
Anonymous Learning via Look-Alike Clustering: A Precise Analysis of Model Generalization

Adel Javanmard
University of Southern California, Google Research
ajavanma@usc.edu

Vahab Mirrokni
Google Research
mirrokni@google.com

Abstract

While personalized recommendations systems have become increasingly popular, ensuring user data protection remains a top concern in the development of these learning systems. A common approach to enhancing privacy involves training models using anonymous data rather than individual data. In this paper, we explore a natural technique called *look-alike clustering*, which involves replacing sensitive features of individuals with the cluster’s average values. We provide a precise analysis of how training models using anonymous cluster centers affects their generalization capabilities. We focus on an asymptotic regime where the size of the training set grows in proportion to the features dimension. Our analysis is based on the Convex Gaussian Minimax Theorem (CGMT) and allows us to theoretically understand the role of different model components on the generalization error. In addition, we demonstrate that in certain high-dimensional regimes, training over anonymous cluster centers acts as a regularization and improves generalization error of the trained models. Finally, we corroborate our asymptotic theory with finite-sample numerical experiments where we observe a perfect match when the sample size is only of order of a few hundreds.

1 Introduction

Look-alike modeling in machine learning encompasses a range of techniques that focus on identifying users who possess similar characteristics, behaviors, or preferences to a specific target individual. This approach primarily relies on the principle that individuals with shared attributes are likely to exhibit comparable interests and behaviors. By analyzing the behavior of these look-alike users, look-alike modeling enables accurate predictions for the target user. This technique has been widely used in various domains, including targeted marketing and personalized recommendations, where it plays a crucial role in enhancing user experiences and driving tailored outcomes [26, 19, 18, 21].

In this paper, we use look-alike clustering for a different purpose, namely to anonymize sensitive information of users. Consider a supervised regression setup where the training set contains n pairs (\mathbf{x}_i, y_i) , for $i \in [n]$, with $y_i \in \mathbb{R}$ denoting the response and $\mathbf{x}_i \in \mathbb{R}^d$ representing a high-dimensional vector of features. We consider two groups of features: sensitive features, which contain some personal information about users and should be protected from the learner, and the non-sensitive features. We assume that the learner has access to a clustering structure on users, which is non-private information (e.g. based on non-sensitive features or other non-sensitive data set on users).

We propose a look-alike clustering approach, where we anonymize the individuals’ sensitive features by replacing them with the cluster’s average values. Only the anonymized dataset will be shared with the learner who then uses it to train a model. We refer to Figure 1 for an illustration of this approach. Note that the learner never gets access to the individuals’ sensitive features and so this approach is safe from re-identification attacks where the learner is given access to the pool of individuals’ sensitive information (up to permutation) and may use the non-sensitive features to re-identify the

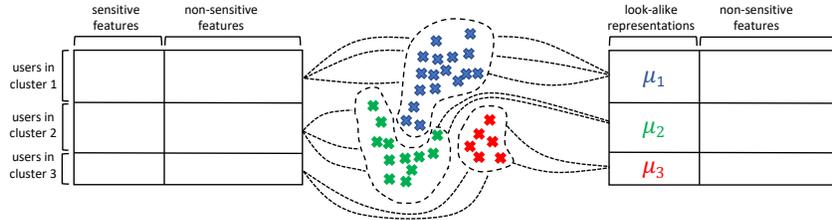


Figure 1: Schematic illustration of look-alike clustering on features data. Within each cluster, the sensitive features of users are replaced by a common look-alike representation (center of the cluster). In this example, μ_1, μ_2, μ_3 represent the average of the sensitive features vectors for users in cluster 1, 2, 3.

users. Also note that since a common representation (average sensitive features) is used for all the users in a cluster, this approach offers m -anonymity provided that each cluster is of size at least m (minimum size clustering).

Minimum size clustering has received an increased attention mainly as a tool for anonymization and when privacy considerations are in place [7, 2, 3]. A particular application is for providing anonymity for user targeting in online advertising with the goal of replacing the use of third-party cookies with a more privacy-respecting entity [10]. There are a variety of approximation algorithms for clustering with minimum size constraint [23, 9, 1, 24], as well as parallel and dynamic implementation [10].

In this paper, we focus on linear regression and derive a precise characterization of model generalization¹ using the look-alike clustering approach, in the so-called *proportional regime* where the size of training set grows in proportion to the number of parameters (which for the linear regression is equal to the number of features). The proportional regime has attracted a significant attention as overparametrized models have become greatly prevalent. It allows to understand the effect under/overparametrization in feature-rich models, providing insights to several intriguing phenomena, including double-descent behavior in the generalization error [22, 8, 14].

Our precise asymptotic theory allows us to demystify the effect of different factors on the model generalization under look-alike clustering, such as the role of cluster size, number of clusters, signal-to-noise ratio of the model as well as the strength of sensitive and non-sensitive features. A key tool in our analysis is a powerful extension of Gordon’s Gaussian process inequality [13] known as the Convex Gaussian Minimax Theorem (CGMT), which was developed in [30] and has been used for studying different learning problems; see e.g. [29, 8, 15, 14, 16].

Initially, it might be presumed that look-alike clustering would hinder model generalization by suppressing sensitive features of individuals, suggesting a possible tradeoff between anonymity (privacy) and model performance. However, our analysis uncovers scenarios in which look-alike clustering actually enhances model generalization! We will develop further insights on these results by arguing that the proposed look-alike clustering can serve as a form of regularization, mitigating model overfitting and consequently improving the model generalization.

Before summarizing our key contributions in this paper, we conclude this section by discussing some of the recent work on the tradeoff between privacy and model generalization at large. An approach to study such potential tradeoff is via the lens of memorization. Modern deep neural networks, with remarkable generalization property, operate in the overparametrized regime where there are more tunable parameters than the number of training samples. Such overparametrized models tend to interpolate the training data and are known to fit well even random labels [34, 33]. Similar phenomenon has been observed in other models, such as random forest [4], Adaboost [25, 32], and kernel methods [5, 17]. Beyond label memorization, [6] studies setting where learning algorithms with near-optimal generalization must encode most of the information about the entire high-dimensional (and high-entropy) covariates of the training examples. Clearly, memorization of training data imposes significant privacy risks when this data contains sensitive personal information, and therefore these results hint to a potential trade-off between privacy protection and model generalization [27, 12, 20]. Lastly, [11] studies settings where data is sampled from a mixture of subpopulations, and shows that label memorization is *necessary* for achieving near-optimal generalization error, whenever the

¹the ability of the model to generalize to new, unseen data from the same distribution as the training data

distribution of subpopulation frequencies is long-tailed. Intuitively, this corresponds to datasets with many small distinct subpopulations. In order to predict more accurately on a subpopulation from which only a very few examples are observed, the algorithm needs to memorize their labels.

1.1 Summary of contributions

We consider a linear regression setting for response variable y given feature \mathbf{x} , and posit a Gaussian Mixture Model on the features to model the clustering structure on the samples. We focus on the high-dimensional asymptotic regime where the number of training samples n , the dimension of sensitive features (p), and the dimension of non-sensitive features ($d - p$) grow in proportion ($p/n \rightarrow \psi_p$ and $d/n \rightarrow \psi_d$, for some constants $0 < \psi_p \leq \psi_d$). Asymptotic analysis in this particular regime, characterized by a fixed sample size to feature size ratio, has recently garnered significant attention due to its relevance to the regime where modern neural networks operate. This analysis allows for the study of various intriguing phenomena related to both statistical properties (such as double-descent) and the tractability of optimizing the learning process in such networks [22, 8, 14], where the population analysis $n/d \rightarrow \infty$ fails to capture. Let $\mathcal{T}^n = \{(\mathbf{x}_i, y_i), i \in [n]\}$ denote the (unanonimized) training set and \mathcal{T}_L^n be the set obtained after replacing the sensitive features with the look-alike representations of clusters. We denote by $\hat{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}_L$ the min-norm estimators fit to \mathcal{T}^n and \mathcal{T}_L^n , respectively. Under this asymptotic setting:

- We provide a precise characterization of the generalization error of $\hat{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}_L$. Despite the randomness in data generating model, we show that in the high-dimensional asymptotic, the generalization errors of these estimators converge in probability to deterministic limits for which we provide explicit expressions.
- Our characterizations reveal several interesting facts about the generalization of the estimators:
 - (i) For the min-norm estimator $\hat{\boldsymbol{\theta}}$ we observe significantly different behavior in the underparametrized regime ($\psi_d \leq 1$) than in the overparametrized regime ($\psi_d > 1$). Note that in the underparametrized regime, the min-norm estimator coincides with the standard least squares estimator. For the look-alike estimator $\hat{\boldsymbol{\theta}}_L$ our analysis identifies the underparametrized regime as $\psi_d - \psi_p \leq 1$ and the overparametrized regime as $\psi_d - \psi_p > 1$.
 - (ii) In the underparametrized regime, our analysis shows that, somewhat surprisingly, the generalization error (for both estimators) does not depend on the number or size of the clusters, nor the scaling of the cluster centers.
 - (iii) In the overparametrized regime, our analysis provides a precise understanding of the role of different factors, including the number of clusters, energy of cluster centers, and the alignment of the model with the constellation of cluster centers, on the generalization error.
- Using our characterizations, we discuss settings where the look-alike estimator $\hat{\boldsymbol{\theta}}_L$ has better generalization than its non-private counterpart $\hat{\boldsymbol{\theta}}$. A relevant quantity that shows up in our analysis is the ratio of the norm of the model component on the sensitive features over the noise in the response, which we refer to as signal-to-noise ratio (SNR). Using our theory, we show that if SNR is below a certain threshold, then look-alike estimator $\hat{\boldsymbol{\theta}}_L$ has lower generalization error than $\hat{\boldsymbol{\theta}}$. This demonstrates scenarios where anonymizing sensitive features via look-alike clustering does ‘not’ hinder model generalization. We give an interpretation for this result, after Theorem 5.1, by arguing that at low-SNR, look-alike clustering acts as a regularization and mitigates overfitting, which consequently improves model generalization.
- In our analysis in the previous parts, we assume that the learner has access to the exact underlying clustering structure on the users, to disentangle the clustering estimation error from look-alike modeling. However, in practice the learner needs to estimate the clustering structure from data. In Section 3.2, we combine our analysis with a perturbation analysis to extend our results to the case of imperfect clustering estimation.

Due to space constraint, we refer to the supplementary material for an overview of our proof techniques as well as proof of theorems and technical lemmas.

2 Model

We consider a linear regression setting, where we are given n i.i.d pairs (\mathbf{x}_i, y_i) , where the response y_i is given by

$$y_i = \langle \mathbf{x}_i, \boldsymbol{\theta}_0 \rangle + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2). \quad (2.1)$$

We assume that there is a clustering structure on features \mathbf{x}_i , $i \in [n]$, independent from the responses. We model this structure via Gaussian-Mixture model.

Gaussian-Mixture Model (GMM) on features. Each example \mathbf{x} belong to cluster $\ell \in [k]$, with probability π_ℓ . We let $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_k] \in \mathbb{R}^k$ with $\boldsymbol{\pi} \geq 0$ and $\mathbf{1}^\top \boldsymbol{\pi} = 1$. The cluster conditional distribution of an example \mathbf{x} in cluster ℓ follows an isotropic Gaussian with mean $\boldsymbol{\mu}_\ell \in \mathbb{R}^d$, namely

$$\mathbf{x} = \boldsymbol{\mu}_\ell + \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}). \quad (2.2)$$

By scaling the model (2.1), without loss of generality we assume $\tau = 1$. Writing in the matrix form, we let

$$\mathbf{X} = [\mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_n] \in \mathbb{R}^{d \times n}, \quad \mathbf{y} = (y_1, \dots, y_n) \in \mathbb{R}^n, \quad \mathbf{M} = [\boldsymbol{\mu}_1 | \boldsymbol{\mu}_2 | \dots | \boldsymbol{\mu}_k] \in \mathbb{R}^{d \times k}. \quad (2.3)$$

It is also convenient to encode the cluster membership as one-hot encoded vectors $\boldsymbol{\lambda}_i \in \mathbb{R}^k$, where $\boldsymbol{\lambda}_i$ is one at entry ℓ (with ℓ being the cluster of example \mathbf{x}_i) and zero everywhere else. The GMM can then be written as

$$\mathbf{X} = \mathbf{M} \boldsymbol{\Lambda} + \mathbf{Z}, \quad (2.4)$$

with $\mathbf{Z} \in \mathbb{R}^{d \times n}$ is a Gaussian matrix with i.i.d $\mathcal{N}(0, 1)$ entries, and $\boldsymbol{\Lambda} \in \mathbb{R}^{k \times n}$ is the matrix obtained by stacking vectors $\boldsymbol{\lambda}_i$ as its column.

Sensitive and non-sensitive features. We assume that some of the features are sensitive for which we have some reservation to share with the learner and some non-sensitive features. Without loss of generality, we write it as $\mathbf{x} = (\mathbf{x}_s, \mathbf{x}_{ns})$, where $\mathbf{x}_s \in \mathbb{R}^p$ representing the sensitive features and $\mathbf{x}_{ns} \in \mathbb{R}^{d-p}$ representing the non-sensitive features. We also decompose the model $\boldsymbol{\theta}_0$ (2.1) as $\boldsymbol{\theta}_0 = (\boldsymbol{\theta}_{0,s}, \boldsymbol{\theta}_{0,ns})$ with $\boldsymbol{\theta}_{0,s} \in \mathbb{R}^p$ and $\boldsymbol{\theta}_{0,ns} \in \mathbb{R}^{d-p}$. Likewise, the cluster mean vector $\boldsymbol{\mu}$ is decomposed as $\boldsymbol{\mu} = (\boldsymbol{\mu}_s, \boldsymbol{\mu}_{ns})$. The idea of look-alike clustering is to replace the sensitive features of an example \mathbf{x}_s with the center of its cluster $\boldsymbol{\mu}_s$. This way, if each cluster is of size at least m , then look-alike clustering offers m -anonymity.

Our goal in this paper is to precisely characterize the effect of look-alike clustering on model generalization. We focus on the high-dimensional asymptotic regime, where the number of training data n , and features sizes d, p grow in proportion.

We formalize the high-dimensional asymptotic setting in the assumption below:

Assumption 1 We assume that the number of clusters k is fixed and focus on the asymptotic regime where $n, d, p \rightarrow \infty$ at a fixed ratio $d/n \rightarrow \psi_d$ and $p/n \rightarrow \psi_p$.

To study the generalization of a model $\boldsymbol{\theta}$ (performance on unseen data) via the *out-of-sample prediction risk* defined as $\text{Risk}(\boldsymbol{\theta}) := \mathbb{E}[(y - \mathbf{x}^\top \boldsymbol{\theta})^2]$, where (y, \mathbf{x}) is generated according to (2.1). Our next lemma characterizes the risk when the feature \mathbf{x} is drawn from GMM.

