta y \Delta y^\top]) - O(n^{\frac{1-t}{2}}) \|\mathbf{A}_{-s}^{-2}\|_2 \sqrt{\log(nk)} \quad (121)$$

$$= \frac{2}{k} \text{Tr}(\mathbf{A}_{-s}^{-2}) - O(n^{\frac{1-t}{2}}) \|\mathbf{A}_{-s}^{-2}\|_2 \sqrt{\log(nk)}, \quad (122)$$

where we have used the fact that Δy is mean zero and $\Delta y[i]^2 \sim \text{Ber}(\frac{2}{k})$.

From Proposition B.4, we obtain the scaling for $\text{Tr}(\mathbf{A}_{-s}^{-2})$ and $\|\mathbf{A}_{-s}^{-2}\|_2$. This implies that with probability at least $1 - O(1/nk)$

$$\langle \mathbf{A}_{-s}^{-1} \Delta y, \mathbf{A}_{-s}^{-1} \Delta y \rangle \geq \Omega(n^{1-t-2p}) - O(n^{\frac{1-t}{2}-2p}) \sqrt{\log(nk)} \quad (123)$$

$$\geq \Omega(n^{1-t-2p}). \quad (124)$$

as $t < 1$.

Bernstein and the union bound implies that with probability at least $1 - O(1/nk)$,

$$\langle \mathbf{Z}_F^\top \mathbf{A}_{-s}^{-1} \Delta y, \mathbf{Z}_F^\top \mathbf{A}_{-s}^{-1} \Delta y \rangle \geq \left(\sum_{j \in [s] \setminus [k]} \Omega(n^{1-t-2p}) \right) - O(n^{\frac{r}{2}+1-t-2p}) \quad (125)$$

$$\geq \Omega(n^{r+1-t-2p}), \quad (126)$$

as $r > 0$.

To wrap up, we will need to upper bound the contribution of the error term in Eq. (118). Its contribution from summing over $j \in [s] \setminus [k]$ is $\min\{1, \mu^{-2}\} o(n^{r+1-t-2p})$, which is negligible compared to the Bernstein term, which as we just proved is $\Omega(n^{r+1-t-2p})$. Hence $\langle \mathbf{Z}_F^\top \mathbf{A}_{-R}^{-1} \Delta y, \mathbf{Z}_F^\top \mathbf{A}_{-R}^{-1} \Delta y \rangle \geq \Omega(n^{r+1-t-2p})$ with high probability, and inserting this back into our lower bound Eq. (112), we see that

$$\text{CN}_{\alpha, \beta, L}^2 + \text{CN}_{\alpha, \beta, F}^2 \geq c_{13} \min\{\mu^2, 1\} \lambda_F^2 \Omega(n^{r+1-t-2p}) \quad (127)$$

$$= c_{13} \mu^{-2} \min\{\mu^2, 1\} \Omega(n^{r-t-1}) \quad (128)$$

$$\geq c_{14}^2 \min\{1, \mu^{-2}\} n^{r-t-1}, \quad (129)$$

where c_{14} is a positive constant. \square

939 E.3 Bounding the unfavored contamination

940 Finally, we wrap up the section by proving matching upper and lower bounds for the unfavored
941 contamination $\text{CN}_{\alpha,\beta,U}$.

942 **Lemma E.3** (Bounding unfavored contamination). *In the same setting as Proposition A.2, if $t > 0$,
943 with probability $1 - O(1/nk)$, the contamination from the unfavored features satisfies*

$$\text{CN}_{\alpha,\beta,U}^2 = c_{15}^2 (1 \pm o(1)) n^{1-t-p},$$

944 where c_{15} is a positive constant.

945 On the other hand, if $t = 0$, then with probability $1 - O(1/nk)$, the unfavored contamination satisfies

$$\text{CN}_{\alpha,\beta,U}^2 \leq c_{16}^2 \min\{1, \mu^{-1}\} n^{1-t-p} \log(nsk).$$

946 *Proof.* By Woodbury, we have

$$\mathbf{A}^{-1} = \mathbf{A}_U^{-1} - \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1}, \quad (130)$$

947 where

$$\mathbf{M}_s \triangleq \mathbf{W}_s (\mathbf{I}_s + \mathbf{H}_s)^{-1} \mathbf{W}_s^\top, \quad (131)$$

948 and $\mathbf{H}_s \triangleq \mathbf{W}_s^\top \mathbf{A}_{-s}^{-1} \mathbf{W}_s$.

949 Now we have

$$\text{CN}_{\alpha,\beta,U}^2 = \Delta y^\top \mathbf{A}^{-1} \mathbf{A}_U \mathbf{A}^{-1} \Delta y \quad (132)$$

$$= \Delta y^\top (\mathbf{A}_U^{-1} - \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1}) \mathbf{A}_U (\mathbf{A}_U^{-1} - \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1}) \Delta y \quad (133)$$

$$= \Delta y^\top (\mathbf{A}_U^{-1} - 2\mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1} + \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1}) \Delta y. \quad (134)$$

950 By Theorem 4.1, we have with probability at least $1 - O(1/nk)$

$$\Delta y^\top \mathbf{A}_U^{-1} \Delta y = c_{15}^2 (1 \pm o(1)) n^{1-t-p}, \quad (135)$$

951 where c_{15} is a positive constant.

952 On the other hand, we have that

$$\begin{aligned} \mathbf{M}_s \mathbf{A}_U^{-1} \mathbf{M}_s &= \mathbf{W}_s (\mathbf{I}_s + \mathbf{H}_s)^{-1} \mathbf{W}_s^\top \mathbf{A}_U^{-1} \mathbf{W}_s (\mathbf{I}_s + \mathbf{H}_s)^{-1} \mathbf{W}_s^\top \\ &= \mathbf{W}_s (\mathbf{I}_s + \mathbf{H}_s)^{-1} \mathbf{H}_s (\mathbf{I}_s + \mathbf{H}_s)^{-1} \mathbf{W}_s^\top \\ &= \mathbf{W}_s (\mathbf{I}_s - (\mathbf{I}_s + \mathbf{H}_s)^{-1}) (\mathbf{I}_s + \mathbf{H}_s)^{-1} \mathbf{W}_s^\top \\ &= \mathbf{W}_s ((\mathbf{I}_s + \mathbf{H}_s)^{-1} - (\mathbf{I}_s + \mathbf{H}_s)^{-2}) \mathbf{W}_s^\top \end{aligned}$$

953 Due to Proposition B.6, $\mu_i((\mathbf{I}_s + \mathbf{H}_s)^{-1}) = \min\{\mu, 1\} (1 \pm o(1))$ for all i with very high
954 probability. Hence to handle the error terms that are not $\Delta y^\top \mathbf{A}_U^{-1} \Delta y$, it suffices to asymp-
955 totically bound $\Delta y^\top \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1} \Delta y$. In turn, we can couple this to the quadratic form
956 $\min\{\mu, 1\} (1 \pm o(1)) \Delta y^\top \mathbf{A}_U^{-1} \mathbf{W}_s \mathbf{W}_s^\top \mathbf{A}_U^{-1} \Delta y$. By Proposition C.4, we have with probability at
957 least $1 - O(1/nk)$

$$\min\{\mu, 1\} (1 \pm o(1)) \Delta y^\top \mathbf{A}_U^{-1} \mathbf{W}_s \mathbf{W}_s^\top \mathbf{A}_U^{-1} \Delta y \quad (136)$$

$$\leq c_8^2 \lambda_F \min\{\mu, 1\} (1 \pm o(1)) (n^{2-2t-2p} + n^{r+1-t-2p}) \log(nsk) \quad (137)$$

$$\leq c_8^2 \min\{1, \mu^{-1}\} (1 \pm o(1)) (n^{1-2t-p} + n^{r-t-p}) \log(nsk) \quad (138)$$

958 For $t > 0$, we claim that the term in Eq. (138) is $o(n^{1-t-p})$, because $1 - 2t - p < 1 - t - p$ and
959 $r - t - p < 1 - t - p$. Hence if $t > 0$ then by union bound we have with probability at least
960 $1 - O(1/nk)$ that

$$\text{CN}_{\alpha,\beta,U}^2 = c_{15}^2 (1 \pm o(1)) n^{1-t-p}, \quad (139)$$

961 as desired.

On the other hand, if $t = 0$, we only have an issue if $q + r < 1$, so that $\min \{1, \mu^{-1}\} = 1$. In this case, the deviation term $n^{1-2t-p} = n^{1-t-p}$. However, this won't affect the fact that the upper bound on contamination will still be $\tilde{O}(n^{1-t-p})$. More precisely, this bound concludes by arguing that

$$\text{CN}_{\alpha,\beta,U}^2 \leq c_{16}^2 \min \{1, \mu^{-1}\} n^{1-t-p} \log(nsk), \quad (140)$$

where c_{16} is an appropriately defined positive constant.

It turns out we don't have to worry about this edge case at all for the lower bound on $\text{CN}_{\alpha,\beta,U}$, because the stated conditions for misclassification imply that $t > 0$ anyway. This completes the proof of the lemma. \square

F Obtaining tight misclassification rate

In this section, we will prove Proposition A.7. Let us restate the main proposition and sketch out its proof more formally.

Proposition A.7 (Correlation bound). *Assume we are in the bi-level ensemble model (Definition 1), the true data generating process is 1-sparse (Assumption 1), and the number of classes scales with n (i.e. $t > 0$). Then for every $\epsilon > 0$, we have*

$$\Pr \left[\max_{\beta \in [k], \beta \neq \alpha} Z^{(\beta)} > n^{-u} \right] \geq 1 - \Theta \left(\frac{1}{k^{1+o(1)}} \right) - \epsilon \quad (49)$$

for sufficiently large n and any $u > 0$.

Proof sketch. Note that the $Z^{(\beta)}$'s that must outcompete the decaying survival to contamination ratio are jointly Gaussian, as they are projections of a standard Gaussian vector $\mathbf{x}_{\text{test}} \in \mathbb{R}^d$. Hence if we want to study the probability that $\max_{\beta} Z^{(\beta)}$ outcompetes n^{-u} , we have to understand the correlation structure of the $Z^{(\beta)}$'s.

We will argue that for $\beta, \gamma \in [k]$ with α, β, γ pairwise distinct, the correlation between $Z^{(\beta)}$ and $Z^{(\gamma)}$ is $\frac{1}{2} \pm o(1)$ with high probability. To that end, we want to look at the correlation (inner product) between the vectors $\{\lambda_j \hat{\mathbf{h}}_{\alpha,\beta}[j]\}$ for $j \notin \{\alpha, \beta\}$ and $\{\lambda_j \hat{\mathbf{h}}_{\alpha,\gamma}[j]\}$ for $j \notin \{\alpha, \gamma\}$. However, note that by independence of the components of \mathbf{x}_{test} from every other random variable and the fact that they are mean zero, we have

$$\mathbb{E}[\hat{\mathbf{h}}_{\alpha,\beta}[\gamma] \mathbf{x}_{\text{test}}[\gamma] \hat{\mathbf{h}}_{\alpha,\gamma}[\beta] \mathbf{x}_{\text{test}}[\beta]] = 0.$$

Hence it suffices to look at the correlation for $j \notin \{\alpha, \beta, \gamma\}$.

We assume WLOG that $\alpha = 1, \beta = 2, \gamma = 3$. Let

$$\Lambda_{\alpha,\beta} \triangleq \text{diag}(1 - \mathbf{1}_{j=\alpha} - \mathbf{1}_{j=\beta})_{j \in [d]} \circ \text{diag}(\lambda_j)_{j \in [d]} \in \mathbb{R}^{d \times d}$$

represent the diagonal matrices containing the squared feature weights with indices α, β zeroed out. Next, let $\mathbf{v}_{\alpha,\beta} \in \mathbb{R}^d$ denote the vector with $\mathbf{v}_{\alpha,\beta}[\alpha] = \mathbf{v}_{\alpha,\beta}[\beta] = 0$ and $\mathbf{v}_{\alpha,\beta}[j] = \lambda_j \hat{\mathbf{h}}_{\alpha,\beta}[j]$ for $j \in [d], j \notin \{\alpha, \beta\}$. Hence $\mathbf{v}_{\alpha,\beta} = \Lambda_{\alpha,\beta}^{1/2}(\hat{\mathbf{f}}_{\alpha} - \hat{\mathbf{f}}_{\beta})$. Since $Z^{(\beta)} = \langle \mathbf{v}_{\alpha,\beta}, \mathbf{x}_{\text{test}} \rangle$, in order to analyze the correlations between $Z^{(\beta)}$ and $Z^{(\gamma)}$, it suffices to analyze $\mathbf{v}_{\alpha,\beta}$. Indeed, we will show that the weighted halfspaces $\Lambda_{\alpha,\beta}^{1/2} \hat{\mathbf{f}}_{\alpha} \in \mathbb{R}^d$ and $\Lambda_{\alpha,\beta}^{1/2} \hat{\mathbf{f}}_{\beta} \in \mathbb{R}^d$ are asymptotically orthogonal.

In other words, we need to show that

$$\frac{\langle \Lambda_{\alpha,\beta}^{1/2} \hat{\mathbf{f}}_{\alpha}, \Lambda_{\alpha,\beta}^{1/2} \hat{\mathbf{f}}_{\beta} \rangle}{\|\Lambda_{\alpha,\beta}^{1/2} \hat{\mathbf{f}}_{\alpha}\|_2 \|\Lambda_{\alpha,\beta}^{1/2} \hat{\mathbf{f}}_{\beta}\|_2} = o(1)$$

with probability at least $1 - O(1/nk)$; we can then union bound against all choices of β . This is the most technically involved part of the proof, and is the content of Proposition F.1.

This in turn will imply (see Lemma F.2) that the maximum (and minimum) correlation between the $\mathbf{v}_{\alpha,\beta}$ for different β is $\frac{1}{2} \pm o(1)$. Let $(\bar{Z}_{\beta})_{\beta \in [k], \beta \neq \alpha}$ be equicorrelated gaussians with correlation

997 $\bar{\rho} = \frac{1}{2} + o(1)$, and $(\underline{Z}_\beta)_{\beta \in [k], \beta \neq \alpha}$ be equicorrelated gaussians with correlation $\underline{\rho} = \frac{1}{2} - o(1)$. By
 998 Slepian's lemma, for any $u > 0$, the probability of $\max_\beta Z^{(\beta)}$ losing to n^{-u} is sandwiched as

$$\Pr \left[\max_\beta \underline{Z}_\beta \leq n^{-u} \right] \leq \Pr \left[\max_\beta Z^{(\beta)} \leq n^{-u} \right] \leq \Pr \left[\max_\beta \bar{Z}_\beta \leq n^{-u} \right],$$

999 where we have adopted the shorthand \max_β to denote $\max_{\beta \in [k], \beta \neq \alpha}$.

1000 Theorem 2.1 of Pinasco et al. (2021) shows that jointly gaussian vectors in \mathbb{R}^k with equicorrelation ρ
 1001 lie in the positive orthant with probability $\Theta(k^{1-1/\rho})$. In particular, applied to \bar{Z}_β , with correlation
 1002 $\bar{\rho} = \frac{1}{2} + o(1)$, we find that

$$\Pr \left[\max_\beta \bar{Z}_\beta \leq 0 \right] = \Theta(k^{-1+o(1)}),$$

1003 and similarly for \underline{Z}_β . Anticoncentration for Gaussian maxima (Chernozhukov et al., 2015,
 1004 Corollary 1) implies that we can transfer over the bound on $\Pr[\max_\beta \bar{Z}_\beta \leq 0]$ to a bound on
 1005 $\Pr[\max_\beta \bar{Z}_\beta \leq n^{-u}]$ to show that for every $\epsilon > 0$, we have

$$\Theta(k^{-1+o(1)}) - \epsilon \leq \Pr \left[\max_\beta Z^{(\beta)} \leq n^{-u} \right] \leq \Theta(k^{-1+o(1)}) + \epsilon \quad (141)$$

1006 for sufficiently large n . Taking the complement of the above event concludes the proof. \square

1007 F.1 Main results for tight misclassification rates

1008 The main result in this section is the following proposition, which states that the halfspace predictions
 1009 are asymptotically orthogonal. Its proof is deferred to the subsequent sections.

1010 **Proposition F.1.** *Assume we are in the bi-level ensemble model (Definition 1), the true data generating
 1011 process is 1-sparse (Assumption 1), and the number of classes scales with n (i.e. $t > 0$).*

1012 *For any distinct $\alpha, \beta \in [k]$, with probability at least $1 - O(1/nk)$, we have*

$$\frac{\langle \Lambda_{\alpha,\beta}^{1/2} \hat{\mathbf{f}}_\alpha, \Lambda_{\alpha,\beta}^{1/2} \hat{\mathbf{f}}_\beta \rangle}{\| \Lambda_{\alpha,\beta}^{1/2} \hat{\mathbf{f}}_\alpha \|_2 \| \Lambda_{\alpha,\beta}^{1/2} \hat{\mathbf{f}}_\beta \|_2} = o(1).$$

1013 Given Proposition F.1, we can show that the $Z^{(\beta)}$ have correlations that approach $\frac{1}{2}$. The intuitive
 1014 reason that this correlation approaches $\frac{1}{2}$ is that the contribution from α is common. The following
 1015 lemma formalizes this intuition.

1016 **Lemma F.2** (Correlation of relative differences of almost orthogonal vectors). *Suppose that we have
 1017 n unit vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ such that $|\langle \mathbf{x}_i, \mathbf{x}_j \rangle| \leq \gamma$ for $\gamma > 0$. Then for any distinct $i, j, k \in [n]$,
 1018 we have*

$$\left| \frac{\langle \mathbf{x}_j - \mathbf{x}_i, \mathbf{x}_i - \mathbf{x}_k \rangle}{\| \mathbf{x}_j - \mathbf{x}_i \| \| \mathbf{x}_i - \mathbf{x}_k \|} - \frac{1}{2} \right| \leq \frac{2\gamma}{1 - \gamma}.$$

1019 *Proof.* For any $i \neq j$, we have $\| \mathbf{x}_i - \mathbf{x}_j \|^2 = 2 - 2 \langle \mathbf{x}_i, \mathbf{x}_j \rangle$. Hence we have

$$2 - 2\gamma \leq \| \mathbf{x}_i - \mathbf{x}_j \|^2 \leq 2 + 2\gamma.$$

1020 Also

$$\begin{aligned} 2 - 2\gamma &\leq \| \mathbf{x}_j - \mathbf{x}_k \|^2 \\ &= \| \mathbf{x}_i - \mathbf{x}_j \|^2 + \| \mathbf{x}_i - \mathbf{x}_k \|^2 - 2 \langle \mathbf{x}_j - \mathbf{x}_i, \mathbf{x}_i - \mathbf{x}_k \rangle \\ &\leq 4 + 4\gamma - 2 \langle \mathbf{x}_j - \mathbf{x}_i, \mathbf{x}_i - \mathbf{x}_k \rangle. \end{aligned}$$

1021 Since $\| \mathbf{x}_i - \mathbf{x}_j \| \geq \sqrt{2 - 2\gamma}$, we can rearrange and obtain that

$$\frac{\langle \mathbf{x}_j - \mathbf{x}_i, \mathbf{x}_i - \mathbf{x}_k \rangle}{\| \mathbf{x}_j - \mathbf{x}_i \| \| \mathbf{x}_i - \mathbf{x}_k \|} \leq \frac{1 + 3\gamma}{2 - 2\gamma}.$$

1022 Similarly we can reverse the inequalities and get

$$2 + 2\gamma \geq 4 - 4\gamma - 2 \langle \mathbf{x}_j - \mathbf{x}_i, \mathbf{x}_i - \mathbf{x}_k \rangle,$$

1023 so

$$\frac{\langle \mathbf{x}_j - \mathbf{x}_i, \mathbf{x}_i - \mathbf{x}_k \rangle}{\|\mathbf{x}_j - \mathbf{x}_i\| \|\mathbf{x}_i - \mathbf{x}_k\|} \geq \frac{1 - 3\gamma}{2 + 2\gamma}.$$

1024

□

1025 Combining Proposition F.1 with Lemma F.2 yields the following formal statement about the correla-
1026 tions between the $Z^{(\beta)}$.

1027 **Lemma F.3** (Asymptotic correlation of relative survivals). *For any distinct $\alpha, \beta, \beta' \in [k]$, under the
1028 same assumptions as Proposition F.1, as $n \rightarrow \infty$, with probability at least $1 - O(1/n)$, we have*

$$\left| \mathbb{E}[Z^{(\beta)} Z^{(\beta')}] - \frac{1}{2} \right| \leq o(1).$$

1029 As a consequence, the asymptotic correlation between the relative survivals approaches $\frac{1}{2}$ at a
1030 polynomial rate.