Lemma 2.1 Under the linear response model (2.1) and a GMM for features \mathbf{x} , the out-of-sample prediction risk of a model $\boldsymbol{\theta}$ is given by

$$\text{Risk}(\boldsymbol{\theta}) = \sigma^2 + \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}\|_{\ell_2}^2 + (\boldsymbol{\theta}_0 - \boldsymbol{\theta})^\top \mathbf{M} \text{diag}(\boldsymbol{\pi}) \mathbf{M}^\top (\boldsymbol{\theta}_0 - \boldsymbol{\theta}).$$

The proof of Lemma 2.1 is deferred to the supplementary.

3 Main results

Consider the minimum ℓ_2 norm (min-norm) least squares regression estimator of \mathbf{y} on \mathbf{X} defined by

$$\widehat{\boldsymbol{\theta}} = (\mathbf{X} \mathbf{X}^\top)^\dagger \mathbf{X} \mathbf{y}, \quad (3.1)$$

where $(\mathbf{X}\mathbf{X}^\top)^\dagger$ denotes the Moore-Penrose pseudoinverse of $\mathbf{X}\mathbf{X}^\top$. This estimator can also be formulated as

$$\widehat{\boldsymbol{\theta}} := \arg \min \left\{ \|\boldsymbol{\theta}\|_{\ell_2} : \boldsymbol{\theta} \text{ minimizes } \|\mathbf{y} - \mathbf{X}^\top \boldsymbol{\theta}\|_{\ell_2} \right\}.$$

We also define the ‘‘look-alike estimator’’ denoted by $\widehat{\boldsymbol{\theta}}_L$, where the sensitive features are first anonymized via look-like modeling, and then the min-norm estimator is computed based on the resulting features. Specifically the sensitive feature \mathbf{x}_s of each sample is replaced by the center of its cluster. In our notation, writing $\mathbf{X}^\top = [\mathbf{X}_s^\top, \mathbf{X}_{ns}^\top]$, we define $\mathbf{X}_L^\top = [(\mathbf{M}_s \boldsymbol{\Lambda})^\top, \mathbf{X}_{ns}^\top]$ the features matrix obtained after look-alike modeling on the sensitive features. The look-alike estimator is then given by

$$\widehat{\boldsymbol{\theta}}_L = (\mathbf{X}_L \mathbf{X}_L^\top)^\dagger \mathbf{X}_L \mathbf{y}, \quad (3.2)$$

Our main result is to provide a precise characterization of the risk of look-alike estimator $\widehat{\boldsymbol{\theta}}_L$ as well as $\widehat{\boldsymbol{\theta}}$ (non-look-alike) in the asymptotic regime, as described in Assumption 1. We then discuss regimes where look-alike clustering offers better generalization.

As our analysis shows there are two majorly different setting in the behavior of the look-alike estimator: (i) $\psi_d - \psi_p \leq 1$, i.e., the sample size n is asymptotically larger than $d - p$, the number of non-sensitive features (referred to as *underparametrized* asymptotics); (ii) $\psi_d - \psi_p \geq 1$, which is referred to as *overparametrized* asymptotics.

Our first theorem is on the risk of look-alike estimator in the underparametrized setting. To present our result, we consider the following singular value decomposition for \mathbf{M}_s , the matrix of cluster centers restricted to sensitive features:

$$\mathbf{M}_s = \mathbf{U}_s \boldsymbol{\Sigma}_s \mathbf{V}_s^\top, \quad \mathbf{U}_s \in \mathbb{R}^{p \times r}, \boldsymbol{\Sigma}_s \in \mathbb{R}^{r \times r}, \mathbf{V}_s \in \mathbb{R}^{k \times r}, \quad \text{with } r = \text{rank}(\mathbf{M}_s) \leq k.$$

Theorem 3.1 (Look-alike estimator, underparametrized regime) *Consider the linear response model (2.1), where the features are coming from the GMM (2.4). Also assume that $\|\boldsymbol{\theta}_{0,s}\| = r_s$ and $\|\mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}\| = \sqrt{\rho} r_s$, for all n, p . Under Assumption 1 with $\psi_d - \psi_p \leq 1$, the out-of-sample prediction risk of look-alike estimator $\widehat{\boldsymbol{\theta}}_L$, defined by (3.2), converges in probability,*

$$\text{Risk}(\widehat{\boldsymbol{\theta}}_L) \xrightarrow{\mathcal{P}} \frac{\sigma^2 + r_s^2}{1 - (\psi_d - \psi_p)} - \rho r_s^2.$$

There are several intriguing observations about this result. In the underparametrized regime:

1. The risk depends on $\boldsymbol{\theta}_{0,s}$ (model component on the sensitive features), only through the norms $\|\boldsymbol{\theta}_{0,s}\| = r_s$ and $\|\mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}\| = \sqrt{\rho} r_s$. Note that $\|\mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}\|$ measures the alignment of the model with the left singular vectors of the cluster centers.
2. The cluster structure on the non-sensitive features plays no role in the risk, nor does $\boldsymbol{\theta}_{0,ns}$ the model component corresponding to the non-sensitive features.
3. The cluster prior probabilities $\boldsymbol{\pi}$ does not impact the risk.

We next proceed to the overparametrized setting. For technical convenience, we make some simplifying assumption, however, we believe a similar derivation can be obtained for the general case, albeit with a more involved analysis.

Assumption 2 *Suppose that there is no cluster structure on the non-sensitive features ($\mathbf{M}_{ns} = \mathbf{0}$). Also, assume orthogonal, equal energy centers for the clusters on the sensitive features ($\mathbf{M}_s = \mu \mathbf{U}_s$ with $\mathbf{U}_s^\top \mathbf{U}_s = \mathbf{I}_k$).*

Our next theorem characterizes the risk of look-alike estimator in the underparametrized regime.

Theorem 3.2 (Look-alike estimator, overparametrized regime) *Consider the linear response model (2.1), where the features are coming from the GMM (2.4). Also assume that $\|\boldsymbol{\theta}_{0,s}\| = r_s$, $\|\boldsymbol{\theta}_{0,ns}\| = r_{ns}$ and $\|\mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}\| = \sqrt{\rho} r_s$, for all n, p, d . Under Assumption 1 with $\psi_d - \psi_p \geq 1$, and Assumption 2, the out-of-sample prediction risk of look-alike estimator $\widehat{\boldsymbol{\theta}}_L$, defined by (3.2), converges in probability,*

$$\text{Risk}(\widehat{\boldsymbol{\theta}}_L) \xrightarrow{\mathcal{P}} \sigma^2 + (1 - \rho)r_s^2 + \gamma_0^2 + \boldsymbol{\alpha}^\top (\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi})) \boldsymbol{\alpha}, \quad (3.3)$$

where $\boldsymbol{\pi} = (\pi_1, \dots, \pi_k)$ encodes the cluster priors and γ_0 and $\boldsymbol{\alpha} \in \mathbb{R}^k$ are given by the following relations:

$$\begin{aligned}\boldsymbol{\alpha} &= \left(\mathbf{I} + \frac{\mu^2 \text{diag}(\boldsymbol{\pi})}{\psi_d - \psi_p - 1} \right)^{-1} \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \\ \gamma_0^2 &= \frac{1}{\psi_d - \psi_p - 1} (\sigma^2 + r_s^2 + \mu^2 \boldsymbol{\alpha}^\top \text{diag}(\boldsymbol{\pi}) \boldsymbol{\alpha}) + \left(1 - \frac{1}{\psi_d - \psi_p} \right) r_{\text{ns}}^2.\end{aligned}$$

As discussed in the introduction, one of the focal interest in this work is to understand cases where look-alike modeling improves generalization. In Section 5 we discuss this by comparing the look-alike estimator $\widehat{\boldsymbol{\theta}}_L$ with the min-norm estimator $\widehat{\boldsymbol{\theta}}$, given by (3.1) which utilizes the full information on the sensitive features. In order to do that, we next derive a precise characterization of the risk of $\widehat{\boldsymbol{\theta}}$ in the asymptotic setting.

Theorem 3.3 (min-norm estimator with no look-alike clustering) *Consider the linear response model (2.1), where the features are coming from the GMM (2.4). Under Assumption 1, the followings hold for the min-norm estimator $\widehat{\boldsymbol{\theta}}$ given by (3.1):*

(a) (underparametrized setting) *If $\psi_d \leq 1$, we have*

$$\text{Risk}(\widehat{\boldsymbol{\theta}}) \xrightarrow{\mathcal{P}} \frac{\sigma^2}{1 - \psi_d}.$$

(b) (overparametrized setting) *If $\psi_d \geq 1$, under Assumption 2, the prediction risk of $\widehat{\boldsymbol{\theta}}$ converges in probability*

$$\text{Risk}(\widehat{\boldsymbol{\theta}}) \xrightarrow{\mathcal{P}} \sigma^2 + \tilde{\gamma}_0^2 + \tilde{\boldsymbol{\alpha}}^\top (\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi})) \tilde{\boldsymbol{\alpha}}, \quad (3.4)$$

where $\tilde{\gamma}_0$ and $\tilde{\boldsymbol{\alpha}}$ are given by the following relations:

$$\begin{aligned}\tilde{\boldsymbol{\alpha}} &= \left(\mathbf{I} + \frac{\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi})}{\psi_d - 1} \right)^{-1} \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \\ \tilde{\gamma}_0^2 &= \frac{1}{\psi_d - 1} (\sigma^2 + \tilde{\boldsymbol{\alpha}}^\top (\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi})) \tilde{\boldsymbol{\alpha}}) + \left(1 - \frac{1}{\psi_d} \right) ((1 - \rho)r_s^2 + r_{\text{ns}}^2).\end{aligned}$$

Example 3.1 (Balanced clusters) *In the case of equal cluster prior ($\pi_1 = \dots = \pi_k = 1/k$), the risk characterization (3.3) depends on $\boldsymbol{\alpha}$ only through $\|\boldsymbol{\alpha}\|_{\ell_2}$ (and likewise, the risk (3.4) depends on $\tilde{\boldsymbol{\alpha}}$ only through its norm). This significantly simplifies these characterizations.*

3.2 Extension to imperfect clustering estimation

In our previous results, we assumed that the underlying cluster memberships of users are known to the learner, so we could concentrate our analysis on the impact of training using anonymous cluster centers. However, in practice, clusters should be estimated from the features and thus includes an estimation error. In our next result, we combine our previous result with a perturbation analysis to bound the risk of the look-alike estimator based on estimated clusters.

Recall matrix $\mathbf{M} \in \mathbb{R}^{d \times k}$ from (2.3), whose columns are the cluster centers. Also, recall the matrix $\boldsymbol{\Lambda} \in \mathbb{R}^{k \times n}$ whose columns are the one-hot encoding of the cluster memberships. We let $\widetilde{\mathbf{M}}$ and $\widetilde{\boldsymbol{\Lambda}}$ indicate the estimated matrices, with the cluster estimation error rate $\delta_n := \frac{1}{\sqrt{n}} \|\mathbf{M}_s \boldsymbol{\Lambda} - \widetilde{\mathbf{M}}_s \widetilde{\boldsymbol{\Lambda}}\|_2$, where $\|\cdot\|_2$ indicates spectral norm. Note that only the cluster estimation error with respect to the sensitive features matters because in the look-alike modeling only those features are anonymized (replaced by the cluster centers).

Proposition 3.4 *Let $\widetilde{\mathbf{X}}^\top := [(\widetilde{\mathbf{M}}_s \widetilde{\boldsymbol{\Lambda}})^\top, \mathbf{X}_{\text{ns}}^\top]$ be the feature matrix after replacing the sensitive features with the estimated cluster centers of users. We also let $\widetilde{\boldsymbol{\theta}}_L = (\widetilde{\mathbf{X}}_L \widetilde{\mathbf{X}}_L^\top)^\dagger \widetilde{\mathbf{X}}_L \mathbf{y}$ be the look-alike estimator based on $\widetilde{\mathbf{X}}_L$. Note that $\widetilde{\boldsymbol{\theta}}_L$ is the counterpart of $\widehat{\boldsymbol{\theta}}_L$ given by (3.2). Define the cluster estimation error rate $\delta_n := \frac{1}{\sqrt{n}} \|\mathbf{M}_s \boldsymbol{\Lambda} - \widetilde{\mathbf{M}}_s \widetilde{\boldsymbol{\Lambda}}\|_2$, and suppose that either of the following conditions hold:*

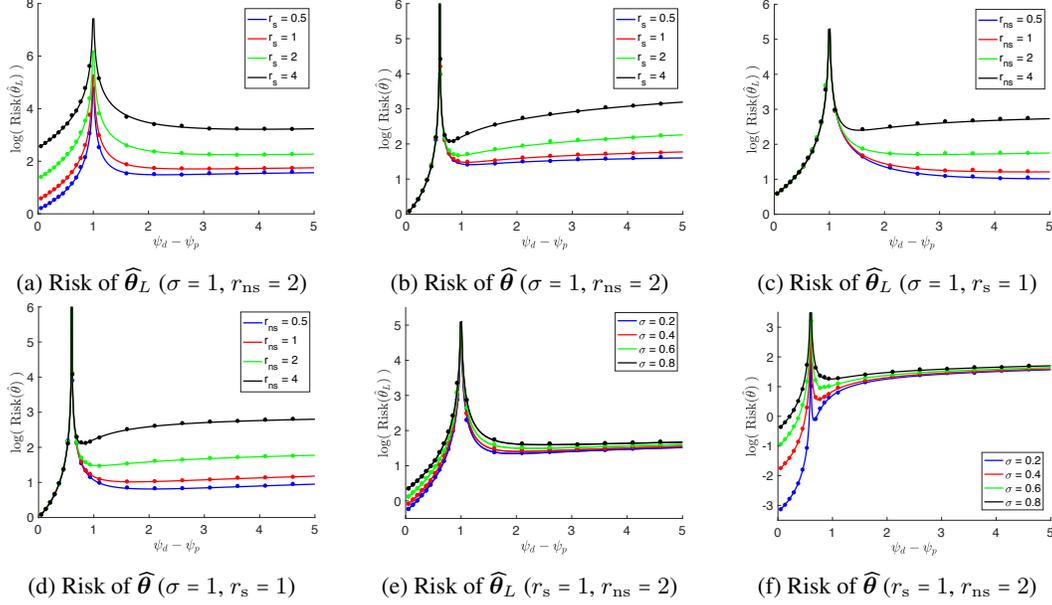


Figure 2: Validation of theoretical characterizations of the risks. Curves correspond to (asymptotic) analytical predictions, and dots to numerical simulations (averaged over 20 realizations). In all the plots, $d = 500, p = 200, \mu = 5, k = 3, \rho = 0.3$. Left panel corresponds to the risk of $\hat{\theta}_L$ and right panel corresponds to the risk of $\hat{\theta}$.