1031 *Proof.* Plugging in the result of Proposition F.1 into Lemma F.2, we obtain the stated result. □

1032 F.2 Lower bounding the denominator

1033 Let us now begin to prove Proposition F.1. The first step is to bound the denominator of the normalized
1034 correlation. Writing out the definitions, we have

$$\begin{aligned} \left\| \Lambda_{\alpha, \beta}^{1/2} \hat{\mathbf{f}}_{\alpha} \right\|^2 &= \sum_{j \notin \{\alpha, \beta\}} \lambda_j^2 \mathbf{y}_{\alpha}^{\top} \mathbf{A}^{-1} \mathbf{z}_j \mathbf{z}_j^{\top} \mathbf{A}^{-1} \mathbf{y}_{\alpha} \\ &= \lambda_F^2 \sum_{j \notin \{\alpha, \beta\}, j \in [s]} \mathbf{y}_{\alpha}^{\top} \mathbf{A}^{-1} \mathbf{z}_j \mathbf{z}_j^{\top} \mathbf{A}^{-1} \mathbf{y}_{\alpha} + \lambda_U^2 \sum_{j > s} \mathbf{y}_{\alpha}^{\top} \mathbf{A}^{-1} \mathbf{z}_j \mathbf{z}_j^{\top} \mathbf{A}^{-1} \mathbf{y}_{\alpha} \end{aligned}$$

1035 Note that these two terms are respectively analogous to $\text{CN}_{\alpha, \beta, L}^2 + \text{CN}_{\alpha, \beta, F}^2$ and $\text{CN}_{\alpha, \beta, U}^2$. In fact, the
1036 proofs of the lower bounds for contamination essentially transfer over verbatim to the lower bounds
1037 on the denominator, because Hanson-Wright implies that we can show that $\|\mathbf{A}_{-s}^{-1} \mathbf{y}_{\alpha}\|_2$ concentrates
1038 the same way that $\|\mathbf{A}_{-s}^{-1} \Delta y\|_2$ does. In essence, we are able to show the following proposition.

1039 **Proposition F.4** (Lower bound on norm of scaled halfspaces). *Under the same assumptions as
1040 Proposition F.1, for any $\alpha, \beta \in [k]$, with $\alpha \neq \beta$, with probability at least $1 - O(1/nk)$, we have*

$$\left\| \Lambda_{\alpha, \beta}^{1/2} \hat{\mathbf{f}}_{\alpha} \right\|^2 \geq \min \{1, \mu^{-2}\} \Omega(n^{r-t-1}) + \Omega(n^{1-t-p}).$$

1041 F.3 Upper bounding the numerator: the unnormalized correlation

1042 We now turn to the more involved part of the bound: proving an upper bound on the numerator. As
1043 before, we can bound the split up the numerator into favored and unfavored terms. For each term, we
1044 will show that it is dominated by the denominator, in the precise sense that each term is

$$o(\min \{1, \mu^{-2}\} n^{r-t-1} + n^{1-t-p}).$$

1045 Now, let's look at the numerator, which is the bilinear form

$$\lambda_F^2 \sum_{j \notin \{\alpha, \beta\}, j \in [s]} \mathbf{y}_{\alpha}^{\top} \mathbf{A}^{-1} \mathbf{z}_j \mathbf{z}_j^{\top} \mathbf{A}^{-1} \mathbf{y}_{\beta} + \lambda_U^2 \sum_{j > s} \mathbf{y}_{\alpha}^{\top} \mathbf{A}^{-1} \mathbf{z}_j \mathbf{z}_j^{\top} \mathbf{A}^{-1} \mathbf{y}_{\beta}. \quad (142)$$

$$\begin{aligned} &= \underbrace{\lambda_F^2 \langle \mathbf{Z}_L^{\top} \mathbf{A}^{-1} \mathbf{y}_{\alpha}, \mathbf{Z}_L^{\top} \mathbf{A}^{-1} \mathbf{y}_{\beta} \rangle}_{\text{cor}_{\alpha, \beta, L}} + \underbrace{\lambda_F^2 \langle \mathbf{Z}_F^{\top} \mathbf{A}^{-1} \mathbf{y}_{\alpha}, \mathbf{Z}_F^{\top} \mathbf{A}^{-1} \mathbf{y}_{\beta} \rangle}_{\text{cor}_{\alpha, \beta, F}} + \underbrace{\lambda_U^2 \sum_{j > s} \mathbf{y}_{\alpha}^{\top} \mathbf{A}^{-1} \mathbf{z}_j \mathbf{z}_j^{\top} \mathbf{A}^{-1} \mathbf{y}_{\beta}}_{\text{cor}_{\alpha, \beta, U}} \end{aligned} \quad (143)$$

1046 We refer to the the first term as the label defining correlation $\text{cor}_{\alpha, \beta, L}$, the second term as the
1047 favored correlation $\text{cor}_{\alpha, \beta, F}$, and the last term as the unfavored correlation $\text{cor}_{\alpha, \beta, U}$. Here, we abuse
1048 terminology slightly and refer to these inner products as *correlations*, even though strictly speaking,
1049 they are unnormalized.

1050 F.3.1 Bounding the favored correlation

1051 We now bound the correlation coming from the favored features; we will ultimately show that its
 1052 contribution is $\min\{1, \mu^{-2}\}o(n^{r-t-1})$. Recall that $\mathbf{W}_R \in \mathbb{R}^{n \times (s-2)}$ is the weighted feature matrix
 1053 for the $s-2$ favored features aside from α and β . Then the label-defining+favored correlation
 1054 $\text{cor}_{\alpha,\beta,L} + \text{cor}_{\alpha,\beta,F}$ can be written succinctly as

$$\lambda_F^2 \langle \mathbf{Z}_R^\top \mathbf{A}^{-1} \mathbf{y}_\alpha, \mathbf{Z}_R^\top \mathbf{A}^{-1} \mathbf{y}_\beta \rangle. \quad (144)$$

1055 Why should we be able to bound this better than Cauchy-Schwarz? Intuitively, although there is a
 1056 mild dependence between \mathbf{y}_α and \mathbf{y}_β , it is not strong enough to cause $\mathbf{Z}_R^\top \mathbf{A}^{-1} \mathbf{y}_\alpha$ and $\mathbf{Z}_R^\top \mathbf{A}^{-1} \mathbf{y}_\beta$ to
 1057 point in the same direction.

1058 To formalize this argument, we will first follow the strategy to bound the favored *contamination*. In
 1059 particular, using the push-through form of Woodbury (Lemma B.2) we see that

$$\langle \mathbf{Z}_R^\top \mathbf{A}^{-1} \mathbf{y}_\alpha, \mathbf{Z}_R^\top \mathbf{A}^{-1} \mathbf{y}_\beta \rangle = \mathbf{y}_\alpha^\top \mathbf{A}_{-R}^{-1} \mathbf{Z}_R (\mathbf{I}_{s-2} + \mathbf{H}_R)^{-2} \mathbf{Z}_R^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\beta. \quad (145)$$

1060 Now, we can apply Proposition B.6 to replace $(\mathbf{I}_{s-2} + \mathbf{H}_R)^{-2}$ with $\min\{\mu^2, 1\}(\mathbf{I}_{s-2} + \mathbf{E})$, where
 1061 $\|\mathbf{E}\|_2 = O(n^{-\kappa_{11}})$ with extremely high probability. Cauchy-Schwarz yields that

$$\langle \mathbf{Z}_R^\top \mathbf{A}^{-1} \mathbf{y}_\alpha, \mathbf{Z}_R^\top \mathbf{A}^{-1} \mathbf{y}_\beta \rangle \quad (146)$$

$$\leq \min\{\mu^2, 1\} \langle \mathbf{Z}_R^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\alpha, \mathbf{Z}_R^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\beta \rangle \quad (147)$$

$$+ \min\{\mu^2, 1\} \|\mathbf{E}\|_2 \|\mathbf{Z}_R^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\beta\|_2 \|\mathbf{Z}_R^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\alpha\|_2. \quad (148)$$

1062 The term in Eq. (148) can be bounded in the same way that we bounded the favored contamination.
 1063 Indeed, since we can swap in \mathbf{y}_α and \mathbf{y}_β with $\Delta \mathbf{y}$, the argument that proved the bounds on
 1064 $\|\mathbf{Z}_R^\top \mathbf{A}_{-R}^{-1} \Delta \mathbf{y}\|_2$ port over immediately. After using the scaling for λ_F and the fact that $\|\mathbf{E}\|_2 =$
 1065 $O(n^{-\kappa_{11}})$, we conclude that this Cauchy-Schwarz error term is at most $\min\{1, \mu^{-2}\}o(n^{r-t-1})$
 1066 with probability at least $1 - O(1/nk)$.

1067 Let us now turn to the term in Eq. (147). As in the proof for the lower bound for favored contamination
 1068 Lemma E.2, to get better concentration than Cauchy-Schwarz, we want to use Bernstein. We can
 1069 rewrite it suggestively as

$$\sum_{j \in [s] \setminus \{\alpha, \beta\}} (\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\alpha) (\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\beta) \quad (149)$$

1070 We cannot immediately power through with the calculation, because \mathbf{A}_{-R}^{-1} is not independent of \mathbf{y}_α
 1071 or \mathbf{y}_β . The main idea is to again use Woodbury and show that the dependent portions contribute
 1072 negligibly to $\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\alpha$. Therefore the dependent contributions get dominated by the lower bound
 1073 on the correlation.

1074 As in the proof for bounding the favored contamination, we can further define $\mathbf{W}_{\alpha,\beta} = [\mathbf{w}_\alpha \ \mathbf{w}_\beta]$
 1075 and $\mathbf{H}_{\alpha,\beta}^{(s)} = \mathbf{W}_{\alpha,\beta}^\top \mathbf{A}_{-s}^{-1} \mathbf{W}_{\alpha,\beta}$. Then we can deduce from another application of Woodbury that

$$\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\alpha = \mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha - \mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{W}_{\alpha,\beta} (\mathbf{I}_2 + \mathbf{H}_{\alpha,\beta})^{-1} \mathbf{W}_{\alpha,\beta}^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha. \quad (150)$$

1076 Again, we can argue that the second term is bounded in magnitude by

$$\min\{1, \mu^{-1}\} n^{\frac{1}{2}-t-p} \log(ns) = \min\{1, \mu^{-1}\} O(n^{\frac{1}{2}-t-p} \cdot n^{-\kappa_6}), \quad (151)$$

1077 because $t > 0$. Since Hanson-Wright (Corollary C.3) implies that $\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha = \tilde{O}(n^{\frac{1}{2}-t-p})$, this
 1078 implies that $\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\alpha = \tilde{O}(n^{\frac{1}{2}-t-p})$, and similarly for β . Hence we have

$$(\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\alpha) (\mathbf{z}_j^\top \mathbf{A}_{-R}^{-1} \mathbf{y}_\beta) \quad (152)$$

$$\leq (\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha + \min\{1, \mu^{-1}\} O(n^{\frac{1}{2}-t-p-\kappa_6})) (\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta + \min\{1, \mu^{-1}\} O(n^{\frac{1}{2}-t-p-\kappa_6})) \quad (153)$$

$$\leq (\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha) (\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta) + \min\{1, \mu^{-1}\} O(n^{\frac{1}{2}-t-p-\kappa_6}) \tilde{O}(n^{\frac{1}{2}-t-p}) \quad (154)$$

$$\leq (\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha) (\mathbf{z}_j^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta) + \min\{1, \mu^{-1}\} o(n^{1-t-2p}). \quad (155)$$

1079 This implies that we can rewrite Eq. (149) as

$$\left(\sum_{j \in [s] \setminus \{\alpha, \beta\}} (z_j^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha) (z_j^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta) \right) \pm \min \{1, \mu^{-1}\} o(n^{r+1-t-2p}). \quad (156)$$

1080 Let us argue that the second term in Eq. (156) will be negligible compared to the denominator, which
 1081 is $\min \{1, \mu^{-2}\} \Omega(n^{r-t-1})$. Tracing back up the stack, we see that its contribution to the favored
 1082 correlation will be at most

$$\lambda_F^2 \min \{\mu^2, 1\} \cdot \min \{1, \mu^{-1}\} o(n^{r+1-t-2p}) \leq \mu^{-2} \min \{\mu^2, \mu^{-1}\} o(n^{r-t-1}) \quad (157)$$

$$\leq \min \{1, \mu^{-3}\} o(n^{r-t-1}) \quad (158)$$

$$\leq \min \{1, \mu^{-2}\} o(n^{r-t-1}). \quad (159)$$

1083 Turning back to the first term, we are now in a position to apply Bernstein. Note that $(\mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha, \mathbf{y}_\beta)$
 1084 are independent of \mathbf{Z}_R . Hence conditioned on $\mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha$, and \mathbf{y}_β , $\langle \mathbf{Z}_R^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha, \mathbf{Z}_R^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta \rangle$ is a sum
 1085 of $s - 2$ subexponential variables, and by Lemma 2.7.7 of Vershynin (2018) each of these random
 1086 variables conditionally has subexponential norm at most $\|\mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha\|_2 \|\mathbf{A}_{-s}^{-1} \mathbf{y}_\beta\|_2$ and conditional
 1087 mean $\langle \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha, \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta \rangle$.

1088 We can use Hanson-Wright (Theorem 4.1) to bound both of these quantities. Indeed, it implies that
 1089 with probability at least $1 - O(1/nk)$,

$$\|\mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha\|_2 \|\mathbf{A}_{-s}^{-1} \mathbf{y}_\beta\|_2 \leq O(n^{1-t-2p}). \quad (160)$$

1090 Let us now compute the Hanson-Wright bound for $\langle \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha, \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta \rangle$. Note that \mathbf{A}_{-s}^{-1} is independent
 1091 of \mathbf{y}_α and \mathbf{y}_β , so we can condition on \mathbf{A}_{-s}^{-1} and conclude that with probability at least $1 - O(1/nk)$

$$\langle \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha, \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta \rangle \leq \mathbb{E}[\langle \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha, \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta \rangle | \mathbf{A}_{-s}^{-1}] + c_6 n^{\frac{1-t}{2}} \|\mathbf{A}_{-s}^{-1}\|_2 \sqrt{\log(nk)}. \quad (161)$$

1092 We can rewrite the expectation as

$$\text{Tr}(\mathbf{A}_s^{-2} \mathbb{E}[\mathbf{y}_\beta \mathbf{y}_\alpha^\top]). \quad (162)$$

1093 Clearly, $\mathbb{E}[\mathbf{y}_\beta \mathbf{y}_\alpha^\top]$ is diagonal, and each diagonal entry is equal to ρ . Let $\rho = \frac{1}{k}$. Then since $\mathbf{y}_\alpha = 1 - \rho$
 1094 implies $\mathbf{y}_\beta = -\rho$ and vice versa, we get

$$\mathbb{E}[\mathbf{y}_\alpha[i] \mathbf{y}_\beta[i]] = 2(1 - \rho)(-\rho) \Pr[\mathbf{y}_\alpha[i] = 1 - \rho] + (-\rho)^2 \Pr[\mathbf{y}_\alpha[i] = \mathbf{y}_\beta[i] = -\rho] \quad (163)$$

$$\leq -2\rho^2(1 - \rho) + \rho^2 \Pr[\mathbf{y}_\alpha[i] = -\rho] \quad (164)$$

$$\leq -\rho^2(1 - \rho). \quad (165)$$

1095 In other words, the expectation is *negative*, so we can neglect it in our upper bound.

1096 On the other hand, the deviation term is with very high probability at most

$$c_6 n^{\frac{1-t}{2}} \|\mathbf{A}_{-s}^{-1}\|_2 \sqrt{\log(nk)} \leq c_6 n^{\frac{1-t}{2}-2p} \sqrt{\log(nk)}. \quad (166)$$

1097 Combining all of our bounds, Bernstein yields

$$\langle \mathbf{Z}_R^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha, \mathbf{Z}_R^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta \rangle \leq \left(\sum_{j \in [s] \setminus \{\alpha, \beta\}} c_6 n^{\frac{1-t}{2}-2p} \sqrt{\log(nk)} \right) + n^{\frac{r}{2}+1-t-2p} \quad (167)$$

$$\leq n^{r+\frac{1-t}{2}-2p} \sqrt{\log(nk)} + n^{\frac{r}{2}-t+1-2p}. \quad (168)$$

1098 Again, let's trace all the way back to Eq. (147) and then the favored correlation bound. We have
 1099 shown that $\langle \mathbf{Z}_R^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\alpha, \mathbf{Z}_R^\top \mathbf{A}_{-s}^{-1} \mathbf{y}_\beta \rangle$'s contribution to the favored correlation is at most

$$c_{17} \lambda_F^2 \min \{\mu^2, 1\} (n^{r+\frac{1-t}{2}-2p} \sqrt{\log(nk)} + n^{\frac{r}{2}-t+1-2p}) \quad (169)$$

$$= c_{17} \mu^{-2} \min \{\mu^2, 1\} (n^{r-\frac{1+t}{2}-1} \sqrt{\log(nk)} + n^{\frac{r}{2}-t-1}) \quad (170)$$

$$\leq c_{17} \min \{1, \mu^{-2}\} (n^{r-\frac{1+t}{2}-1} \sqrt{\log(nk)} + n^{\frac{r}{2}-t-1}) \quad (171)$$

$$\leq c_{17} \min \{1, \mu^{-2}\} o(n^{r-t-1}), \quad (172)$$

1100 where the last line follows because $0 < t < r < 1$, and c_{17} is a positive constant.

1101 **F.3.2 Bounding the unfavored correlation**

1102 Now, let us show that the unfavored correlation $\text{cor}_{\alpha,\beta,U}$ is negligible; more precisely, we'll show
 1103 that it's $\min\{1, \mu^{-2}\}o(n^{r-t-1}) + o(n^{1-t-p})$. We can rewrite $\text{cor}_{\alpha,\beta,U}$ as

$$\lambda_U \mathbf{y}_\alpha^\top \mathbf{A}^{-1} \mathbf{A}_U \mathbf{A}^{-1} \mathbf{y}_\beta,$$

1104 and play the same game with using Woodbury to replace \mathbf{A}^{-1} with $\mathbf{A}^{-1} - \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1}$, where we
 1105 recall that

$$\mathbf{M}_s \triangleq \mathbf{W}_s (\mathbf{I}_s + \mathbf{H}_s)^{-1} \mathbf{W}_s^\top \in \mathbb{R}^{n \times n}.$$

1106 This yields

$$\mathbf{y}_\alpha^\top \mathbf{A}_U^{-1} \mathbf{y}_\beta - 2 \mathbf{y}_\alpha^\top \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1} \mathbf{y}_\beta + \mathbf{y}_\alpha^\top \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1} \mathbf{y}_\beta. \quad (173)$$

1107 Let us first focus on the first term of Eq. (173). Hanson-Wright implies that with probability at least
 1108 $1 - O(1/n)$,

$$|\mathbf{y}_\alpha^\top \mathbf{A}_U^{-1} \mathbf{y}_\beta - \mathbb{E}[\mathbf{y}_\alpha^\top \mathbf{A}_U^{-1} \mathbf{y}_\beta | \mathbf{A}_U^{-1}]| \leq n^{\frac{1-t}{2}} \|\mathbf{A}_U^{-1}\|_2 \sqrt{\log n}. \quad (174)$$

1109 Also, $\mathbb{E}[\mathbf{y}_\alpha^\top \mathbf{A}_U^{-1} \mathbf{y}_\beta | \mathbf{A}_U^{-1}] = \text{Tr}(\mathbf{A}_U^{-1} \mathbb{E}[\mathbf{y}_\beta \mathbf{y}_\alpha^\top]) = \Theta(n^{1-2t-p})$ with high probability, and
 1110 $\|\mathbf{A}_U^{-1}\|_2 \leq n^{-p}$ with extremely high probability. Hence we see that $\mathbf{y}_\alpha^\top \mathbf{A}_U^{-1} \mathbf{y}_\beta \leq O(n^{\frac{1-t}{2}-p}) \leq$
 1111 $o(n^{1-t-p})$ as $t > 0$.