- (i) $\psi_d - \psi_p < 0.5$ and $\delta < \sqrt{1 - (\psi_d - \psi_p)} - \sqrt{\psi_d - \psi_p}$.
- (ii) $\psi_d - \psi_p > 2$ and $\delta < \sqrt{\psi_d - \psi_p} - 1$.

Then,

$$\text{Risk}(\tilde{\theta}_L) \leq \text{Risk}(\hat{\theta}_L) + C\delta,$$

for some constant C depending on the problem parameters.

4 Numerical experiments

In this section, we validate our theory with numerical experiments. We consider GMM with k clusters, where the centers of clusters are given by $\mu \mathbf{u}_\ell$, for $\ell \in [k]$, where $\mathbf{u}_\ell \in \mathbb{R}^d$ are of unit ℓ_2 -norm. Also the vectors \mathbf{u}_ℓ are non-zero only on the first p entries, and their restriction to these entries form a random orthogonal constellation. Therefore, defining $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_k]$, we have $\mathbf{U} = \begin{bmatrix} \mathbf{U}_s \\ \mathbf{0} \end{bmatrix}$, with $\mathbf{U}_s^\top \mathbf{U}_s = \mathbf{I}_k$. In this setting there is no cluster structure on the non-sensitive features and the cluster centers on the sensitive features are orthogonal and of same norm.

Recall the decomposition of the model $\theta_0 = (\theta_{0,s}, \theta_{0,ns})$, with θ_0 the true underlying model (2.1) and $\theta_{0,s}, \theta_{0,ns}$ the components corresponding to sensitive and non-sensitive features. We generate $\theta_{0,ns} \in \mathbb{R}^{d-p}$ to have i.i.d standard normal entries and then normalize it to have $\|\theta_{0,ns}\|_{\ell_2} = r_{ns}$. For $\theta_{0,s}$, we generate $\mathbf{Z}_1, \mathbf{Z}_2 \sim \mathcal{N}(\mathbf{0}_p, \mathbf{I}_p)$, independently and let

$$\theta_{0,s} = r_s \sqrt{\rho} \frac{\mathbf{P}_{\mathbf{U}_s} \mathbf{z}_1}{\|\mathbf{P}_{\mathbf{U}_s} \mathbf{z}_1\|_{\ell_2}} + r_s \sqrt{1 - \rho} \frac{(\mathbf{I} - \mathbf{P}_{\mathbf{U}_s}) \mathbf{z}_2}{\|(\mathbf{I} - \mathbf{P}_{\mathbf{U}_s}) \mathbf{z}_2\|_{\ell_2}},$$

where $\mathbf{P}_{\mathbf{U}_s} := \mathbf{U}_s \mathbf{U}_s^\top$ is the projection onto column space of \mathbf{U}_s . Therefore, $\|\theta_{0,s}\|_{\ell_2} = r_s$ and $\|\mathbf{U}_s^\top \theta_{0,s}\|_{\ell_2} = \sqrt{\rho} r_s$. Note that ρ quantifies the alignment of the model with the cluster centers, confined to the sensitive features.)

We will vary the values of r_s and r_{ns} in our experiments. We also consider the case of balanced clusters, so the cluster prior probabilities are all equal, $\pi_\ell = 1/k$, for $\ell \in [k]$. We set the number of

cluster $k = 3$, dimension of sensitive features $p = 200$ and the dimension of entire features vector $d = 500$. We also set $\mu = 5$ and $\rho = 0.3$. In our experiments, we vary the sample size n and plot the risk of $\widehat{\theta}_L$ and $\widehat{\theta}$ versus $\psi_d - \psi_p = (d - p)/n$. We consider different settings, where we vary r_s , r_{ns} and σ (noise variance in model (2.1)). In Figure 2, we report the results. Curves correspond to our asymptotic theory and dots to the numerical simulations. (Each dot is obtained by averaging over 20 realizations of that configuration.) As we observe, in all scenarios our theoretical predictions are a perfect match to the empirical performance.

5 When does look-alike clustering improve generalization?

In Section 3, we provided a precise characterization of the risk of look-alike estimator $\widehat{\theta}_L$ and its counterpart, the min-norm estimator $\widehat{\theta}$ which utilizes the full information on the sensitive features. By virtue of these characterizations, we would like to understand regimes where the look-alike clustering helps with the model generalization, and the role of different problem parameters in achieving this improvement. Notably, since the look-alike estimator offers m -anonymity on the sensitive features (with m the minimum size of clusters), our discussion here points out instances where data anonymization and model generalization are not in-conflict.

We define the gain of look-alike estimator as $\Delta := \text{Risk}(\widehat{\theta})/\text{Risk}(\widehat{\theta}_L)$ to indicate the gain obtained in generalization via look-alike clustering. For ease in presentation, we focus on the case of balanced clusters (equal priors $\pi_1 = \dots = \pi_k = 1/k$), and consider three cases:

- **Case 1** ($\psi_d \leq 1$): In this case, both $\widehat{\theta}_L$ and $\widehat{\theta}$ are in the underparametrized regime and Theorems 3.1 and 3.3 (a) provide simple closed-form characterization of the risks of $\widehat{\theta}_L$ and $\widehat{\theta}$, by which we obtain

$$\Delta \xrightarrow{\mathcal{P}} \frac{(1 - \psi_d)^{-1}}{(1 + r_s^2/\sigma^2)(1 - \psi_d + \psi_p)^{-1} - \rho r_s^2/\sigma^2}.$$

Define the signal-to-noise ratio $\text{SNR} = r_s/\sigma$. Since $\rho \leq 1$, it is easy to see that Δ is decreasing in the SNR. In particular, as $\text{SNR} \rightarrow 0$, we have $\Delta \rightarrow (1 - \psi_d + \psi_p)/(1 - \psi_d) > 1$, which means the look-alike estimator $\widehat{\theta}_L$ achieves lower risk compared to $\widehat{\theta}$. In Figure 3a we plot $\log(\Delta)$ versus SNR, for several values of ψ_p . Here we set $\psi_d = 0.9$ and $\rho = 0.3$. As we observe in low SNR, the look-alike estimator has lower risk. Specifically, for each curve there is a threshold for the SNR, below which $\log(\Delta) > 0$. Furthermore, this threshold increases with ψ_p , covering a larger range of SNR where $\widehat{\theta}_L$ has better generalization.

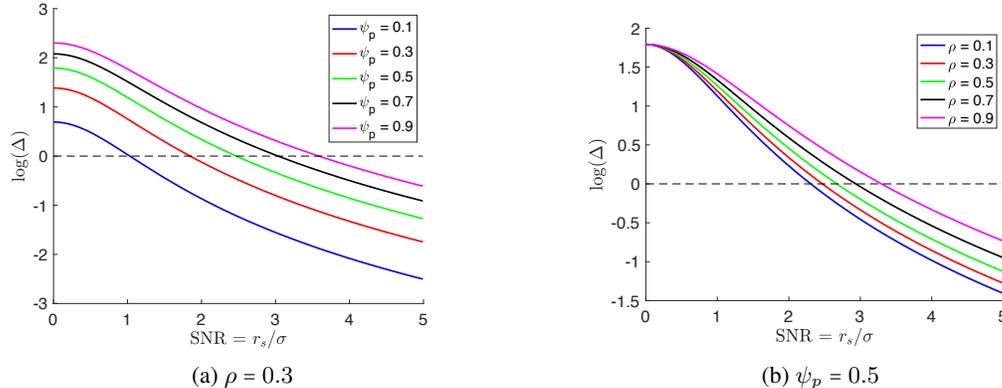


Figure 3: Behavior of gain Δ in the generalization of the look-alike estimator $\widehat{\theta}_L$ over min-norm estimator $\widehat{\theta}$ as we vary $\text{SNR} = r_s/\sigma$. Here, $\psi_d = 0.9$, $\sigma = 1$, and we are in the underparametrized regime for both $\widehat{\theta}_L$ and $\widehat{\theta}$.

In Figure 3b we report similar curves, where this time $\psi_p = 0.5$ and we consider several values of ρ . As we observe, at fixed SNR the gain Δ is increasing in ρ . This is expected since ρ measures the alignment of the underlying model θ_0 with the (left eigenvectors of) cluster centers and so higher ρ is to advantage of the look-alike estimator which uses the cluster centers instead of individuals' sensitive features.

• **Case 2** ($\psi_d \geq 1, \psi_d - \psi_p \leq 1$): In this case, the look-alike estimator $\widehat{\theta}_L$ is in the underparametrized regime, while the min-norm $\widehat{\theta}$ is in the overparametrized regime. The following theorem uses the characterizations in Theorem 3.1 and and 3.3 (b), and shows that in the low SNR= r_s/σ , the look-alike estimator $\widehat{\theta}_L$ has a positive gain. It further shows the monotonicity of the gain with respect to different problem parameters.

Theorem 5.1 *Suppose that $\psi_d \geq 1$ and $\psi_d - \psi_p \leq 1$, and consider the case of equal cluster priors. The gain Δ is increasing in r_{ns} and ρ , and is decreasing in μ^2/k . Furthermore, under the following condition*

$$\text{SNR}^2 := \left(\frac{r_s}{\sigma}\right)^2 \leq \frac{1 + (\psi_d - 1)^{-1} - (1 - \psi_d + \psi_p)^{-1}}{(1 - \psi_d + \psi_p)^{-1} + \psi_d^{-1} - 1}, \quad (5.1)$$

we have $\Delta \geq 1$, for all values of other parameters (μ, k, ρ, r_{ns}).

An interpretation based on regularization: We next provide an argument to build further insight on the result of Theorem 5.1. Recall the data model (2.1), where substituting from (2.2) and decomposing over sensitive and non-sensitive features we arrive at

$$\begin{aligned} y &= \langle \mathbf{x}_s, \boldsymbol{\theta}_s \rangle + \langle \mathbf{x}_{ns}, \boldsymbol{\theta}_{ns} \rangle + \varepsilon \\ &= \langle \boldsymbol{\mu}_s, \boldsymbol{\theta}_s \rangle + \langle \mathbf{z}_s, \boldsymbol{\theta}_s \rangle + \langle \mathbf{x}_{ns}, \boldsymbol{\theta}_{ns} \rangle + \varepsilon. \end{aligned}$$

Note that $\langle \mathbf{z}_s, \boldsymbol{\theta}_s \rangle \sim \text{N}(0, \|\boldsymbol{\theta}_s\|^2)$. At low SNR, this term is of order of the noise term $\varepsilon \sim \text{N}(0, \sigma^2)$. Recall that the look-alike clustering approach replaces the sensitive feature \mathbf{x}_s by the cluster center $\boldsymbol{\mu}_s$, and therefore drops the term $\langle \mathbf{z}_s, \boldsymbol{\theta}_s \rangle$ from the model during the training process. In other words, look-alike clustering acts as a form of regularization which prevents overfitting to the noisy component $\langle \mathbf{z}_s, \boldsymbol{\theta}_s \rangle$, and this will help with the model generalization, together with anonymizing the sensitive features.

In Figure 4a we plot $\log(\Delta)$ versus μ for several values of r_{ns} . Here, $\psi_d = 2, \psi_p = 1.7, \sigma = 1, r_s = 0.5$ and so condition (5.1) holds. As we observe $\log(\Delta)$ is positive, decreasing in μ and also at any fixed μ , it is increasing in r_{ns} , all of which are consistent with the Theorem 5.1. In Figure 4b, we plot similar curves where this time $r_{ns} = 0.2$ and we try several values of ρ . As we see the look-alike estimator has larger gain Δ at larger values of ρ .

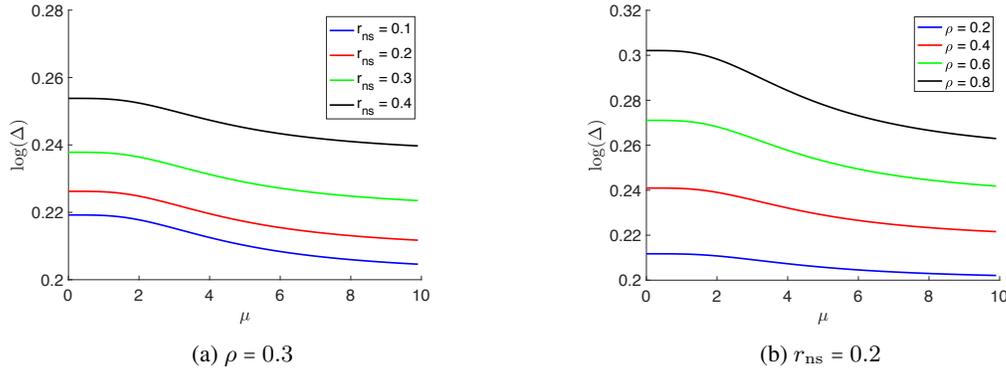


Figure 4: Behavior of gain Δ in the generalization of the look-alike estimator $\widehat{\theta}_L$ over min-norm estimator $\widehat{\theta}$ as we vary μ the energy of cluster centers.

• **Case 3** ($\psi_d - \psi_p \geq 1$): In this case, both $\widehat{\theta}_L$ and $\widehat{\theta}$ are in the overparametrized regime. Let us first focus on r_{ns} , the energy of the model on the non-sensitive features. Invoking the equations (3.3) and (3.4) and hiding the terms that do not depend on r_{ns} in constants C_1, C_2 we arrive at

$$\Delta \stackrel{(P)}{\rightarrow} \frac{C_1 + (1 - \frac{1}{\psi_d})r_{ns}^2}{C_2 + (1 - \frac{1}{\psi_d - \psi_p})r_{ns}^2}.$$

Therefore, $\lim_{r_{ns} \rightarrow \infty} \Delta = (1 - \psi_d^{-1}) / (1 - (\psi_d - \psi_p)^{-1}) > 1$, indicating a gain for the look-alike estimator over $\widehat{\theta}$. In Figure 5a, we plot $\log(\Delta)$ versus r_{ns} for several values of ψ_p . As we observe, when r_{ns} is large enough we always have a gain, which is increasing in ψ_p .