1112 Next, let's turn to the second and third terms of Eq. (173). We claim that only the second term will be
 1113 relevant to bound asymptotically, and moreover that they are both $\min\{1, \mu^{-2}\}o(n^{r-t-1})$. Since

$$\begin{aligned} \mathbf{M}_s \mathbf{A}_U^{-1} \mathbf{M}_s &= \mathbf{W}_s (\mathbf{I}_s + \mathbf{H}_s)^{-1} \mathbf{H}_s (\mathbf{I}_s + \mathbf{H}_s)^{-1} \mathbf{W}_s^\top \\ &= \mathbf{W}_s ((\mathbf{I}_s + \mathbf{H}_s)^{-1} - (\mathbf{I}_s + \mathbf{H}_s)^{-2}) \mathbf{W}_s^\top, \end{aligned}$$

1114 the second and third term can be rewritten as

$$\begin{aligned} &- 2 \mathbf{y}_\alpha^\top \mathbf{A}_U^{-1} \mathbf{W}_s (\mathbf{I}_s + \mathbf{H}_s)^{-1} \mathbf{W}_s^\top \mathbf{A}_U^{-1} \mathbf{y}_\beta \\ &+ \mathbf{y}_\alpha^\top \mathbf{A}_U^{-1} \mathbf{W}_s ((\mathbf{I}_s + \mathbf{H}_s)^{-1} - (\mathbf{I}_s + \mathbf{H}_s)^{-2}) \mathbf{W}_s^\top \mathbf{A}_U^{-1} \mathbf{y}_\beta. \end{aligned}$$

1115 As we are going to use Hanson-Wright to bound the entries of $\mathbf{Z}_s^\top \mathbf{A}_U^{-1} \mathbf{y}_\alpha$, it follows that only the
 1116 second term of Eq. (173) is relevant asymptotically.

1117 To bound the second term, we will use Cauchy-Schwarz. We see that

$$\mathbf{y}_\alpha^\top \mathbf{A}_U^{-1} \mathbf{M}_s \mathbf{A}_U^{-1} \mathbf{y}_\beta \leq \lambda_F \|(\mathbf{I}_s + \mathbf{H}_s)^{-1}\|_2 \|\mathbf{Z}_s^\top \mathbf{A}_U^{-1} \mathbf{y}_\alpha\|_2 \|\mathbf{Z}_s^\top \mathbf{A}_U^{-1} \mathbf{y}_\beta\|_2 \quad (175)$$

$$\leq \lambda_F \min\{\mu, 1\} O(n^{r-t+1-2p}) \log(ns) \quad (176)$$

$$\leq \mu^{-1} \min\{\mu, 1\} O(n^{r-t-p}) \log(ns) \quad (177)$$

$$\leq \min\{1, \mu^{-1}\} O(n^{r-t-p}) \log(ns), \quad (178)$$

1118 where in the second line we have used Proposition B.6.

1119 Now, note that if regression works, this yields an upper bound of $O(n^{r-t-p}) \log(ns)$. But since
 1120 $p > 1$, this is $o(n^{r-t-1})$, which means this contribution is dominated by the denominator.

1121 On the other hand, if regression fails, then the upper bound is now $\mu^{-1} O(n^{r-t-p}) \log(ns)$, which
 1122 we claim is $o(\mu^{-2} n^{r-t-1})$. Indeed, from the definition of the bi-level ensemble Definition 1, we
 1123 have $p > q + r$, so

$$\begin{aligned} \min\{1, \mu^{-1}\} n^{r-t-p} &\leq \mu^{-1} n^{r-t-1} \cdot n^{1-p} \\ &\leq \mu^{-1} o(n^{r-t-1} \cdot n^{1-q-r}) \\ &= \mu^{-2} o(n^{r-t-1}), \end{aligned}$$

1124 as desired.

1125 Let us now go back to Eq. (173) and combine our two bounds. Since $\lambda_U = O(1)$, we have just shown
 1126 that

$$\text{cor}_{\alpha,\beta,U} \leq \min\{1, \mu^{-2}\} o(n^{r-t-1}) + o(n^{1-t-p}), \quad (179)$$

1127 as desired.

1128 G A new variant of the Hanson-Wright inequality

1129 In this section, we prove Theorem 4.1. First, we outline a high level idea of the proof. The starting
1130 point of the proof is to explicitly decompose the quadratic form into diagonal and off-diagonal terms

$$\mathbf{x}^\top \mathbf{M} \mathbf{y} - \mathbb{E}[\mathbf{x}^\top \mathbf{M} \mathbf{y}] = \sum_{i,j} m_{ij} X_i Y_j - \sum_i m_{ii} \mathbb{E}[X_i Y_i] \quad (180)$$

$$= \underbrace{\sum_i m_{ii} (X_i Y_i - \mathbb{E}[X_i Y_i])}_{\triangleq S_{\text{diag}}} + \underbrace{\sum_{i \neq j} m_{ij} X_i Y_j}_{\triangleq S_{\text{offdiag}}} \quad (181)$$

1131 where in the first line we have used the fact that for $i \neq j$, X_i and Y_j are independent and mean zero
1132 to conclude that $\mathbb{E}[X_i Y_j] = 0$.

1133 We can start with the upper tail inequality $\mathbb{P}[\mathbf{x}^\top \mathbf{M} \mathbf{y} - \mathbb{E}[\mathbf{x}^\top \mathbf{M} \mathbf{y}] > t]$ and conclude the lower tail
1134 inequality by replacing \mathbf{M} with $-\mathbf{M}$. To bound S_{diag} and S_{offdiag} , we will proceed by explicitly
1135 bounding the MGF and applying Chernoff's inequality.

1136 G.1 Diagonal terms

1137 For the diagonal terms, we want to bound the MGF of $S_{\text{diag}} = \sum_i m_{ii} (X_i Y_i - \mathbb{E}[X_i Y_i])$. For
1138 $\lambda^2 < \frac{1}{2C_1 K^2 \max_i m_{ii}^2}$, we obtain

$$\exp(\lambda S_{\text{diag}}) = \prod_{i=1}^n \mathbb{E}_{X_i, Y_i} \exp(\lambda m_{ii} (X_i Y_i - \mathbb{E}[X_i Y_i])) \quad (182)$$

$$\leq \prod_{i=1}^n \mathbb{E}_{Y_i} \mathbb{E}_{X_i} [\exp(\lambda m_{ii} Y_i (X_i - \mathbb{E}[X_i | Y_i])) | Y_i] \quad (183)$$

$$\leq \prod_{i=1}^n \mathbb{E}_{Y_i} \exp(C_1 \lambda^2 m_{ii}^2 K^2 Y_i^2) \quad (184)$$

1139 where we have applied Jensen's inequality in the second line and the subgaussian assumption on X_i
1140 conditioned on Y_i in the last line. Here, C_1 is a universal positive constant relating the equivalent
1141 formulations of subgaussianity [Vershynin \(2018\)](#). Continuing with our calculation, we have

$$\exp(\lambda S_{\text{diag}}) \leq \prod_{i=1}^n \mathbb{E}_{Y_i} [1 + 2C_1 \lambda^2 m_{ii}^2 K^2 Y_i^2] \quad (185)$$

$$\leq \prod_{i=1}^n (1 + 2C_1 \pi \lambda^2 K^2 m_{ii}^2) \quad (186)$$

$$\leq \exp\left(2C_1 \pi \lambda^2 K^2 \sum_{i=1}^n m_{ii}^2\right). \quad (187)$$

1142 where in the first line we have used the inequality $\exp(x) \leq 1 + 2x$ valid for $x \leq \frac{1}{2}$, in the second
1143 line we have used the soft sparsity assumption on Y_i , and in the last line we have used the inequality
1144 $1 + x \leq \exp(x)$, valid for all x .

1145 Now Markov's inequality yields for $\epsilon > 0$ that

$$\Pr[S_{\text{diag}} > \epsilon] \leq \frac{\mathbb{E} \exp(\lambda S_{\text{diag}})}{\exp(\lambda \epsilon)} \quad (188)$$

$$\leq \exp\left(-\lambda \epsilon + 2\pi C_1 K^2 \lambda^2 \sum_{i=1}^n m_{ii}^2\right), \quad (189)$$

1146 and optimizing λ in the region $\lambda^2 \leq \frac{1}{2C_1 K^2 \max_i m_{ii}^2}$ yields

$$\lambda = \min \left\{ \frac{\epsilon}{2C_1 K^2 \pi \sum_{i=1}^n m_{ii}^2}, \frac{1}{2C_1 K \max_i |m_{ii}|} \right\}. \quad (190)$$

1147 Plugging in this value of λ into the Markov calculation yields the desired upper tail bound. We can
 1148 repeat the argument with $-M$ to get the lower tail bound. A union bound completes the proof.

1149 G.2 Offdiagonal terms

1150 Following Rudelson and Vershynin (2013), for the offdiagonal terms we can decouple the terms in the
 1151 sum. More precisely, the terms in S_{offdiag} involving indices i and j are precisely $m_{ij}X_iY_j + m_{ji}Y_iX_j$.
 1152 The issue is that Y_i can be correlated with X_i , which complicates the behavior of this random variable.
 1153 Decoupling ensures that for any $j \in [n]$ we will have exactly one term which involves either X_j or
 1154 Y_j , so in particular we will regain independence of the terms, allowing us to bound the MGF more
 1155 easily.

1156 Let $\{\delta_i\}_{i \in [n]}$ denote iid Bernoulli's with parameter $1/2$, which are independent of all other random
 1157 variables.

1158 Let

$$S_\delta \triangleq \sum_{i \neq j} m_{ij} \delta_i (1 - \delta_j) X_i Y_j.$$

1159 Since $\mathbb{E}[\delta_i(1 - \delta_j)] = \frac{1}{4}$, we have

$$S_{\text{offdiag}} = 4\mathbb{E}_\delta[S_\delta],$$

1160 Hence, Jensen's inequality yields

$$\mathbb{E}_{\mathbf{x}, \mathbf{y}} \exp(\lambda S_{\text{offdiag}}) \leq \mathbb{E}_{\mathbf{x}, \mathbf{y}, \delta} \exp(4\lambda S_\delta),$$

1161 where we have used the independence of δ and all other random variables. It follows that it suffices
 1162 to upper bound the MGF of S_δ .

1163 Define the random set $\Lambda_\delta = \{i \in [n] : \delta_i = 1\}$ to denote the indices selected by δ . For a vector
 1164 $\mathbf{u} \in \mathbb{R}^n$ we also introduce the shorthand $\mathbf{u}_{\Lambda_\delta}$ to denote the subvector of \mathbf{u} where $\delta_i = 1$ and $\mathbf{u}_{\Lambda_\delta^c}$ to
 1165 denote the subvector of \mathbf{u} where $\delta_i = 0$.

1166 Hence, we can rewrite $S_\delta \triangleq \sum_{i \in \Lambda_\delta, j \in \Lambda_\delta^c} m_{ij} X_i Y_j$. For $|\lambda| \leq \frac{1}{2C_1 K \|\mathbf{M}\|_2}$, we have

$$\mathbb{E} \exp(\lambda S_{\text{offdiag}}) \leq \mathbb{E} \exp(4\lambda S_\delta) \tag{191}$$

$$\leq \mathbb{E}_\delta \prod_{i \in \Lambda_\delta, j \in \Lambda_\delta^c} \mathbb{E}_{\mathbf{x}_{\Lambda_\delta}, \mathbf{y}_{\Lambda_\delta^c}} [\exp(\lambda m_{ij} X_i Y_j)] \tag{192}$$

1167 Now we can use the fact that the X_i and Y_j are mean zero and independent because $i \in \Lambda_\delta$ and
 1168 $j \in \Lambda_\delta^c$, to show that

$$\prod_{i \in \Lambda_\delta, j \in \Lambda_\delta^c} \mathbb{E}_{\mathbf{x}_{\Lambda_\delta}, \mathbf{y}_{\Lambda_\delta^c}} [\exp(\lambda m_{ij} X_i Y_j)] \leq \prod_{i \in \Lambda_\delta, j \in \Lambda_\delta^c} \mathbb{E}_{\mathbf{y}_{\Lambda_\delta^c}} [\exp(C_1 \lambda^2 K^2 m_{ij}^2 Y_j^2)] \tag{193}$$

$$\leq \prod_{i \in \Lambda_\delta, j \in \Lambda_\delta^c} \mathbb{E}_{\mathbf{y}_{\Lambda_\delta^c}} [1 + 2C_1 \lambda^2 K^2 m_{ij}^2 Y_j^2] \tag{194}$$

$$\leq \prod_{i \in \Lambda_\delta, j \in \Lambda_\delta^c} (1 + 2\pi C_1 \lambda^2 K^2 m_{ij}^2 Y_j^2) \tag{195}$$

$$\leq \prod_{i \in \Lambda_\delta, j \in \Lambda_\delta^c} \exp(2\pi C_1 \lambda^2 K^2 m_{ij}^2 Y_j^2) \tag{196}$$

$$\leq \exp\left(2\pi C_1 \lambda^2 K^2 \|\mathbf{M}\|_F^2\right). \tag{197}$$

1169 In the first line, we have used the subgaussianity of X_i ; in the second line, we have used the
 1170 assumption on λ , in the third line, we have used the variance bound on Y_j .

1171 Again, we can apply Markov's inequality and to see that for $\epsilon > 0$,

$$\Pr[S_{\text{diag}} > \epsilon] \leq \frac{\mathbb{E} \exp(\lambda S_{\text{diag}})}{\exp(\lambda \epsilon)} \tag{198}$$

$$\leq \exp\left(-\lambda \epsilon + 2\pi C_1 K^2 \lambda^2 \|\mathbf{M}\|_F^2\right), \tag{199}$$

1172 Picking

$$\lambda = \min \left\{ \frac{\epsilon}{2C_1 K^2 \pi \|M\|_F^2}, \frac{1}{2C_1 K \|M\|_2} \right\} \quad (200)$$

1173 yields the desired result.

1174 H Proofs of main lemmas for concentration of spectrum

1175 The goal of this section is ultimately to prove Proposition B.6, which asserts that for valid (T, S) ,
 1176 the hat matrix $H_{T,S}$ is a flat matrix whose spectrum is $\min \{\mu, 1\} (1 + o(1))$ with extremely high
 1177 probability. First, let us recall some notation. For any $\emptyset \neq T \subseteq S \subseteq [s]$, we can define the (T, S)
 1178 hat matrix as $H_{T,S} \triangleq W_T^\top A_{-S}^{-1} W_T$. Here, W_T is the $n \times |T|$ matrix of weighted features in T , and
 1179 $A_{-T} = A - W_T W_T^\top$ is the leave- T -out Gram matrix.

1180 First, Wishart concentration applied to $W_T^\top W_T$ yields the following result.

1181 **Lemma H.1.** Recall that $\mu \triangleq n^{q+r-1}$ and $W_T^\top W_T \in \mathbb{R}^{|T| \times |T|}$. For any nonempty $T \subseteq [s]$, with
 1182 probability at least $1 - 2e^{-\sqrt{n}}$ we have that for all $i \in [|T|]$,

$$\mu_i(W_T^\top W_T) = \left(1 \pm c_T \sqrt{\frac{|T|}{n}}\right) \mu^{-1} n^p \quad (201)$$

1183 *Proof.* We can apply Lemma B.3 with $M = Z_T \in \mathbb{R}^{n \times |T|}$, with $M = n$, $m = |T| = o(n)$, and
 1184 $\epsilon = n^{\frac{1}{4}} = o(\sqrt{n|T|})$. Hence we have

$$n - 2\sqrt{n|T|} + o(\sqrt{n|T|}) \leq \mu_{|T|}(Z_T^\top Z_T) \leq \mu_1(Z_T^\top Z_T) \leq n + 2\sqrt{n|T|} + o(\sqrt{n|T|}).$$

1185 Plugging in the scaling $\lambda_F = n^{p-q-r}$ and dividing through by n yields the desired result. Here, we
 1186 define c_T to be an appropriately defined positive constant which only depends on $|T|$ (as the favored
 1187 features are identically distributed). \square

1188 Next, we can use Wishart concentration to bound the spectrum of A_U .

1189 **Lemma H.2** (Concentration of spectrum for unfavored Gram matrix). *Throughout this theo-*
 1190 *rem, assume we are in the bi-level model (Definition 1). Define $\mu \triangleq n^{q+r-1}$. Recall $A_U =$*
 1191 *$\lambda_U \sum_{j>s} z_j z_j^\top \in \mathbb{R}^{n \times n}$. With probability at least $1 - 2e^{-n}$, for $i \in [n]$ we have*

$$\mu_i(A_U) = (1 \pm c_{18} n^{\kappa_7}) n^p, \quad (202)$$

1192 *where c_{18} and κ_7 are positive constants. In other words, the spectrum of the unfavored Gram matrix*
 1193 *A_U is flat.*

1194 *Proof.* Note that $A_U = \lambda_U \sum_{j>s} z_j z_j^\top$. Under the bi-level model, $\lambda_U = 1 + o(1)$. Now we can
 1195 apply Lemma B.3 with $M = \sum_{j>s} z_j z_j^\top$, $M = d - s = n^p - n^r$, $m = n = o(d)$, and $\epsilon = \sqrt{2n}$ to
 1196 conclude that with probability at least $1 - 2e^{-n}$, we have

$$d - 2\sqrt{dn} + n - \sqrt{2n} \leq \mu_n\left(\sum_{j>s} z_j z_j^\top\right) \leq \mu_1\left(\sum_{j>s} z_j z_j^\top\right) \leq d + 2\sqrt{dn} + n + \sqrt{2n}. \quad (203)$$

1197 We can obtain the spectrum of A_U by multiplying through by

$$\lambda_U = \frac{(1-a)d}{d-s} \quad (204)$$

$$= 1 + n^{\max\{-q, r-p\}} + o(n^{\max\{-q, r-p\}}), \quad (205)$$

1198 where in the last line we have used the power series expansion for $\frac{1}{1-x} = 1 + x + o(x)$. Preserving
 1199 only first order terms for λ_U and the spectrum of $\sum_{j>s} z_j z_j^\top$ in Eq. (203) yields

$$\mu_i(A_U) = (1 \pm c_{18} n^{\max\{\frac{1-p}{2}, r-p, -q\}}) n^p. \quad (206)$$

1200 In fact, we know $\frac{1-p}{2} > 1-p > r-p$, since $r < 1$ and $p > 1$. This means we can neglect the $r-p$
 1201 term in the max, define $\kappa_7 = \min \left\{ \frac{p-1}{2}, -q \right\} > 0$ and c_{18} to be an appropriately defined positive
 1202 constant. \square

1203 Since $\mathbf{A}_{-T} = \mathbf{A}_U + \mathbf{W}_{[s]\setminus T} \mathbf{W}_{[s]\setminus T}^\top$, we can apply Lemmas H.1 and H.2 to control the spectrum
 1204 of \mathbf{A}_{-T} . We show that there is a (potentially) spiked portion of the spectrum corresponding to the
 1205 $s - |T|$ favored features which were not taken out, whereas the rest of the $n - s + |T|$ eigenvalues
 1206 are flat.

1207 **Lemma H.3.** *Recall that $\mathbf{A}_{-T} \in \mathbb{R}^{n \times n}$. For any nonempty $T \subseteq [s]$, with probability at least*
 1208 *$1 - 2e^{-\sqrt{n}} - 2e^{-n}$, we have that for all $i \in [s - |T|]$,*

$$\mu_i(\mathbf{A}_{-T}) = \left(1 \pm c_T \sqrt{\frac{|T|}{n}}\right) \mu^{-1} n^p + (1 \pm c_{18} n^{-\kappa_7}) n^p. \quad (207)$$

1209 For all $i \in [n] \setminus [s - |T|]$, we have

$$\mu_i(\mathbf{A}_{-T}) = (1 \pm c_{18} n^{-\kappa_7}) n^p. \quad (208)$$

1210 *Proof.* We can write

$$\mathbf{A}_{-T} = \mathbf{W}_{[s]\setminus T} \mathbf{W}_{[s]\setminus T}^\top + \mathbf{A}_U. \quad (209)$$

1211 Weyl's inequality (Horn and Johnson, 2012, Corollary 4.3.15) implies that for any $i \in [n]$, we have

$$\mu_i(\mathbf{W}_{[s]\setminus T} \mathbf{W}_{[s]\setminus T}^\top) + \mu_n(\mathbf{A}_U) \leq \mu_i(\mathbf{A}_{-T}) \leq \mu_i(\mathbf{W}_{[s]\setminus T} \mathbf{W}_{[s]\setminus T}^\top) + \mu_1(\mathbf{A}_U). \quad (210)$$

1212 Then applying Lemmas H.1 and H.2, for $i \in [s - |T|]$ we conclude that

$$\mu_i(\mathbf{A}_{-T}) = \left(1 \pm c_T \sqrt{\frac{|T|}{n}}\right) \mu^{-1} n^p + (1 \pm c_{18} n^{-\kappa_7}) n^p.. \quad (211)$$

1213 which proves Eq. (207).