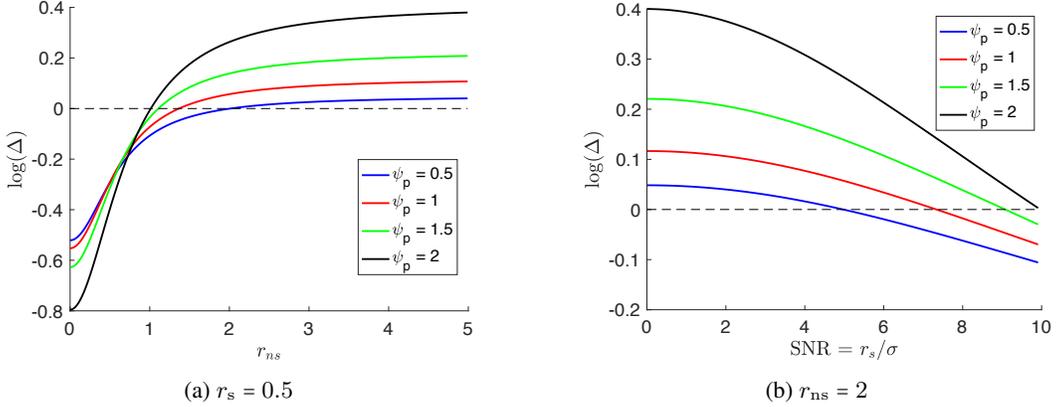


Figure 5: Behavior of gain Δ versus r_{ns} and $SNR:=r_s/\sigma$ for several values of ψ_p . Here, $\psi_d = 4$, $\sigma = 0.1$, $\rho = 0.3$, $\mu = 5$, $k = 5$. Here, we are in the overparametrized regime for both $\widehat{\theta}_L$ and $\widehat{\theta}$.

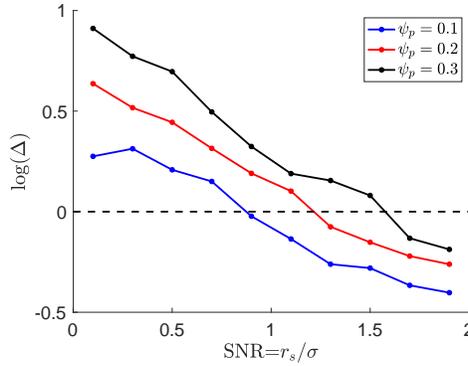


Figure 6: Behavior of gain Δ versus SNR for the nonlinear model described in Section 6. At small SNR, we observe a positive gain (lower risk of look-alike estimator $\widehat{\theta}_L$ compared to $\widehat{\theta}$).

We next consider the effect of $SNR = r_s/\sigma$. In Figure 5b we $\log(\Delta)$ versus SNR, for several values of ψ_p . Similar to the underparametrized regime, we observe that in low SNR, the look-alike estimator has better generalization ($\log(\Delta) > 0$).

6 Beyond linear models

In previous section, we used our theory for linear models to show that at low SNR, look-alike modeling improves model generalization. We also provided an insight for this phenomenon by arguing that look-alike modeling acts as a form of regularization and avoids over-fitting at low SNR regime. In this section we show empirically that this phenomenon also extends to non-linear models.

Consider the following data generative model:

$$y \sim \text{Binomial}(N, p_{\mathbf{x}}), \quad p_{\mathbf{x}} = \frac{1}{1 + \exp(-\langle \mathbf{x}, \boldsymbol{\theta}_0 \rangle + \varepsilon)},$$

where $\varepsilon \sim N(0, \sigma^2)$. We construct the model $\boldsymbol{\theta}_0 = (\boldsymbol{\theta}_{0,s}, \boldsymbol{\theta}_{0,ns})$ similar to the setup in Section 4. We set $n = 200$, $d = 180$, $k = 3$, $\mu = 5$, $\sigma = 1$, $\rho = 0.3$, $r_{ns} = 2$ and $N = 1000$ (number of trials in Binomial distribution). We vary SNR by changing r_s in the set $\{0.1, 0.3, \dots, 1.9\}$. The estimators $\widehat{\theta}$ and $\widehat{\theta}_L$ are obtained by fitting a GLM with logit link function and binomial distribution. We compute the risks of $\widehat{\theta}$ and $\widehat{\theta}_L$ by averaging over a test set of size $50K$. In Figure 6, we plot the gain $\log(\Delta)$ versus r_s where each data point is by averaging over 50 different realizations of data. As we observe at low SNR, $\log(\Delta) > 0$ indicating that the look-alike estimator $\widehat{\theta}_L$ obtains a lower risk than the min-norm estimator.

Acknowledgement

This work is supported in part by the NSF CAREER Award DMS-1844481 and the NSF Award DMS-2311024.

References

- [1] F. Abu-Khzam, C. Bazgan, K. Casel, and H. Fernau. Building clusters with lower-bounded sizes. In *27th International Symposium on Algorithms and Computation (ISAAC 2016)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2016.
- [2] G. Aggarwal, T. Feder, K. Kenthapadi, R. Motwani, R. Panigrahy, D. Thomas, and A. Zhu. Approximation algorithms for k-anonymity. *Journal of Privacy Technology*, 2005112001:400, 2005.
- [3] G. Aggarwal, R. Panigrahy, T. Feder, D. Thomas, K. Kenthapadi, S. Khuller, and A. Zhu. Achieving anonymity via clustering. *ACM Transactions on Algorithms (TALG)*, 6(3):1–19, 2010.
- [4] P. Bartlett, Y. Freund, W. S. Lee, and R. E. Schapire. Boosting the margin: A new explanation for the effectiveness of voting methods. *The annals of statistics*, 26(5):1651–1686, 1998.
- [5] M. Belkin, S. Ma, and S. Mandal. To understand deep learning we need to understand kernel learning. In *International Conference on Machine Learning*, pages 541–549. PMLR, 2018.
- [6] G. Brown, M. Bun, V. Feldman, A. Smith, and K. Talwar. When is memorization of irrelevant training data necessary for high-accuracy learning? In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 123–132, 2021.
- [7] J.-W. Byun, A. Kamra, E. Bertino, and N. Li. Efficient k-anonymization using clustering techniques. In *Advances in Databases: Concepts, Systems and Applications: 12th International Conference on Database Systems for Advanced Applications, DASFAA 2007, Bangkok, Thailand, April 9-12, 2007. Proceedings 12*, pages 188–200. Springer, 2007.
- [8] Z. Deng, A. Kammoun, and C. Thrampoulidis. A model of double descent for high-dimensional binary linear classification. *Information and Inference: A Journal of the IMA*, 11(2):435–495, 2022.
- [9] H. Ding, L. Hu, L. Huang, and J. Li. Capacitated center problems with two-sided bounds and outliers. In *Algorithms and Data Structures: 15th International Symposium, WADS 2017, St. John's, NL, Canada, July 31–August 2, 2017, Proceedings 15*, pages 325–336. Springer, 2017.
- [10] A. Epasto, M. Mahdian, V. Mirrokni, and P. Zhong. Massively parallel and dynamic algorithms for minimum size clustering. In *Proceedings of the 2022 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1613–1660. SIAM, 2022.
- [11] V. Feldman. Does learning require memorization? a short tale about a long tail. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pages 954–959, 2020.
- [12] V. Feldman and C. Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. *Advances in Neural Information Processing Systems*, 33:2881–2891, 2020.
- [13] Y. Gordon. On milman’s inequality and random subspaces which escape through a mesh in r^n . In *Geometric aspects of functional analysis*, pages 84–106. Springer, 1988.
- [14] H. Hassani and A. Javanmard. The curse of overparametrization in adversarial training: Precise analysis of robust generalization for random features regression. *arXiv preprint arXiv:2201.05149*, 2022.
- [15] A. Javanmard and M. Soltanolkotabi. Precise statistical analysis of classification accuracies for adversarial training. *The Annals of Statistics*, 50(4):2127–2156, 2022.

- [16] A. Javanmard, M. Soltanolkotabi, and H. Hassani. Precise tradeoffs in adversarial training for linear regression. In *Conference on Learning Theory*, pages 2034–2078. PMLR, 2020.
- [17] T. Liang and A. Rakhlin. Just interpolate: Kernel “ridgeless” regression can generalize. *The Annals of Statistics*, 48(3):1329–1347, 2020.
- [18] Q. Ma, E. Wagh, J. Wen, Z. Xia, R. Ormandi, and D. Chen. Score look-alike audiences. In *2016 IEEE 16th International Conference on Data Mining Workshops (ICDMW)*, pages 647–654. IEEE, 2016.
- [19] Q. Ma, M. Wen, Z. Xia, and D. Chen. A sub-linear, massive-scale look-alike audience extension system a massive-scale look-alike audience extension. In *Workshop on Big Data, Streams and Heterogeneous Source Mining: Algorithms, Systems, Programming Models and Applications*, pages 51–67. PMLR, 2016.
- [20] M. Malekzadeh, A. Borovykh, and D. Gündüz. Honest-but-curious nets: Sensitive attributes of private inputs can be secretly coded into the classifiers’ outputs. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pages 825–844, 2021.
- [21] A. Mangalampalli, A. Ratnaparkhi, A. O. Hatch, A. Bagherjeiran, R. Parekh, and V. Pudi. A feature-pair-based associative classification approach to look-alike modeling for conversion-oriented user-targeting in tail campaigns. In *Proceedings of the 20th international conference companion on World wide web*, pages 85–86, 2011.
- [22] S. Mei and A. Montanari. The generalization error of random features regression: Precise asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*, 75(4):667–766, 2022.
- [23] H. Park and K. Shim. Approximate algorithms for k-anonymity. In *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*, pages 67–78, 2007.
- [24] A. Sarker, W.-K. Sung, and M. S. Rahman. A linear time algorithm for the r-gathering problem on the line. *Theoretical Computer Science*, 866:96–106, 2021.
- [25] R. E. Schapire. Explaining adaboost. In *Empirical inference*, pages 37–52. Springer, 2013.
- [26] J. Shen, S. C. Geyik, and A. Dasdan. Effective audience extension in online advertising. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2099–2108, 2015.
- [27] C. Song and V. Shmatikov. Overlearning reveals sensitive attributes. *arXiv preprint arXiv:1905.11742*, 2019.
- [28] G. W. Stewart. On the perturbation of pseudo-inverses, projections and linear least squares problems. *SIAM review*, 19(4):634–662, 1977.
- [29] C. Thrampoulidis, E. Abbasi, and B. Hassibi. Precise error analysis of regularized m -estimators in high dimensions. *IEEE Transactions on Information Theory*, 64(8):5592–5628, 2018.
- [30] C. Thrampoulidis, S. Oymak, and B. Hassibi. Regularized linear regression: A precise analysis of the estimation error. In *Conference on Learning Theory*, pages 1683–1709, 2015.
- [31] S. Tu. On the smallest singular value of non-centered gaussian designs, 2020.
- [32] A. J. Wyner, M. Olson, J. Bleich, and D. Mease. Explaining the success of adaboost and random forests as interpolating classifiers. *The Journal of Machine Learning Research*, 18(1):1558–1590, 2017.
- [33] C. Zhang, S. Bengio, M. Hardt, M. C. Mozer, and Y. Singer. Identity crisis: Memorization and generalization under extreme overparameterization. In *International Conference on Learning Representations*, 2020.
- [34] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.

Supplementary Material

In this supplementary, we first provide an overview of our proof techniques in Appendix A and then in Appendix B provide the proofs of theorems and technical lemmas stated in the main paper.

A Overview of proof techniques

Our analysis of the generalization error is based on an extension of Gordon’s Gaussian process inequality [13], called Convex-Gaussian Minimax Theorem (CGMT) [30]. Here, we outline the general steps of this framework and refer to the supplementary for complete details and derivations.

Consider the following two Gaussian processes:

$$\begin{aligned} X_{\mathbf{u},\mathbf{v}} &:= \mathbf{u}^\top \mathbf{G} \mathbf{v} + \psi(\mathbf{u}, \mathbf{v}), \\ Y_{\mathbf{u},\mathbf{v}} &:= \|\mathbf{u}\|_{\ell_2} \mathbf{g}^\top \mathbf{v} + \|\mathbf{v}\|_{\ell_2} \mathbf{h}^\top \mathbf{u} + \psi(\mathbf{u}, \mathbf{v}), \end{aligned}$$

where $\mathbf{G} \in \mathbb{R}^{n \times d}$, $\mathbf{g} \in \mathbb{R}^n$ and $\mathbf{h} \in \mathbb{R}^d$, all have i.i.d standard normal entries. Further, $\psi : \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuous function, which is convex in the first argument and concave in the second argument.

Given the above two processes, consider the following min-max optimization problems, which are respectively referred to as the *Primary Optimization (PO)* and the *Auxiliary Optimization (AO)* problems:

$$\Phi_{\text{PO}}(\mathbf{G}) := \min_{\mathbf{u} \in \mathcal{S}_u} \max_{\mathbf{v} \in \mathcal{S}_v} X_{\mathbf{u},\mathbf{v}}, \tag{A.1}$$

$$\Phi_{\text{AO}}(\mathbf{g}, \mathbf{h}) := \min_{\mathbf{u} \in \mathcal{S}_u} \max_{\mathbf{v} \in \mathcal{S}_v} Y_{\mathbf{u},\mathbf{v}}. \tag{A.2}$$

The main result of CGMT is to connect the above two random optimization problems. As shown in [30](Theorem 3), if \mathcal{S}_u and \mathcal{S}_v are compact and convex then, for any $\lambda \in \mathbb{R}$ and $t > 0$,

$$\mathbb{P}(|\Phi_{\text{PO}}(\mathbf{G}) - \lambda| > t) \leq 2\mathbb{P}(|\Phi_{\text{AO}}(\mathbf{g}, \mathbf{h}) - \lambda| > t).$$

An immediate corollary of this result (by choosing $\lambda = \mathbb{E}[\Phi_{\text{AO}}(\mathbf{g}, \mathbf{h})]$) is that if the optimal cost of AO problem concentrates in probability, then the optimal cost of the corresponding PO problem also concentrates, in probability, around the same value. In addition, as shown in part (iii) of [30](Theorem 3), concentration of the optimal solution of the AO problem implies concentration of the optimal solution of the PO around the same value. Therefore, the two optimization are intimately connected and by analyzing the AO problem, which is substantially simpler, one can derive corresponding properties of the PO problem.

The CGMT framework has been used to infer statistical properties of estimators in certain high-dimensional asymptotic regime. The intermediate steps in the CGMT framework can be summarized as follows: First form an PO problem in the form of (A.1) and construct the corresponding AO problem. Second, derive the point-wise limit of the AO objective in terms of a convex-concave optimization problem, over only few scalar variables. This step is called ‘scalarization’. Next, it is possible to establish uniform convergence of the scalarized AO to the (deterministic) min-max optimization problem using convexity conditions. Finally, by analyzing the latter deterministic problem, one can derive the desired asymptotic characterizations.

Of course implementing the above steps involved problem-specific intricate calculations. Our proofs of Theorems 3.1, 3.2, 3.3 in the supplementary follow this general strategy.