1214 For $i > s - |T|$, applying Lemma H.2 and the fact that $\mu_i(\mathbf{W}_{[s]\setminus T} \mathbf{W}_{[s]\setminus T}^\top) = 0$ to Eq. (210) yields

$$\mu_i(\mathbf{A}_{-T}) = (1 \pm c_{18} n^{-\kappa_7}) n^p. \quad (212)$$

1215 which proves Eq. (208). \square

1216 By inverting the bounds proved above, we can also control the spectrum of \mathbf{A}_{-T}^{-1} .

1217 **Corollary H.4.** *Recall that $\mathbf{A}_{-T} \in \mathbb{R}^{n \times n}$. For any nonempty $T \subseteq [s]$, with probability at least*
 1218 *$1 - 2e^{-\sqrt{n}} - 2e^{-n}$, we have that for all $i \in [n - s + |T|]$,*

$$\mu_i(\mathbf{A}_{-T}^{-1}) = (1 \pm c_{19} n^{-\kappa_7}) n^{-p} \quad (213)$$

1219 For all $i \in [n] \setminus [n - s + |T|]$, we have

$$\mu_i(\mathbf{A}_{-T}^{-1}) = \min \{\mu, 1\} (1 \pm c_{20} n^{-\kappa_8}) n^{-p}. \quad (214)$$

1220 where κ_8 is a positive constant depending on $|T|$.

1221 *Proof.* By inverting the bounds in Lemma H.3, using the fact that $\mu_i(\mathbf{A}_{-T}^{-1}) = \frac{1}{\mu_{n-i+1}(\mathbf{A}_{-T})}$ we see
 1222 that for $i \in [n - s + |T|]$,

$$\mu_i(\mathbf{A}_{-T}^{-1}) = \frac{1}{1 \pm c_{18} n^{-\kappa_7}} n^{-p} \quad (215)$$

$$= (1 \pm c_{19} n^{-\kappa_7}) n^{-p}, \quad (216)$$

1223 where we have used the power series expansion $\frac{1}{1-x} = 1 + x + o(x^2)$ and c_{19} is a positive constant.

1224 On the other hand, for $i > n - s + |T|$, we get

$$\mu_i(\mathbf{A}_{-T}^{-1}) = \frac{1}{\left(1 \pm c_T \sqrt{\frac{|T|}{n}}\right) \mu^{-1} + (1 \pm c_{18} n^{-\kappa_7})} n^{-p} \quad (217)$$

$$= \min\{\mu, 1\} (1 \pm c_{20} n^{-\kappa_8}) n^{-p}, \quad (218)$$

1225 where c_{20} and κ_8 are positive constants defined as follows. If $q + r < 1$, i.e. regression works,

1226 then $\mu^{-1} = \omega(1)$, so the denominator becomes $\mu^{-1} \left(1 \pm c_T \sqrt{\frac{|T|}{n}} + \mu(1 \pm c_{18} n^{-\kappa_7})\right)$. Then, since

1227 $|T| \leq s = n^r$, we see that we can pick

$$\kappa_8 = \min\left\{\frac{1-r}{2}, 1-q-r\right\}.$$

1228 On the other hand, if $q + r > 1$, i.e. regression fails, then $\mu^{-1} = o(1)$, and so we can define

$$\kappa_8 = \min\{\kappa_7, q + r - 1\}.$$

1229 Hence to cover both cases we can pick

$$\kappa_8 = \min\left\{\frac{1-r}{2}, \kappa_7, |1-q-r|\right\}.$$

1230 The choice of c_{20} is picked by again using the power series expansion for $\frac{1}{1-x}$. \square

1231 Note that Corollary H.4 immediately implies Proposition B.4, with κ_9 defined based on picking
1232 $T = [k]$. We are now in a position to prove that the generalized hat matrices $\mathbf{H}_{T,S}$, and hence the
1233 Woodbury terms $(\mathbf{I}_{|T|} + \mathbf{H}_{T,S})^{-1}$ have a flat spectrum as well.

1234 **Proposition B.6** (Generalized hat matrices are flat). *Assume we are in the bi-level ensemble Definition 1. For any nonempty $T \subseteq S \subseteq [s]$, with probability at least $1 - 2e^{-\sqrt{n}} - 2e^{-n}$, we have all the*
1235 *eigenvalues tightly controlled:*
1236

$$\mu_i((\mathbf{I}_{|T|} + \mathbf{H}_{T,S})^{-1}) = \min\{\mu, 1\} (1 \pm c_{T,S} n^{-\kappa_{11}}). \quad (74)$$

1237 where $c_{T,S}$ and κ_{11} are positive constants that depend on $|T|$ and $|S|$.

1238 *Proof.* We seek to control the spectrum of the hat matrix $\mathbf{H}_{T,S} = \mathbf{W}_T^\top \mathbf{A}_{-S}^{-1} \mathbf{W}_T$. We cannot directly
1239 use naive eigenvalue bounds to bound the minimum and maximum eigenvalue, as this does not rule
1240 out the possibility that $\mathbf{H}_{T,S}$ has a spike. Instead, we control the spectrum from first principles.

1241 **The spectrum of $\mathbf{H}_{T,S}$ is flat:** By the rotational invariance of the distribution of \mathbf{W}_T and the fact
1242 that \mathbf{A}_{-S}^{-1} is independent of \mathbf{W}_T (as $T \subseteq S$), we can assume WLOG that the symmetric matrix \mathbf{A}_{-S}^{-1}
1243 is diagonal and equal to

$$\mathbf{D} \triangleq \begin{bmatrix} \mathbf{D}_{\text{flat}} & \\ & \mathbf{D}_{\text{spiked}} \end{bmatrix} \in \mathbb{R}^{n \times n} = \begin{bmatrix} \mu_1(\mathbf{A}_{-S}^{-1}) & & \\ & \ddots & \\ & & \mu_n(\mathbf{A}_{-S}^{-1}) \end{bmatrix}, \quad (219)$$

1244 where $\mathbf{D}_{\text{flat}} \in \mathbb{R}^{(n-s+|T|) \times (n-s+|T|)}$ and $\mathbf{D}_{\text{spiked}} \in \mathbb{R}^{(s-|T|) \times (s-|T|)}$ correspond to the flat and
1245 (downwards) spiked portions of the spectrum of \mathbf{A}_{-S}^{-1} . We can also correspondingly decompose

$$\mathbf{Z}_T = \begin{bmatrix} \mathbf{B}_T \\ \mathbf{C}_T \end{bmatrix}, \quad (220)$$

1246 where $\mathbf{B}_T \in \mathbb{R}^{(n-s+|T|) \times |T|}$ and $\mathbf{C}_T \in \mathbb{R}^{(s-|T|) \times |T|}$. Note that each entry of these matrices are
1247 i.i.d. $N(0, 1)$ variables.

1248 By direct computation we have

$$\mathbf{Z}_T^\top \mathbf{D} \mathbf{Z}_T = \mathbf{B}_T^\top \mathbf{D}_{\text{flat}} \mathbf{B}_T + \mathbf{C}_T^\top \mathbf{D}_{\text{spiked}} \mathbf{C}_T \quad (221)$$

1249 We thus have by using standard eigenvalue inequalities that

$$\mu_{|T|}(\mathbf{Z}_T^\top \mathbf{D} \mathbf{Z}_T) \geq \mu_{|T|}(\mathbf{B}_T^\top \mathbf{D}_{\text{flat}} \mathbf{B}_T) + \mu_{|T|}(\mathbf{C}_T^\top \mathbf{D}_{\text{spiked}} \mathbf{C}_T) \quad (222)$$

$$\geq \mu_{|T|}(\mathbf{B}_T^\top \mathbf{B}_T) \mu_{n-s+|T|}(\mathbf{A}_{-S}^{-1}) + \mu_{|T|}(\mathbf{C}_T^\top \mathbf{C}_T) \mu_n(\mathbf{A}_{-S}^{-1}) \quad (223)$$

$$\geq \mu_{|T|}(\mathbf{B}_T^\top \mathbf{B}_T) \mu_{n-s+|T|}(\mathbf{A}_{-S}^{-1}), \quad (224)$$

1250 where in the last line we have used $\mu_n(\mathbf{A}_{-S}^{-1}) \geq 0$.

1251 Since $n > s$, we have $n - s + |T| > |T|$, so we can apply Wishart concentration (Lemma B.3) to

1252 $\mathbf{B}_T^\top \mathbf{B}_T$ to obtain that with probability at least $1 - 2e^{-\sqrt{n}}$ we have

$$\mu_{|T|}(\mathbf{B}_T^\top \mathbf{B}_T) \geq n - s + |T| - 2\sqrt{(n - s + |T|)|T|} + o(\sqrt{(n - s + |T|)|T|}) \quad (225)$$

$$\geq n(1 - n^{r-1} - c_{21}\sqrt{\frac{|T|}{n}}), \quad (226)$$

1253 where c_{21} is a positive constant.