B Proof of theorems and technical lemmas

B.1 Proof of Lemma 2.1

By substituting for y from (2.1) in the definition of risk we obtain

$$\begin{aligned}
\text{Risk}(\boldsymbol{\theta}) &= \mathbb{E}[(y - \mathbf{x}^\top \boldsymbol{\theta})^2] \\
&= \mathbb{E}[(\mathbf{x}^\top (\boldsymbol{\theta}_0 - \boldsymbol{\theta}))^2] + \mathbb{E}[\varepsilon^2] \\
&\stackrel{(a)}{=} \sum_{\ell \in [k]} \pi_\ell \mathbb{E}[(\boldsymbol{\mu}_\ell + \mathbf{z})^\top (\boldsymbol{\theta}_0 - \boldsymbol{\theta})^2] + \mathbb{E}[\varepsilon^2] \\
&= \sum_{\ell \in [k]} \pi_\ell \mathbb{E}[(\boldsymbol{\mu}_\ell^\top (\boldsymbol{\theta}_0 - \boldsymbol{\theta}))^2] + \sum_{\ell \in [k]} \pi_\ell \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_{\ell_2}^2 + \sigma^2 \\
&\stackrel{(b)}{=} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \mathbf{M} \text{diag}(\boldsymbol{\pi}) \mathbf{M}^\top (\boldsymbol{\theta}_0 - \boldsymbol{\theta}) + \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}\|_{\ell_2} + \sigma^2,
\end{aligned}$$

where (a) follows from the Gaussian-Mixture model (2.2) and (b) holds since $\sum_{\ell \in [k]} \pi_\ell = 1$.

B.2 Proof of Theorem 3.1 and Theorem 3.2

Recall that the look-alike estimator is defined as the min-norm estimator over the feature matrix \mathbf{X}_L , where the look-alike representations are used instead of individual sensitive features; see (3.2).

To analyze risk of $\widehat{\boldsymbol{\theta}}_L$, we consider the ridge regression estimator given by

$$\widehat{\boldsymbol{\theta}}_\lambda = \arg \min_{\boldsymbol{\theta}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}_L^\top \boldsymbol{\theta}\|_{\ell_2}^2 + \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2.$$

The minimum-norm estimator is given by $\widehat{\boldsymbol{\theta}}_L = \lim_{\lambda \rightarrow 0^+} \widehat{\boldsymbol{\theta}}_\lambda$.

We follow the CGMT framework explained in Section A. Recall that

$$\mathbf{X}_L = \begin{bmatrix} \mathbf{M}_s \boldsymbol{\Lambda} \\ \mathbf{M}_{ns} \boldsymbol{\Lambda} + \mathbf{Z}_{ns} \end{bmatrix},$$

and therefore by substituting for \mathbf{y} , \mathbf{X} , and \mathbf{X}_L , we get

$$\begin{aligned}
\frac{1}{2n} \|\mathbf{y} - \mathbf{X}_L^\top \boldsymbol{\theta}\|_{\ell_2}^2 &= \frac{1}{2n} \|\varepsilon + \mathbf{X}^\top \boldsymbol{\theta}_0 - \mathbf{X}_L^\top \boldsymbol{\theta}\|_{\ell_2}^2 \\
&= \frac{1}{2n} \|\varepsilon + \boldsymbol{\Lambda}^\top \mathbf{M}_s^\top (\boldsymbol{\theta}_{0,s} - \boldsymbol{\theta}_s) + \mathbf{Z}_s^\top \boldsymbol{\theta}_{0,s} + (\boldsymbol{\Lambda}^\top \mathbf{M}_{ns}^\top + \mathbf{Z}_{ns}^\top) (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns})\|_{\ell_2}^2.
\end{aligned}$$

We define the primary optimization loss as follows:

$$\mathcal{L}_{PO}(\boldsymbol{\theta}_s, \boldsymbol{\theta}_{ns}) := \frac{1}{2n} \|\varepsilon + \boldsymbol{\Lambda}^\top \mathbf{M}_s^\top (\boldsymbol{\theta}_{0,s} - \boldsymbol{\theta}_s) + \mathbf{Z}_s^\top \boldsymbol{\theta}_{0,s} + (\boldsymbol{\Lambda}^\top \mathbf{M}_{ns}^\top + \mathbf{Z}_{ns}^\top) (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns})\|_{\ell_2}^2 + \lambda \|\boldsymbol{\theta}_s\|_{\ell_2}^2 + \lambda \|\boldsymbol{\theta}_{ns}\|_{\ell_2}^2$$

We continue by deriving the auxiliary optimization (AO) problem. By duality, we have

$$\begin{aligned}
\mathcal{L}_{PO}(\boldsymbol{\theta}_s, \boldsymbol{\theta}_{ns}) &= \max_{\mathbf{v}} \frac{1}{n} \left(\mathbf{v}^\top \varepsilon + \mathbf{v}^\top \boldsymbol{\Lambda}^\top \mathbf{M}_s^\top (\boldsymbol{\theta}_{0,s} - \boldsymbol{\theta}_s) + \mathbf{v}^\top \mathbf{Z}_s^\top \boldsymbol{\theta}_{0,s} + \mathbf{v}^\top (\boldsymbol{\Lambda}^\top \mathbf{M}_{ns}^\top + \mathbf{Z}_{ns}^\top) (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}) - \frac{\|\mathbf{v}\|_{\ell_2}^2}{2} \right) \\
&\quad + \lambda \|\boldsymbol{\theta}_s\|_{\ell_2}^2 + \lambda \|\boldsymbol{\theta}_{ns}\|_{\ell_2}^2
\end{aligned}$$

Note that the above is jointly convex in $(\boldsymbol{\theta}_s, \boldsymbol{\theta}_{ns})$ and concave in \mathbf{v} , and the Gaussian matrix \mathbf{Z} is independent of everything else. Therefore, the AO problem reads:

$$\begin{aligned}
\mathcal{L}_{AO}(\boldsymbol{\theta}_s, \boldsymbol{\theta}_{ns}) &= \max_{\mathbf{v}} \frac{1}{n} \left(\mathbf{v}^\top \varepsilon + \mathbf{v}^\top \boldsymbol{\Lambda}^\top \mathbf{M}_s^\top (\boldsymbol{\theta}_{0,s} - \boldsymbol{\theta}_s) \right. \\
&\quad + \|\boldsymbol{\theta}_{0,s}\|_{\ell_2} \mathbf{g}_s^\top \mathbf{v} + \|\mathbf{v}\|_{\ell_2} \mathbf{h}_s^\top \boldsymbol{\theta}_{0,s} \\
&\quad + \|\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}\|_{\ell_2} \mathbf{g}_{ns}^\top \mathbf{v} + \|\mathbf{v}\|_{\ell_2} \mathbf{h}_{ns}^\top (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}) \\
&\quad \left. + \mathbf{v}^\top \boldsymbol{\Lambda}^\top \mathbf{M}_{ns}^\top (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}) - \frac{\|\mathbf{v}\|_{\ell_2}^2}{2} \right) + \lambda \|\boldsymbol{\theta}_s\|_{\ell_2}^2 + \lambda \|\boldsymbol{\theta}_{ns}\|_{\ell_2}^2,
\end{aligned}$$

where $\mathbf{g}_s, \mathbf{g}_{ns} \in \mathbb{R}^n$ and $\mathbf{h}_s \in \mathbb{R}^p$, $\mathbf{h}_{ns} \in \mathbb{R}^{d-p}$ are independent Gaussian random vectors with i.i.d $\mathcal{N}(0, 1)$ entries.

We next fix norm of $\|\mathbf{v}\|_{\ell_2} = \beta$, and maximize over its direction to obtain

$$\begin{aligned} \mathcal{L}_{AO}(\boldsymbol{\theta}_s, \boldsymbol{\theta}_{ns}) &= \max_{\beta \geq 0} \frac{1}{n} \left(\beta \left\| \boldsymbol{\varepsilon} + \boldsymbol{\Lambda}^\top \mathbf{M}_s^\top (\boldsymbol{\theta}_{0,s} - \boldsymbol{\theta}_s) + \|\boldsymbol{\theta}_{0,s}\|_{\ell_2} \mathbf{g}_s + \|\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}\|_{\ell_2} \mathbf{g}_{ns} + \boldsymbol{\Lambda}^\top \mathbf{M}_{ns}^\top (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}) \right\|_{\ell_2} \right. \\ &\quad \left. + \beta \mathbf{h}_s^\top \boldsymbol{\theta}_{0,s} + \beta \mathbf{h}_{ns}^\top (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}) - \frac{\beta^2}{2} \right) + \lambda \|\boldsymbol{\theta}_s\|_{\ell_2}^2 + \lambda \|\boldsymbol{\theta}_{ns}\|_{\ell_2}^2 \\ &= \max_{\beta \geq 0} \frac{1}{n} \left(\beta \left\| \boldsymbol{\varepsilon} + \boldsymbol{\Lambda}^\top \mathbf{M}_s^\top (\boldsymbol{\theta}_{0,s} - \boldsymbol{\theta}_s) + \boldsymbol{\Lambda}^\top \mathbf{M}_{ns}^\top (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}) + \sqrt{\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \|\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}\|_{\ell_2}^2} \mathbf{g} \right\|_{\ell_2} \right. \\ &\quad \left. + \beta \mathbf{h}_s^\top \boldsymbol{\theta}_{0,s} + \beta \mathbf{h}_{ns}^\top (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}) - \frac{\beta^2}{2} \right) + \lambda \|\boldsymbol{\theta}_s\|_{\ell_2}^2 + \lambda \|\boldsymbol{\theta}_{ns}\|_{\ell_2}^2, \end{aligned}$$

where we used that $\mathbf{g}_s, \mathbf{g}_{ns} \in \mathbb{R}^n$ have independent Gaussian entries. Here, $\mathbf{g} \in \mathbb{R}^n$ has i.i.d entries from $\mathcal{N}(0, 1)$. Next, note that the above optimization over β has a closed form. Using the identity $\max_{\beta \geq 0} (\beta x - \beta^2/2) = x_+^2/2$, with $x_+ = \max(x, 0)$, we get

$$\begin{aligned} \mathcal{L}_{AO}(\boldsymbol{\theta}_s, \boldsymbol{\theta}_{ns}) &= \frac{1}{2n} \left(\left\| \boldsymbol{\varepsilon} + \boldsymbol{\Lambda}^\top \mathbf{M}_s^\top (\boldsymbol{\theta}_{0,s} - \boldsymbol{\theta}_s) + \boldsymbol{\Lambda}^\top \mathbf{M}_{ns}^\top (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}) + \sqrt{\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \|\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}\|_{\ell_2}^2} \mathbf{g} \right\|_{\ell_2} \right. \\ &\quad \left. + \mathbf{h}_s^\top \boldsymbol{\theta}_{0,s} + \mathbf{h}_{ns}^\top (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}) \right)_+^2 + \lambda \|\boldsymbol{\theta}_s\|_{\ell_2}^2 + \lambda \|\boldsymbol{\theta}_{ns}\|_{\ell_2}^2. \end{aligned} \quad (\text{B.1})$$

Scalarization of the auxiliary optimization (AO) problem. We next proceed to scalarize the AO problem. Consider the singular value decomposition

$$\mathbf{M}_s = \mathbf{U}_s \boldsymbol{\Sigma}_s \mathbf{V}_s^\top,$$

with $\mathbf{U}_s \in \mathbb{R}^{p \times r}$, $\boldsymbol{\Sigma}_s \in \mathbb{R}^{r \times r}$, $\mathbf{V}_s \in \mathbb{R}^{k \times r}$, where $r = \text{rank}(\mathbf{M}_s) \leq k$. Decompose $\mathbf{q}_s := \boldsymbol{\theta}_{0,s} - \boldsymbol{\theta}_s$ in its projections onto the space spanned by the columns $\mathbf{u}_{1,s}, \dots, \mathbf{u}_{r,s}$ of \mathbf{U}_s , and the orthogonal component:

$$\mathbf{q}_s = \sum_{i=1}^r \alpha_i \mathbf{u}_{i,s} + \alpha_0 \mathbf{q}_s^\perp,$$

where $\|\mathbf{q}_s^\perp\|_{\ell_2} = 1$, $\alpha_0 \geq 0$, and $\mathbf{U}_s^\top \mathbf{q}_s^\perp = \mathbf{0}$. Using the shorthand $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_r)$, we write

$$\boldsymbol{\Lambda}^\top \mathbf{M}_s^\top (\boldsymbol{\theta}_{0,s} - \boldsymbol{\theta}_s) = \boldsymbol{\Lambda}^\top \mathbf{V}_s \boldsymbol{\Sigma}_s \mathbf{U}_s^\top \mathbf{q}_s = \boldsymbol{\Lambda}^\top \mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha}.$$

In addition,

$$\begin{aligned} \|\boldsymbol{\theta}_s\|_{\ell_2}^2 &= \|\boldsymbol{\theta}_{0,s} - (\boldsymbol{\theta}_{0,s} - \boldsymbol{\theta}_s)\|_{\ell_2}^2 \\ &= \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \|\mathbf{q}_s\|_{\ell_2}^2 - 2\langle \boldsymbol{\theta}_{0,s}, \mathbf{q}_s \rangle \\ &= \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \|\mathbf{q}_s\|_{\ell_2}^2 - 2\langle \boldsymbol{\theta}_{0,s}, \mathbf{U}_s \boldsymbol{\alpha} \rangle - 2\alpha_0 \langle \boldsymbol{\theta}_{0,s}, \mathbf{q}_s^\perp \rangle \\ &= \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \|\mathbf{q}_s\|_{\ell_2}^2 - 2\langle \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \boldsymbol{\alpha} \rangle - 2\alpha_0 \langle \boldsymbol{\theta}_{0,s}, \mathbf{q}_s^\perp \rangle \\ &= \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + (\alpha_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2) - 2\langle \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \boldsymbol{\alpha} \rangle - 2\alpha_0 \langle \boldsymbol{\theta}_{0,s}, \mathbf{q}_s^\perp \rangle. \end{aligned} \quad (\text{B.2})$$