1254 On the other hand, we can deduce that

$$\mu_1(\mathbf{Z}_T^\top \mathbf{D} \mathbf{Z}_T) \leq \mu_1(\mathbf{Z}_T^\top \mathbf{Z}_T) \mu_1(\mathbf{A}_{-S}^{-1}). \quad (227)$$

1255 Lemma H.1 implies that with probability at least $1 - 2e^{-\sqrt{n}}$

$$\mu_1(\mathbf{Z}_T^\top \mathbf{Z}_T) \leq n(1 + c_T\sqrt{\frac{|T|}{n}}) \quad (228)$$

1256 Similarly, Corollary H.4 implies that with probability at least $1 - 2e^{-n} - 2e^{-\sqrt{n}}$, $\mu_1(\mathbf{A}_{-S}^{-1})$ and

1257 $\mu_{n-s+|T|}(\mathbf{A}_{-S}^{-1})$ are both $(1 \pm c_{20}n^{\kappa_8})n^{-p}$. Since $|T| \leq s = n^r$, Eqs. (226) and (228) together

1258 demonstrate that for all $i \in [|T|]$,

$$\mu_i(\mathbf{Z}_T^\top \mathbf{D} \mathbf{Z}_T) = n^{1-p}(1 \pm c_{22}n^{-\kappa_{10}}). \quad (229)$$

1259 Here, c_{22} and κ_{10} are positive constants defined as follows. Since $|T| \leq s = n^r$. Then $\kappa_{10} =$

1260 $\min\{1 - r, \frac{1-r}{2}, \kappa_8\}$, and c_{22} is a constant chosen appropriately based on c_{20} and c_{21} . Plugging in

1261 the scaling $\lambda_F = n^{p-q-r}$, we conclude that with extremely high probability, for all $i \in [|T|]$,

$$\mu_i(\mathbf{H}_{T,S}) = \mu^{-1}(1 \pm c_{22}n^{-\kappa_{10}}). \quad (230)$$

1262 From here, it is easy to compute the spectrum of $(\mathbf{I}_{|T|} + \mathbf{H}_{T,S})^{-1}$. Indeed, reading off our result

1263 from Eq. (230) yields

$$\mu_i((\mathbf{I}_{|T|} + \mathbf{H}_{T,S})^{-1}) = \frac{1}{1 + \mu_{n-i+1}(\mathbf{H}_{T,S})} \quad (231)$$

$$= \min\{\mu, 1\}(1 \pm c_{T,S}n^{-\kappa_{11}}). \quad (232)$$

1264 Here, the positive constant $c_{T,S}$ is picked appropriately and $\kappa_{11} = \min\{\kappa_{10}, |1 - q - r|\} > 0$. This

1265 completes the proof. \square

1266 I Miscellaneous lemmas

1267 **Lemma I.1** (Coupling of quadratic forms). *Let $\mathbf{B} \in \mathbb{R}^{n \times m}$ be an arbitrary real matrix and*

1268 *$\mathbf{M} \in \mathbb{R}^{n \times n}$ be a PSD matrix. Then for any vector $\mathbf{x} \in \mathbb{R}^m$, we have*

$$\lambda_n(\mathbf{M})\mathbf{x}^\top \mathbf{B}^\top \mathbf{B} \mathbf{x} \leq \mathbf{x}^\top \mathbf{B}^\top \mathbf{M} \mathbf{B} \mathbf{x} \leq \lambda_1(\mathbf{M})\mathbf{x}^\top \mathbf{B}^\top \mathbf{B} \mathbf{x}. \quad (233)$$

1269 *Proof.* For any PSD matrix \mathbf{C} , the matrix $\mathbf{B}^\top \mathbf{C} \mathbf{B}$ is PSD. In particular, \mathbf{C} has a unique square root

1270 $\mathbf{C}^{1/2} \in \mathbb{R}^{n \times n}$ with $\mathbf{C}^{1/2} \mathbf{C}^{1/2} = (\mathbf{C}^{1/2})^\top \mathbf{C}^{1/2} = \mathbf{C}$. We thus have

$$\mathbf{x}^\top \mathbf{B}^\top \mathbf{C} \mathbf{B} \mathbf{x} = \mathbf{x}^\top \mathbf{B}^\top \mathbf{C}^{1/2} \mathbf{C}^{1/2} \mathbf{B} \mathbf{x} \quad (234)$$

$$= \left\| \mathbf{C}^{1/2} \mathbf{B} \mathbf{x} \right\|_2^2 \geq 0. \quad (235)$$

1271 Hence

$$\lambda_1(\mathbf{M})\mathbf{x}^\top \mathbf{B}^\top \mathbf{B}\mathbf{x} - \mathbf{x}^\top \mathbf{B}^\top \mathbf{M}\mathbf{B}\mathbf{x} = \mathbf{x}^\top \mathbf{B}^\top (\lambda_1(\mathbf{M})\mathbf{I}_n - \mathbf{M})\mathbf{B}\mathbf{x}. \quad (236)$$

1272 Since $\mathbf{M} \preceq \lambda_1(\mathbf{M})\mathbf{I}_n$ by definition, $\lambda_1(\mathbf{M})\mathbf{I}_n - \mathbf{M}$ is a PSD matrix. Hence by applying Eq. (235),
1273 we conclude that

$$\lambda_1(\mathbf{M})\mathbf{x}^\top \mathbf{B}^\top \mathbf{B}\mathbf{x} - \mathbf{x}^\top \mathbf{B}^\top \mathbf{M}\mathbf{B}\mathbf{x} \geq 0, \quad (237)$$

1274 which gives the upper bound in Eq. (233).

1275 Similarly, $\mathbf{M} \succeq \lambda_n(\mathbf{M})\mathbf{I}_n$, so an analogous argument

$$\lambda_n(\mathbf{M})\mathbf{x}^\top \mathbf{B}^\top \mathbf{B}\mathbf{x} - \mathbf{x}^\top \mathbf{B}^\top \mathbf{M}\mathbf{B}\mathbf{x} \leq 0, \quad (238)$$

1276 which gives the lower bound in Eq. (233). \square

1277 Next, we prove the elementary anti-concentration result that we will need.

1278 **Proposition 1.2** (Gaussian anticoncentration). *Let $\mathbf{x} \sim N(0, \mathbf{I}_d)$ be a standard Gaussian vector, and*
1279 *let $\mathbf{v} \in \mathbb{R}^d$ be arbitrary deterministic vector. Then*

$$\Pr[|\langle \mathbf{x}, \mathbf{v} \rangle| \leq \epsilon] \leq \frac{2\epsilon}{\sqrt{2\pi}\|\mathbf{v}\|_2}.$$

1280 *Proof.* Note that $\langle \mathbf{x}, \mathbf{v} \rangle$ is a linear projection of a standard multivariate Gaussian, so it is itself a
1281 one-dimensional Gaussian. It is also clearly zero mean, and its variance is just give by the squared
1282 norm of \mathbf{v} . So $\langle \mathbf{x}, \mathbf{v} \rangle \sim N(0, \|\mathbf{v}\|_2^2)$. Now we have

$$\Pr[|\langle \mathbf{x}, \mathbf{v} \rangle| \leq \epsilon] = \frac{1}{\sqrt{2\pi}\|\mathbf{v}\|_2} \int_{-\epsilon}^{\epsilon} \exp\left(-\frac{x^2}{\|\mathbf{v}\|_2^2}\right) dx \leq \frac{2\epsilon}{\sqrt{2\pi}\|\mathbf{v}\|_2}.$$

1283 \square

1284 J Comparison to the straightforward non-interpolative scheme

1285 In this section, we quickly give calculations for how well a straightforward non-interpolating scheme
1286 for learning classifiers can work asymptotically. However, a similar analysis using the tools developed
1287 to prove our main results should give a rigorous proof of the below derivation.

1288 This scheme simply uses the sum/average of all positive training examples of a class as the vector we
1289 take an inner-product with to generate scores for classifying test points. For $m \in [k]$, define

$$\hat{\mathbf{f}}_m = \sum_{i:\ell_i=m} \mathbf{x}_i^w. \quad (239)$$

1290 To understand how well this will do asymptotically, it is easy to see that the for the true label-
1291 defining direction, the positive exemplars in the bi-level model will be tightly concentrating around
1292 $\sqrt{2\log k}\sqrt{\lambda_F}$ which, keeping only the polynomial-order scaling, will be like $n^{\frac{p-q-r}{2}}$. There will be
1293 roughly $\frac{n}{k} = n^{1-t}$ positive examples for every class with high probability. For simplicity, let us just
1294 look at $m = 1$ and consider $\frac{k}{n}\hat{\mathbf{f}}_1 = n^{t-1}\hat{\mathbf{f}}_1$. We see

$$n^{t-1}\hat{\mathbf{f}}_1[1] \approx n^{\frac{p-q-r}{2}}. \quad (240)$$

1295 For the other directions that are not true-label defining, we will just have random Gaussians. The
1296 favored directions will be Gaussian with variance $\lambda_F = n^{p-q-r}$ while the unfavored directions will
1297 essentially be Gaussian with unit variance. By averaging over n^{1-t} examples, those variances will
1298 be reduced by that factor. This means that for the $s = n^r$ favored directions, the variance of the
1299 average will be $n^{p-q-r-(1-t)}$ each and for the essentially n^p unfavored directions, the variance of
1300 the average will be n^{t-1} each.

1301 On a test point, we are going to take the inner product of $n^{t-1}\hat{\mathbf{f}}_m$ with an independent random
1302 draw of $\mathbf{x}_{\text{test}}^w$. For classification to succeed, we need this inner product to be dominated by the true

1303 m -th feature-defining direction. When that happens, the correct label will win the comparison. One
 1304 can easily see that the contribution from the true feature-defining direction will be a Gaussian with
 1305 mean 0 and variance $\lambda_F \cdot (n^{\frac{p-q-r}{2}})^2 = \lambda_F^2 = n^{2p-2q-2r}$. Meanwhile, the s favored features will
 1306 have their scaled variances sum up in the score to give a total variance of $n^r \cdot \lambda_F \cdot n^{p-q-r-(1-t)} =$
 1307 $n^{2p-2q-r-(1-t)}$. And finally, the unfavored features will also have their variances sum up in the
 1308 score to give a total variance of $n^p \cdot 1 \cdot n^{t-1} = n^{p+t-1}$.

1309 For the true-feature-defining direction to dominate the contamination from other favored directions,
 1310 we need

$$2p - 2q - 2r > 2p - 2q - r - (1 - t) \quad (241)$$

1311 which immediately gives the condition $t < 1 - r$.

1312 For the true-feature-defining direction to dominate the contamination from other unfavored directions,
 1313 we need

$$2p - 2q - 2r > p + t - 1 \quad (242)$$

1314 which gives the condition $t < p + 1 - 2(q + r)$.

1315 Here, there is no difference between regimes in which regression works or does not work. The
 1316 condition for classification to asymptotically succeed is $t < \min(1 - r, p + r - 2(q + r))$.

1317 Notice that when MNI regression does not work $q + r > 1$, this is identical to the tight characterization
 1318 given for MNI classification in (13). But in the regime where MNI regression *does* work $q + r < 1$,
 1319 this is different. For MNI classification, (13) tells us that we require $t < \min(1 - r, p - 1)$. Consider
 1320 $q = 0.1, r = 0.5$ and $p = 1.1$. MNI classification can only allow $t < 0.1$. Meanwhile, the non-
 1321 interpolating average-of-positive-examples classifier will work as long as $t < 0.5$. This demonstrates
 1322 the potential for significant suboptimality (in terms of the number of distinct classes that can be
 1323 learned) of MNI classifiers in this regime of benign overfitting for regression.