Similarly, we define $\mathbf{q}_{ns} = \boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}$ and consider the singular value decomposition

$$\mathbf{M}_{ns} = \mathbf{U}_{ns} \boldsymbol{\Sigma}_{ns} \mathbf{V}_{ns}^\top,$$

with $\mathbf{U}_{ns} \in \mathbb{R}^{(d-p) \times t}$, $\boldsymbol{\Sigma}_{ns} \in \mathbb{R}^{t \times t}$, $\mathbf{V}_{ns} \in \mathbb{R}^{k \times t}$, where $t = \text{rank}(\mathbf{M}_{ns}) \leq k$. Decomposing \mathbf{q}_{ns} in its projections on the orthogonal columns $\mathbf{u}_{1,ns}, \dots, \mathbf{u}_{t,ns}$ of \mathbf{U}_{ns} , and the orthogonal component we write

$$\mathbf{q}_{ns} = \sum_{i=1}^t \gamma_i \mathbf{u}_{i,ns} + \gamma_0 \mathbf{q}_{ns}^\perp,$$

with $\|\mathbf{q}_{ns}^\perp\|_{\ell_2} = 1$, $\gamma_0 \geq 0$, and $\mathbf{U}_{ns}^\top \mathbf{q}_{ns}^\perp = \mathbf{0}$. Define $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_t)$. In this notation, we have

$$\boldsymbol{\Lambda}^\top \mathbf{M}_{ns}^\top (\boldsymbol{\theta}_{0,ns} - \boldsymbol{\theta}_{ns}) = \boldsymbol{\Lambda}^\top \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \mathbf{U}_{ns}^\top \mathbf{q}_{ns} = \boldsymbol{\Lambda}^\top \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma}.$$

Also, $\|\boldsymbol{\theta}_{0,\text{ns}} - \boldsymbol{\theta}_{\text{ns}}\|_{\ell_2} = \|\mathbf{q}_{\text{ns}}\|_{\ell_2} = \sqrt{\gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2}$. In addition,

$$\mathbf{h}_{\text{ns}}^\top (\boldsymbol{\theta}_{0,\text{ns}} - \boldsymbol{\theta}_{\text{ns}}) = \mathbf{h}_{\text{ns}}^\top \mathbf{q}_{\text{ns}} = \sum_{i=1}^t \gamma_i \mathbf{h}_{\text{ns}}^\top \mathbf{u}_{i,\text{ns}} + \gamma_0 \mathbf{h}_{\text{ns}}^\top \mathbf{q}_{\text{ns}}^\perp.$$

Using the above identities in (B.1), we have

$$\begin{aligned} \mathcal{L}_{AO}(\boldsymbol{\theta}_s, \boldsymbol{\theta}_{\text{ns}}) &= \frac{1}{2n} \left(\left\| \boldsymbol{\varepsilon} + \boldsymbol{\Lambda}^\top \mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \boldsymbol{\Lambda}^\top \mathbf{V}_{\text{ns}} \boldsymbol{\Sigma}_{\text{ns}} \boldsymbol{\gamma} + \sqrt{\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2} \mathbf{g} \right\|_{\ell_2} \right. \\ &\quad \left. + \mathbf{h}_s^\top \boldsymbol{\theta}_{0,s} + \sum_{i=1}^t \gamma_i \mathbf{h}_{\text{ns}}^\top \mathbf{u}_{i,\text{ns}} + \gamma_0 \mathbf{h}_{\text{ns}}^\top \mathbf{q}_{\text{ns}}^\perp \right)_+^2 \\ &\quad + \lambda \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \lambda(\alpha_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2) - 2\lambda \langle \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \boldsymbol{\alpha} \rangle - 2\lambda \alpha_0 \langle \boldsymbol{\theta}_{0,s}, \mathbf{q}_s^\perp \rangle \\ &\quad + \lambda \|\boldsymbol{\theta}_{0,\text{ns}}\|_{\ell_2}^2 + \lambda(\gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2) - 2\lambda \langle \mathbf{U}_{\text{ns}}^\top \boldsymbol{\theta}_{0,\text{ns}}, \boldsymbol{\gamma} \rangle - 2\lambda \gamma_0 \langle \boldsymbol{\theta}_{0,\text{ns}}, \mathbf{q}_{\text{ns}}^\perp \rangle. \end{aligned} \quad (\text{B.3})$$

By the above characterization, minimization over $\boldsymbol{\theta}_s$ and $\boldsymbol{\theta}_{\text{ns}}$ reduces to minimization over $\alpha_0, \gamma_0, \boldsymbol{\alpha}, \boldsymbol{\gamma}, \mathbf{q}_s^\perp$ and $\mathbf{q}_{\text{ns}}^\perp$. Further, these variables are free from each other and can be optimized over separately. For \mathbf{q}_s^\perp , there is only one term involving this variable and therefore, minimization over it reduces to

$$\min_{\mathbf{q}_s^\perp, \|\mathbf{q}_s^\perp\|_{\ell_2}=1} -\langle \boldsymbol{\theta}_{0,s}, \mathbf{q}_s^\perp \rangle = \min_{\mathbf{q}_s^\perp, \|\mathbf{q}_s^\perp\|_{\ell_2}=1} -\langle \mathbf{U}_s^\perp (\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s}, \mathbf{q}_s^\perp \rangle = -\|(\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s}\|_{\ell_2}.$$

For $\mathbf{q}_{\text{ns}}^\perp$, we note that there are two terms involving this variable, namely $\langle \frac{\mathbf{h}_{\text{ns}}}{\sqrt{n}}, \mathbf{q}_{\text{ns}}^\perp \rangle$ and $\langle (\mathbf{U}_{\text{ns}}^\perp)^\top \boldsymbol{\theta}_{0,\text{ns}}, \mathbf{q}_{\text{ns}}^\perp \rangle$. Since $\|\mathbf{q}_{\text{ns}}^\perp\|_{\ell_2} = 1$, it is easy to see that the optimal $\mathbf{q}_{\text{ns}}^\perp$ should be in the span of $\mathbf{h}_{\text{ns}}^\perp$ and $(\mathbf{U}_{\text{ns}}^\perp)^\top \boldsymbol{\theta}_{0,\text{ns}}$. In addition,

$$\left\langle \frac{\mathbf{h}_{\text{ns}}^\perp}{\sqrt{n}}, (\mathbf{U}_{\text{ns}}^\perp)^\top \boldsymbol{\theta}_{0,\text{ns}} \right\rangle \xrightarrow{(p)} 0,$$

by the law of large numbers. In words, these two vectors are asymptotically orthogonal. Hence, we can consider the following decomposition of the optimal $\mathbf{q}_{\text{ns}}^\perp$:

$$\mathbf{q}_{\text{ns}}^\perp = -\xi \frac{\mathbf{h}_{\text{ns}}^\perp}{\|\mathbf{h}_{\text{ns}}^\perp\|_{\ell_2}} + \sqrt{1 - \xi^2} \frac{\mathbf{U}_{\text{ns}}^\perp (\mathbf{U}_{\text{ns}}^\perp)^\top \boldsymbol{\theta}_{0,\text{ns}}}{\|(\mathbf{U}_{\text{ns}}^\perp)^\top \boldsymbol{\theta}_{0,\text{ns}}\|_{\ell_2}},$$

where $\xi \geq 0$ and $\mathbf{h}_{\text{ns}}^\perp$ denotes the projection of \mathbf{h}_{ns} onto the (left) null space of \mathbf{U}_{ns} . This brings us to

$$\begin{aligned} \min_{\boldsymbol{\theta}_s, \boldsymbol{\theta}_{\text{ns}}} \mathcal{L}_{AO}(\boldsymbol{\theta}_s, \boldsymbol{\theta}_{\text{ns}}) &= \min_{\alpha_0, \gamma_0 \geq 0, \boldsymbol{\alpha}, \boldsymbol{\gamma}} \frac{1}{2} \left(\frac{1}{\sqrt{n}} \left\| \boldsymbol{\varepsilon} + \boldsymbol{\Lambda}^\top \mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \boldsymbol{\Lambda}^\top \mathbf{V}_{\text{ns}} \boldsymbol{\Sigma}_{\text{ns}} \boldsymbol{\gamma} + \sqrt{\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2} \mathbf{g} \right\|_{\ell_2} \right. \\ &\quad \left. + \frac{\mathbf{h}_s^\top \boldsymbol{\theta}_{0,s}}{\sqrt{n}} + \sum_{i=1}^t \gamma_i \frac{\mathbf{h}_{\text{ns}}^\top \mathbf{u}_{i,\text{ns}}}{\sqrt{n}} - \gamma_0 \xi \frac{\|\mathbf{h}_{\text{ns}}^\perp\|_{\ell_2}}{\sqrt{n}} \right)_+^2 \\ &\quad + \lambda \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \lambda(\alpha_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2) - 2\lambda \langle \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \boldsymbol{\alpha} \rangle - 2\lambda \alpha_0 \|(\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s}\|_{\ell_2} \\ &\quad + \lambda \|\boldsymbol{\theta}_{0,\text{ns}}\|_{\ell_2}^2 + \lambda(\gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2) - 2\lambda \langle \mathbf{U}_{\text{ns}}^\top \boldsymbol{\theta}_{0,\text{ns}}, \boldsymbol{\gamma} \rangle - 2\lambda \gamma_0 \sqrt{1 - \xi^2} \|(\mathbf{U}_{\text{ns}}^\perp)^\top \boldsymbol{\theta}_{0,\text{ns}}\|_{\ell_2}. \end{aligned} \quad (\text{B.4})$$

Note that at this stage, the AO problem is reduced to an optimization over $r + t + 3$ scalar variables ($\alpha_0, \gamma_0 \geq 0, 0 \leq \xi \leq 1$ and $\boldsymbol{\alpha} \in \mathbb{R}^r, \boldsymbol{\gamma} \in \mathbb{R}^t$).

Convergence of the auxiliary optimization problem. We next continue to derive the point-wise in-probability limit of the AO problem.

First observe that since $\boldsymbol{\varepsilon}$ and \mathbf{g} are independent with i.i.d $\mathcal{N}(0, 1)$ entries, we have

$$\boldsymbol{\varepsilon} + \sqrt{\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2} \mathbf{g} \stackrel{(d)}{=} \sqrt{\sigma^2 + 2(\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2)} \tilde{\mathbf{g}},$$

where $\tilde{\mathbf{g}} \in \mathbb{R}^n$ has i.i.d $\mathcal{N}(0, 1)$ entries.

Second, by construction $\boldsymbol{\Lambda} \boldsymbol{\Lambda}^\top = \text{diag}(n_1, \dots, n_k) \in \mathbb{R}^{k \times k}$, where n_ℓ denotes the number of examples from cluster ℓ . Hence,

$$\begin{aligned} \frac{1}{n} \left\| \boldsymbol{\Lambda}^\top \mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \boldsymbol{\Lambda}^\top \mathbf{V}_{\text{ns}} \boldsymbol{\Sigma}_{\text{ns}} \boldsymbol{\gamma} \right\|_{\ell_2}^2 &= (\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{\text{ns}} \boldsymbol{\Sigma}_{\text{ns}} \boldsymbol{\gamma})^\top \text{diag}\left(\frac{n_1}{n}, \dots, \frac{n_k}{n}\right) (\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{\text{ns}} \boldsymbol{\Sigma}_{\text{ns}} \boldsymbol{\gamma}) \\ &\stackrel{(p)}{\rightarrow} (\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{\text{ns}} \boldsymbol{\Sigma}_{\text{ns}} \boldsymbol{\gamma})^\top \text{diag}(\boldsymbol{\pi}) (\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{\text{ns}} \boldsymbol{\Sigma}_{\text{ns}} \boldsymbol{\gamma}) \end{aligned}$$

Next, by using concentration of Lipschitz functions of Gaussian vectors, we obtain

$$\begin{aligned} & \frac{1}{\sqrt{n}} \left\| \boldsymbol{\varepsilon} + \boldsymbol{\Lambda}^\top \mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \boldsymbol{\Lambda}^\top \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma} + \sqrt{\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2} \mathbf{g} \right\|_{\ell_2} \\ & \xrightarrow{p} \sqrt{(\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma})^\top \text{diag}(\boldsymbol{\pi})(\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma}) + \sigma^2 + (\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2)} \end{aligned}$$

Also, since $\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}$ is bounded and $\|\mathbf{u}_{i,s}\|_{\ell_2} = 1$, we get

$$\frac{\mathbf{h}_s^\top \boldsymbol{\theta}_{0,s}}{\sqrt{n}}, \frac{\mathbf{h}_{ns}^\top \mathbf{u}_{i,ns}}{\sqrt{n}} \xrightarrow{p} 0.$$

In addition, $\|\mathbf{h}_{ns}^\perp\|_{\ell_2}$ concentrates around $\sqrt{d-p-t}$ and $(d-p-t)/n \rightarrow \psi_d - \psi_p$, because $t \leq k$ remains bounded as n diverges, and so

$$\frac{\|\mathbf{h}_{ns}^\perp\|_{\ell_2}}{\sqrt{n}} \xrightarrow{p} \sqrt{\psi_d - \psi_p}.$$

Using the above limits, the objective in (B.4) converges in-probability to

$$\mathcal{D}(\alpha_0, \gamma_0, \xi, \boldsymbol{\alpha}, \boldsymbol{\gamma}) :=$$

$$\begin{aligned} & \frac{1}{2} \left(\sqrt{(\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma})^\top \text{diag}(\boldsymbol{\pi})(\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma}) + \sigma^2 + (\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2)} - \gamma_0 \xi \sqrt{\psi_d - \psi_p} \right)_+^2 \\ & + \lambda \|\boldsymbol{\theta}_0\|_{\ell_2}^2 + \lambda(\alpha_0^2 + \gamma_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2) \\ & - 2\lambda \left(\langle \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \boldsymbol{\alpha} \rangle + \alpha_0 \left\| (\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2} + \langle \mathbf{U}_{ns}^\top \boldsymbol{\theta}_{0,ns}, \boldsymbol{\gamma} \rangle + \gamma_0 \sqrt{1 - \xi^2} \left\| (\mathbf{U}_{ns}^\perp)^\top \boldsymbol{\theta}_{0,ns} \right\|_{\ell_2} \right) \quad (\text{B.5}) \end{aligned}$$

We are now ready to prove the theorems.

B.2.1 Proof of Theorem 3.1

Using Lemma 2.1, we have

$$\begin{aligned} \text{Risk}(\widehat{\boldsymbol{\theta}}_L) &= \sigma^2 + \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}\|_{\ell_2}^2 + (\boldsymbol{\theta}_0 - \boldsymbol{\theta})^\top \mathbf{M} \text{diag}(\boldsymbol{\pi}) \mathbf{M}^\top (\boldsymbol{\theta}_0 - \boldsymbol{\theta}) \\ &= \sigma^2 + \|\mathbf{q}_s\|_{\ell_2}^2 + \|\mathbf{q}_{ns}\|_{\ell_2}^2 + \mathbf{q}_s^\top \mathbf{M}_s \text{diag}(\boldsymbol{\pi}) \mathbf{M}_s^\top \mathbf{q}_s + \mathbf{q}_{ns}^\top \mathbf{M}_{ns} \text{diag}(\boldsymbol{\pi}) \mathbf{M}_{ns}^\top \mathbf{q}_{ns} \\ &= \sigma^2 + (\alpha_0^2 + \gamma_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2) \\ & \quad + \boldsymbol{\alpha}^\top \boldsymbol{\Sigma}_s \mathbf{V}_s^\top \text{diag}(\boldsymbol{\pi}) \mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \boldsymbol{\gamma}^\top \boldsymbol{\Sigma}_{ns} \mathbf{V}_{ns}^\top \text{diag}(\boldsymbol{\pi}) \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma}. \quad (\text{B.6}) \end{aligned}$$

Since $\psi_d - \psi_p \leq 1$, we are in the over- determined (a.k.a underparametrized) regime. As $\lambda \rightarrow 0^+$, the terms involving λ become negligible compared to the first term in (B.5) except those that include α_0 , as α_0 is not present in the first term. Since $(x)_+^2$ is increasing, and

$$(\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma})^\top \text{diag}(\boldsymbol{\pi})(\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma}) + \|\boldsymbol{\gamma}\|_{\ell_2}^2 \geq 0,$$

the minimum over $\boldsymbol{\alpha}$ and $\boldsymbol{\gamma}$ is achieved for $\boldsymbol{\alpha} = \mathbf{0} \in \mathbb{R}^r$ and $\boldsymbol{\gamma} = \mathbf{0} \in \mathbb{R}^t$. The optimization (B.5) then reduces to

$$\min_{\alpha_0, \gamma_0 \geq 0, 0 \leq \xi \leq 1} \frac{1}{2} \left(\sqrt{\sigma^2 + (\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2)} - \gamma_0 \xi \sqrt{\psi_d - \psi_p} \right)_+^2 + \lambda \alpha_0^2 - 2\lambda \alpha_0 \left\| (\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2}. \quad (\text{B.7})$$

The optimal ξ is given by $\xi = 1$. Also, setting derivative with respect to α_0 to zero we obtain the optimal $\alpha_0 = \left\| (\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2}$. Next, by setting derivative with respect to γ_0 we arrive at

$$\gamma_0^2 = (\sigma^2 + \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2) \frac{\psi_d - \psi_p}{1 - (\psi_d - \psi_p)}.$$

Using the optimal variables in (B.6) we obtain the risk of minimum-norm estimator as

$$\begin{aligned} \text{Risk}(\widehat{\boldsymbol{\theta}}_L) &= \sigma^2 + \left\| (\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2}^2 + (\sigma^2 + \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2) \frac{\psi_d - \psi_p}{1 - (\psi_d - \psi_p)} \\ &= (\sigma^2 + \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2) \frac{1}{1 - (\psi_d - \psi_p)} - \left\| \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2}^2. \end{aligned}$$

Recall that by assumption, $r_s = \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}$ and $\left\| \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2} = \sqrt{\rho} r_s$, which completes the proof.

B.2.2 Proof of Theorem 3.2

We continue from (B.5). In the case of $\psi_d - \psi_p \leq 1$, it is easy to see that the derivative of the first term of (B.5), in the active region is decreasing in γ_0 . With the consideration $\lambda \rightarrow 0^+$, minimizing over γ_0 will push us into the non-active region. Therefore the optimization problem (B.5) reduces to

$$\begin{aligned} \text{minimize} \quad & \|\boldsymbol{\theta}_0\|_{\ell_2}^2 + \alpha_0^2 + \gamma_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2 \\ & - 2 \left(\langle \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \boldsymbol{\alpha} \rangle + \alpha_0 \left\| (\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2} + \langle \mathbf{U}_{ns}^\top \boldsymbol{\theta}_{0,ns}, \boldsymbol{\gamma} \rangle + \gamma_0 \sqrt{1 - \xi^2} \left\| (\mathbf{U}_{ns}^\perp)^\top \boldsymbol{\theta}_{0,ns} \right\|_{\ell_2} \right) \end{aligned}$$

subject to

$$(\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma})^\top \text{diag}(\boldsymbol{\pi}) (\mathbf{V}_s \boldsymbol{\Sigma}_s \boldsymbol{\alpha} + \mathbf{V}_{ns} \boldsymbol{\Sigma}_{ns} \boldsymbol{\gamma}) + \sigma^2 + (\|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2) \leq \gamma_0^2 \xi^2 (\psi_d - \psi_p) \quad (\text{B.8})$$

By Assumption 2, $\boldsymbol{\Sigma}_s = \mu \mathbf{I}_k$, $\mathbf{V}_s = \mathbf{I}_{k \times k}$, and $\boldsymbol{\Sigma}_{ns} = \mathbf{0}$, $\mathbf{U}_{ns} = \mathbf{0}$ (no cluster structure on non-sensitive features and an orthogonal, equal energy cluster centers on the sensitive features). Therefore, by fixing $\boldsymbol{\gamma} := \|\boldsymbol{\gamma}\|_{\ell_2}$, the optimization problem (B.8) becomes:

$$\text{minimize} \quad \alpha_0^2 + \gamma_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2 + \gamma^2 - 2 \left(\langle \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \boldsymbol{\alpha} \rangle + \alpha_0 \left\| (\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2} + \gamma_0 \sqrt{1 - \xi^2} \|\boldsymbol{\theta}_{0,ns}\|_{\ell_2} \right)$$

subject to

$$\mu^2 \boldsymbol{\alpha}^\top \text{diag}(\boldsymbol{\pi}) \boldsymbol{\alpha} + \sigma^2 + \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2 + \gamma^2 \leq \gamma_0^2 \xi^2 (\psi_d - \psi_p). \quad (\text{B.9})$$

Since α_0 does not appear in the constraint, it is easy to see that its optimal value is given by $\alpha_0 = \left\| (\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2}$. Also, note that by decreasing γ the objective value decreases and also by the constraint on the other variables become more relaxed. Consequently, the optimal value of γ is $\gamma = 0$. Removing α_0 from the objective function, we are left with

$$\begin{aligned} \text{minimize} \quad & \gamma_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2 - 2 \left(\langle \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \boldsymbol{\alpha} \rangle + \gamma_0 \sqrt{1 - \xi^2} \|\boldsymbol{\theta}_{0,ns}\|_{\ell_2} \right) \\ \text{subject to} \quad & \mu^2 \boldsymbol{\alpha}^\top \text{diag}(\boldsymbol{\pi}) \boldsymbol{\alpha} + \sigma^2 + \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2 \leq \gamma_0^2 \xi^2 (\psi_d - \psi_p). \end{aligned} \quad (\text{B.10})$$

Optimal choice of ξ results in the constraint to become equality. Solving for ξ , the optimization reduces to

$$\text{minimize} \quad \gamma_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2 - 2 \left(\langle \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}, \boldsymbol{\alpha} \rangle + \sqrt{\gamma_0^2 - \frac{\mu^2 \boldsymbol{\alpha}^\top \text{diag}(\boldsymbol{\pi}) \boldsymbol{\alpha} + \sigma^2 + \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2}{\psi_d - \psi_p}} \|\boldsymbol{\theta}_{0,ns}\|_{\ell_2} \right)$$

Setting derivative with respect to γ_0 to zero, we obtain

$$\sqrt{\gamma_0^2 - \frac{\mu^2 \boldsymbol{\alpha}^\top \text{diag}(\boldsymbol{\pi}) \boldsymbol{\alpha} + \sigma^2 + \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \gamma_0^2}{\psi_d - \psi_p}} = \left(1 - \frac{1}{\psi_d - \psi_p} \right) \|\boldsymbol{\theta}_{0,ns}\|_{\ell_2}. \quad (\text{B.11})$$

Setting derivative with respect to $\boldsymbol{\alpha}$ to zero and using the previous stationary equation, we get

$$\boldsymbol{\alpha} = \left(\mathbf{I} + \frac{\mu^2 \text{diag}(\boldsymbol{\pi})}{\psi_d - \psi_p - 1} \right)^{-1} \mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}. \quad (\text{B.12})$$

We next square both sides of (B.12) and rearrange the terms to get

$$\begin{aligned} \gamma_0^2 &= \frac{1}{\psi_d - \psi_p - 1} \left(\sigma^2 + \|\boldsymbol{\theta}_{0,s}\|_{\ell_2}^2 + \mu^2 \boldsymbol{\alpha}^\top \text{diag}(\boldsymbol{\pi}) \boldsymbol{\alpha} \right) + \left(1 - \frac{1}{\psi_d - \psi_p} \right) \|\boldsymbol{\theta}_{0,ns}\|_{\ell_2}^2 \\ &= \frac{1}{\psi_d - \psi_p - 1} \left(\sigma^2 + r_s^2 + \mu^2 \boldsymbol{\alpha}^\top \text{diag}(\boldsymbol{\pi}) \boldsymbol{\alpha} \right) + \left(1 - \frac{1}{\psi_d - \psi_p} \right) r_{ns}^2, \end{aligned}$$

which are the same expressions for $\boldsymbol{\alpha}$ and γ_0 given in the theorem statement.

The final step is to write the risk of estimator in terms of $\boldsymbol{\alpha}$, γ_0 . Invoke equation (B.6), and recall that in the current case, $\boldsymbol{\Sigma}_{ns} = \mathbf{0}$, $\boldsymbol{\Sigma}_s = \mu \mathbf{I}$. Also, as we showed in our derivation, $\boldsymbol{\gamma} = \|\boldsymbol{\gamma}\|_{\ell_2} = 0$, $\alpha_0 = \left\| (\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2}$, by which we arrive at

$$\begin{aligned} \text{Risk}(\widehat{\boldsymbol{\theta}}_L) &= \mu^2 \boldsymbol{\alpha}^\top \text{diag}(\boldsymbol{\pi}) \boldsymbol{\alpha} + \sigma^2 + \left(\left\| (\mathbf{U}_s^\perp)^\top \boldsymbol{\theta}_{0,s} \right\|_{\ell_2}^2 + \gamma_0^2 + \|\boldsymbol{\alpha}\|_{\ell_2}^2 \right) \\ &= \sigma^2 + (1 - \rho) r_s^2 + \gamma_0^2 + \boldsymbol{\alpha}^\top \left(\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi}) \right) \boldsymbol{\alpha}. \end{aligned} \quad (\text{B.13})$$

This concludes the proof.

B.3 Proof of Theorem 3.3

We follow the proof strategy used for Theorem 3.1-3.2. Here, we would like to characterize the risk of min-norm estimator $\widehat{\boldsymbol{\theta}}$. The features matrix has a clustering structure, but the learner is not using that (no look-alike clustering) and is just compute the min-norm estimator for fitting the responses to individual features. Therefore, one can think of this setting as a special case of our previous analysis when there is no sensitive features (so $\psi_p = 0$).

(a) By setting $\psi_p = 0$ and $r_s = 0$ in the result of Theorem 3.1, we get that when $\psi_d \leq 1$,

$$\text{Risk}(\widehat{\boldsymbol{\theta}}) = \frac{\sigma^2}{1 - \psi_d}.$$

(b) In this case, we specialize the proof of Theorem 3.2 to the case that $\psi_p = 0$. Continuing from (B.8), and removing the terms corresponding to sensitive features, we arrive at

$$\begin{aligned} & \text{minimize} \quad \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2 - 2 \left(\langle \mathbf{U}_{\text{ns}}^\top \boldsymbol{\theta}_{0,\text{ns}}, \boldsymbol{\gamma} \rangle + \gamma_0 \sqrt{1 - \xi^2} \|\mathbf{U}_{\text{ns}}^\perp \boldsymbol{\theta}_{0,\text{ns}}\|_{\ell_2} \right) \\ & \text{subject to} \\ & \quad (\mathbf{V}_{\text{ns}} \boldsymbol{\Sigma}_{\text{ns}} \boldsymbol{\gamma})^\top \text{diag}(\boldsymbol{\pi}) (\mathbf{V}_{\text{ns}} \boldsymbol{\Sigma}_{\text{ns}} \boldsymbol{\gamma}) + \sigma^2 + \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2 \leq \gamma_0^2 \xi^2 \psi_d \end{aligned} \quad (\text{B.14})$$

We drop the index ‘ns’ as it is not relevant in this case. Also by Assumption 2, $\boldsymbol{\Sigma}_{\text{ns}} = \mu \mathbf{I}_d$, $\mathbf{V}_{\text{ns}} = \mathbf{I}_d$. Therefore, the above optimization can be written as

$$\begin{aligned} & \text{minimize} \quad \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2 - 2 \left(\langle \mathbf{U}^\top \boldsymbol{\theta}_0, \boldsymbol{\gamma} \rangle + \gamma_0 \sqrt{1 - \xi^2} \|(\mathbf{U}^\perp)^\top \boldsymbol{\theta}_0\|_{\ell_2} \right) \\ & \text{subject to} \quad \boldsymbol{\gamma}^\top (\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi})) \boldsymbol{\gamma} + \sigma^2 + \gamma_0^2 \leq \gamma_0^2 \xi^2 \psi_d. \end{aligned} \quad (\text{B.15})$$

Optimal ξ makes the constraint equality. Solving for ξ , the above optimization can be written as so we have

$$\text{minimize} \quad \gamma_0^2 + \|\boldsymbol{\gamma}\|_{\ell_2}^2 - 2 \left(\langle \mathbf{U}^\top \boldsymbol{\theta}_0, \boldsymbol{\gamma} \rangle + \sqrt{\gamma_0^2 - \frac{\boldsymbol{\gamma}^\top (\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi})) \boldsymbol{\gamma} + \sigma^2 + \gamma_0^2}{\psi_d}} \|(\mathbf{U}^\perp)^\top \boldsymbol{\theta}_0\|_{\ell_2} \right).$$

Setting the derivative with respect to γ_0 to zero, we get

$$\sqrt{\gamma_0^2 - \frac{\boldsymbol{\gamma}^\top (\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi})) \boldsymbol{\gamma} + \sigma^2 + \gamma_0^2}{\psi_d}} = \left(1 - \frac{1}{\psi_d}\right) \|(\mathbf{U}^\perp)^\top \boldsymbol{\theta}_0\|_{\ell_2}. \quad (\text{B.16})$$

Setting derivative with respect to $\boldsymbol{\gamma}$ to zero and using the above equation, we obtain

$$\boldsymbol{\gamma} = \left(\mathbf{I} + \frac{\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi})}{\psi_d - 1} \right)^{-1} \mathbf{U}^\top \boldsymbol{\theta}_0. \quad (\text{B.17})$$

We next square both sides of equation (B.16), and rearrange the terms to get:

$$\gamma_0^2 = \frac{1}{\psi_d - 1} \left(\sigma^2 + \boldsymbol{\gamma}^\top (\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi})) \boldsymbol{\gamma} \right) + \left(1 - \frac{1}{\psi_d}\right) \|(\mathbf{U}^\perp)^\top \boldsymbol{\theta}_0\|_{\ell_2}^2.$$

Under the simplifying Assumption 2, there is no cluster structure on the non-sensitive features and so $\mathbf{U}_{\text{ns}} = 0$. Therefore,

$$\begin{aligned} \|\mathbf{U}^\top \boldsymbol{\theta}_0\|_{\ell_2} &= \|\mathbf{U}_s^\top \boldsymbol{\theta}_{0,s}\|_{\ell_2} = \sqrt{\rho} r_s, \\ \|(\mathbf{U}^\perp)^\top \boldsymbol{\theta}_0\|_{\ell_2}^2 &= \|\boldsymbol{\theta}_0\|_{\ell_2}^2 - \|\mathbf{U}^\top \boldsymbol{\theta}_0\|_{\ell_2}^2 = (1 - \rho) r_s^2 + r_{\text{ns}}^2. \end{aligned}$$

We next proceed to compute the risk of estimator in terms of $\boldsymbol{\gamma}$, γ_0 . We use equation (B.6), which for the min-norm estimator with no look-alike clustering, reduces to

$$\text{Risk}(\widehat{\boldsymbol{\theta}}) = \sigma^2 + \gamma_0^2 + \boldsymbol{\gamma}^\top (\mathbf{I} + \mu^2 \text{diag}(\boldsymbol{\pi})) \boldsymbol{\gamma}. \quad (\text{B.18})$$

This concludes the proof. Note that in the theorem statement we made the change of variables $\gamma_0 \rightarrow \tilde{\gamma}_0$ and $\boldsymbol{\gamma} \rightarrow \tilde{\boldsymbol{\alpha}}$, for an easier comparison with the risk of look-alike estimator.)

B.4 Proof of Proposition 3.4

Consider singular value decompositions $\mathbf{X}_L = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ and $\tilde{\mathbf{X}}_L = \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}\tilde{\mathbf{V}}^T$. We then can write the estimators $\hat{\boldsymbol{\theta}}_L$ and $\tilde{\boldsymbol{\theta}}_L$ as follows:

$$\hat{\boldsymbol{\theta}}_L = \mathbf{U}\mathbf{\Sigma}^{-1}\mathbf{V}^T\mathbf{y}, \quad \tilde{\boldsymbol{\theta}}_L = \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}^{-1}\tilde{\mathbf{V}}^T\mathbf{y}.$$

We first bound $\|\hat{\boldsymbol{\theta}}_L - \tilde{\boldsymbol{\theta}}_L\|$. We write

$$\|\hat{\boldsymbol{\theta}}_L - \tilde{\boldsymbol{\theta}}_L\| \leq \|\mathbf{U}\mathbf{\Sigma}^{-1}\mathbf{V}^T - \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}^{-1}\tilde{\mathbf{V}}^T\|\|\mathbf{y}\|. \quad (\text{B.19})$$

We have

$$\|\mathbf{y}\| = \|\mathbf{X}^T\boldsymbol{\theta}_0 + \boldsymbol{\varepsilon}\| = \|\boldsymbol{\Lambda}^T\mathbf{M}^T\boldsymbol{\theta}_0 + \mathbf{Z}^T\boldsymbol{\theta}_0 + \boldsymbol{\varepsilon}\|.$$

Note that $\mathbf{Z}^T\boldsymbol{\theta}_0 + \boldsymbol{\varepsilon} \stackrel{(d)}{=} \sqrt{\|\boldsymbol{\theta}_0\|^2 + \sigma^2}\mathbf{g}$ where $\mathbf{g} \sim \mathcal{N}(0, I_n)$. In addition,

$$\begin{aligned} \frac{1}{n}\|\boldsymbol{\Lambda}^T\mathbf{M}^T\boldsymbol{\theta}_0\|^2 &= \frac{1}{n}\boldsymbol{\theta}_0^T\mathbf{M}\boldsymbol{\Lambda}\boldsymbol{\Lambda}^T\mathbf{M}^T\boldsymbol{\theta}_0 \\ &= \boldsymbol{\theta}_0^T\mathbf{M}\text{diag}\left(\frac{n_1}{n}, \dots, \frac{n_k}{n}\right)\mathbf{M}^T\boldsymbol{\theta}_0 \xrightarrow{p} \boldsymbol{\theta}_0^T\mathbf{M}\text{diag}(\pi_1, \dots, \pi_k)\mathbf{M}^T\boldsymbol{\theta}_0. \end{aligned}$$

Therefore by using concentration of Lipschitz functions of Gaussian vectors, we get

$$\frac{1}{\sqrt{n}}\|\mathbf{y}\| \xrightarrow{p} \sqrt{\boldsymbol{\theta}_0^T\mathbf{M}\text{diag}(\boldsymbol{\pi})\mathbf{M}^T\boldsymbol{\theta}_0 + \|\boldsymbol{\theta}_0\|^2 + \sigma^2}.$$

This shows that

$$\frac{1}{\sqrt{n}}\|\mathbf{y}\| \rightarrow C \leq \sqrt{(\mu+1)(r_s^2 + r_{ns}^2) + \sigma^2}. \quad (\text{B.20})$$

We next use the result of [28, Theorem 3.3], by which we obtain

$$\|\mathbf{U}\mathbf{\Sigma}^{-1}\mathbf{V}^T - \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}^{-1}\tilde{\mathbf{V}}^T\| \leq \frac{1+\sqrt{5}}{2} \max\left(\frac{1}{\sigma_{\min}(\boldsymbol{\Sigma})^2}, \frac{1}{\sigma_{\min}(\tilde{\boldsymbol{\Sigma}})^2}\right) \|\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T - \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}\tilde{\mathbf{V}}^T\|. \quad (\text{B.21})$$

Note that

$$\|\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T - \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}\tilde{\mathbf{V}}^T\| = \|\mathbf{X}_L - \tilde{\mathbf{X}}_L\| = \|\mathbf{M}_s\boldsymbol{\Lambda} - \tilde{\mathbf{M}}_s\tilde{\boldsymbol{\Lambda}}\| \leq \delta\sqrt{n}, \quad (\text{B.22})$$

by the assumption of the theorem statement. We next lower bound $\sigma_{\min}(\boldsymbol{\Sigma}) = \sigma_{\min}(\mathbf{X}_L)$. Recall that $\mathbf{X}_L^T = (\mathbf{M}\boldsymbol{\Lambda})^T + [\mathbf{0}_{n \times p}, \mathbf{Z}_{n \times (d-p)}]$, with \mathbf{Z} having i.i.d $\mathcal{N}(0, 1)$ entries.

Next suppose that Condition (i) holds true, namely $\delta < \sqrt{1 - (\psi_d - \psi_p)} - \sqrt{\psi_d - \psi_p}$, with $\psi_d - \psi_p < 0.5$. Using the result of [31, Theorem 2.1], we have with probability at least $1 - n^{-1}$,

$$\sigma_{\min}(\mathbf{X}_L) \geq \sqrt{n} \left(\sqrt{\psi_d - \psi_p - 1} - 1 - \sqrt{\frac{2 \log n}{n}} \right).$$

Furthermore,

$$\begin{aligned} \sigma_{\min}(\tilde{\mathbf{X}}_L) &\geq \sigma_{\min}(\mathbf{X}_L) - \|\mathbf{X}_L - \tilde{\mathbf{X}}_L\| \\ &\geq \sqrt{n} \left(\sqrt{1 - (\psi_d - \psi_p)} - \sqrt{\psi_d - \psi_p} - \sqrt{\frac{2 \log n}{n}} - \delta \right) \\ &\geq c' \sqrt{n} \left(\sqrt{1 - (\psi_d - \psi_p)} - \sqrt{\psi_d - \psi_p} \right), \end{aligned}$$

using the assumption on the estimation error rate δ . Therefore, using the above bound along with (B.22) in (B.21) we get

$$\|\mathbf{U}\mathbf{\Sigma}^{-1}\mathbf{V}^T - \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}^{-1}\tilde{\mathbf{V}}^T\| \leq \frac{1+\sqrt{5}}{2c'^2} \frac{1}{\sqrt{n} \left(\sqrt{1 - (\psi_d - \psi_p)} - \sqrt{\psi_d - \psi_p} \right)^2} \delta.$$

Combining the above bound with (B.20), we get

$$\|\widehat{\boldsymbol{\theta}}_L - \widetilde{\boldsymbol{\theta}}_L\| \leq \frac{1 + \sqrt{5}}{2c^2} \frac{C}{\left(\sqrt{1 - (\psi_d - \psi_p)} - \sqrt{\psi_d - \psi_p}\right)^2} \delta. \quad (\text{B.23})$$

We next note that by triangle inequality, the above bound implies that

$$\|\widetilde{\boldsymbol{\theta}}_L - \boldsymbol{\theta}_0\| - \|\widehat{\boldsymbol{\theta}}_L - \boldsymbol{\theta}_0\| \leq \|\widehat{\boldsymbol{\theta}}_L - \widetilde{\boldsymbol{\theta}}_L\| = O(\delta).$$

Therefore, by invoking Lemma 2.1, we obtain the desired result on $\text{Risk}(\widehat{\boldsymbol{\theta}}_L)$.

Next suppose that Condition (ii) holds, namely $\delta < \sqrt{\psi_d - \psi_p - 1} - 1$ with $\psi_d - \psi_p > 2$. Using the result of [31, Theorem 2.1] for \mathbf{X}^\top , we have with probability at least $1 - n^{-1}$,

$$\sigma_{\min}(\mathbf{X}_L) \geq \sqrt{n} \left(\sqrt{\psi_d - \psi_p - 1} - 1 - \sqrt{\frac{2 \log n}{n}} \right).$$

By following a similar argument we prove the claim under Condition (ii).

B.5 Proof of Theorem 5.1

We use Theorem 3.3 (b) to characterize $\text{Risk}(\widehat{\boldsymbol{\theta}})$ in the regime of $\psi_d \geq 1$. Specializing to the case of balanced cluster priors, the risk depends on $\tilde{\alpha}$ only through its norm $\tilde{\alpha} := \|\tilde{\boldsymbol{\alpha}}\|_{\ell_2}$, and is given by

$$\begin{aligned} \text{Risk}(\widehat{\boldsymbol{\theta}}) &\stackrel{\mathcal{P}}{\rightarrow} \sigma^2 + \tilde{\gamma}_0^2 + \left(\frac{\mu^2}{k} + 1 \right) \tilde{\alpha}^2 \\ &= \frac{\psi_d}{\psi_d - 1} \left(\sigma^2 + \left(\frac{\mu^2}{k} + 1 \right) \tilde{\alpha}^2 \right) + \left(1 - \frac{1}{\psi_d} \right) \left((1 - \rho)r_s^2 + r_{\text{ns}}^2 \right), \end{aligned}$$

with

$$\tilde{\alpha} = \left(1 + \frac{\frac{\mu^2}{k} + 1}{\psi_d - 1} \right)^{-1} \sqrt{\rho} r_s.$$

In addition, by Theorem 3.1 we have

$$\text{Risk}(\widehat{\boldsymbol{\theta}}_L) \stackrel{\mathcal{P}}{\rightarrow} \frac{\sigma^2 + r_s^2}{1 - \psi_d + \psi_p} - \rho r_s^2.$$

Note that $\text{Risk}(\widehat{\boldsymbol{\theta}}_L)$ in this regime does not depend on μ^2/k . Also, it is easy to verify that $\text{Risk}(\widehat{\boldsymbol{\theta}})$ is decreasing in μ^2/k . Therefore the gain Δ is decreasing in μ^2/k .

Also observe that $\text{Risk}(\widehat{\boldsymbol{\theta}})$ is increasing in r_{ns} , while $\text{Risk}(\widehat{\boldsymbol{\theta}}_L)$ does not depend on r_{ns} . Therefore, the gain Δ is increasing in r_{ns} .

To understand the dependence of Δ on ρ , we write

$$\begin{aligned} \Delta - 1 &= \frac{\text{Risk}(\widehat{\boldsymbol{\theta}})}{\text{Risk}(\widehat{\boldsymbol{\theta}}_L)} - 1 \\ &= \frac{\frac{\psi_d}{\psi_d - 1} \left(\sigma^2 + \left(\frac{\mu^2}{k} + 1 \right) \left(1 + \frac{\mu^2/k + 1}{\psi_d - 1} \right)^{-2} \rho r_s^2 \right) + \left(1 - \frac{1}{\psi_d} \right) \left((1 - \rho)r_s^2 + r_{\text{ns}}^2 \right)}{\frac{\sigma^2 + r_s^2}{1 - \psi_d + \psi_p} - \rho r_s^2} - 1 \\ &= \frac{\frac{\psi_d}{\psi_d - 1} \left(\sigma^2 + \left(\frac{\mu^2}{k} + 1 \right) \left(1 + \frac{\mu^2/k + 1}{\psi_d - 1} \right)^{-2} \rho r_s^2 \right) + \left(1 - \frac{1}{\psi_d} \right) \left(r_s^2 + r_{\text{ns}}^2 \right) - \frac{\sigma^2 + r_s^2}{1 - \psi_d + \psi_p} + \frac{\rho r_s^2}{\psi_d}}{\frac{\sigma^2 + r_s^2}{1 - \psi_d + \psi_p} - \rho r_s^2} \end{aligned}$$

As we see the numerator is increasing in ρ and denominator is decreasing in ρ , which implies that the gain Δ is increasing in ρ .

We next show that $\Delta \geq 1$ if condition (5.1) holds. Since Δ is decreasing in μ^2/k and increasing in ρ , it suffices to show the claim assuming $\mu^2/k \rightarrow \infty$ and $\rho = 0$. In this case we have $(\frac{\mu^2}{k} + 1)\tilde{\alpha}^2 \rightarrow 0$ and so

$$\begin{aligned} \Delta &\rightarrow \frac{\frac{\sigma^2 \psi_d}{\psi_d - 1} + \left(1 - \frac{1}{\psi_d}\right)(r_s^2 + r_{\text{ns}}^2)}{\frac{\sigma^2 + r_s^2}{1 - \psi_d + \psi_p}} \\ &\geq \frac{\frac{\sigma^2 \psi_d}{\psi_d - 1} + \left(1 - \frac{1}{\psi_d}\right)r_s^2}{\frac{\sigma^2 + r_s^2}{1 - \psi_d + \psi_p}} \\ &= \frac{\frac{\psi_d}{\psi_d - 1} + \left(1 - \frac{1}{\psi_d}\right)\text{SNR}^2}{\frac{1 + \text{SNR}^2}{1 - \psi_d + \psi_p}} \geq 1, \end{aligned}$$

where the last step follows from condition (5.1).