# High Dimensional Data Enrichment: Interpretable, Fast, and Data-Efficient[*]

Amir Asiaee[†], Samet Oymak[‡], Kevin R. Coombes[§], and Arindam Banerjee[¶]

**Abstract.** Given samples from a set of groups, a data-enriched model describes observations by a common and per-group individual parameters. In high-dimensional regime, each parameter has its own structure such as sparsity or group sparsity. In this paper, we consider the general form of data enrichment where data comes in a fixed but arbitrary number of groups $G$ and any convex function, e.g., norm, can characterize the structure of both common and individual parameters. We propose an estimator for the high-dimensional data enriched model and investigate its statistical properties. We delineate sample complexity of our estimator and provide high probability non-asymptotic bound for estimation error of all parameters under a condition weaker than the state-of-the-art. We propose an iterative estimation algorithm with a geometric convergence rate and supplement our theoretical analysis with synthetic and real experimental results. In particular, we show the predictive power of data-enriched model along with its interpretable results in anticancer drug sensitivity analysis. Overall, we present a first through statistical and computational analysis of inference in the data enriched model.

**Key words.** example, LaTeX

**AMS subject classifications.** 68Q25, 68R10, 68U05

## 1. Introduction.
Over the past two decades, major advances have been made in estimating structured parameters, e.g., sparse, low-rank, etc., in high-dimensional small sample problems [11, 18, 19]. Such estimators consider a suitable (semi) parametric model of the response: $y = \phi(\mathbf{x}, \boldsymbol{\beta}^*) + w$ based on $n$ samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and the parameter of interest, $\boldsymbol{\beta}^* \in \mathbb{R}^p$. The unique aspect of such high-dimensional regime of $n \ll p$ is that the structure of $\boldsymbol{\beta}^*$ makes the estimation possible for large enough samples $n = m$ known as the sample complexity [9, 10, 36]. While the earlier developments in such high-dimensional estimation problems had focused on parametric linear models, the results have been widely extended to non-linear models, e.g., generalized linear models [1, 28], broad families of semi-parametric and single-index models [7, 33], non-convex models [5, 23], etc.

In several real world problems, the assumption that one global model parameter $\boldsymbol{\beta}_0^*$ is suitable for the entire population is unrealistic. We consider the more general setting where the population consists of sub-populations (groups) which are similar is many aspects but have unique differences. For example, in the context of anti-cancer drug sensitivity prediction where the goal is to predict responses of different tumor cells to a drug, using a same prediction model across cancer types (groups) ignores the unique properties of each cancer and leads to an uninterpretable global model. Alternatively, in such a setting, one can assume a separate model for each group $g$ as $y = \phi(\mathbf{x}, \boldsymbol{\beta}_g^*) + w$ based on a group specific parameter $\boldsymbol{\beta}_g^*$. Such a modeling choice fails to leverage the similarities across the sub-populations, and can only be estimated when sufficient number of samples are available for each

[†]Mathematical Biosciences Institute, The Ohio State University, Columbus, OH (asiaeetaheri.1@osu.edu).
[‡]Department of Electrical and Computer Engineering, UC Riverside, Riverside, CA (oymak@ece.ucr.edu).
[§]Departments of Biomedical Informatics, The Ohio State University, Columbus, OH (coombes.3@osu.edu).
[¶]Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN (banerjee@cs.umn.edu).
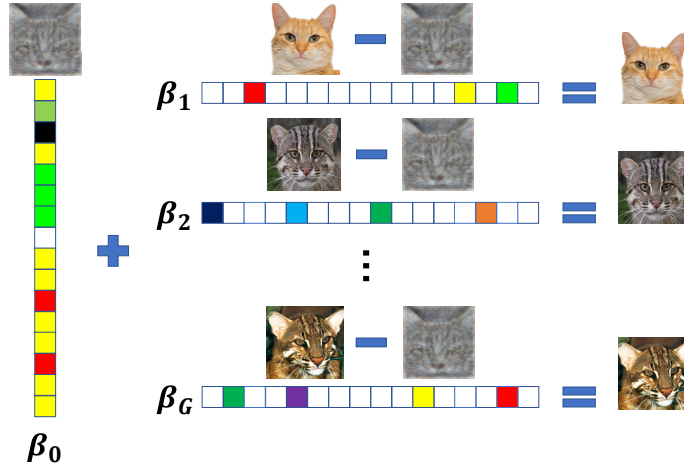
Figure 1: A conceptual illustration of data enrichment model for learning representation of cat species. The common parameter $\beta_0$ captures a *generic cat* which consists of shared features among all cats.

35 group which is not the case in several problems, e.g., anti-cancer drug sensitivity prediction [3, 22].
36      The middle ground model for such a scenario is the *superposition* of common and individual
37 parameters $\beta_0^* + \beta_g^*$ which has been of recent interest [?, 21]. Such a collection of *coupled* superposition
38 models is known by multiple names in the statistical machine learning community. It is a form of
39 multi-task learning [24, 42] when we consider regression in each group as a task. It is also called data
40 sharing [20] since information contained in different groups is shared through the common parameter
41 $\beta_0^*$. And finally, it has been called data enrichment [14] because we enrich our data set with pooling
42 multiple samples from different but related sources.
43      Following the successful application of such a modeling scheme in recent years [17, 20, 30, 31],
44 we consider the below *data enrichment* (DE) model:

45  (1.1) $$y_{gi} = \phi(\mathbf{x}_{gi}, (\beta_0^* + \beta_g^*)) + w_{gi}, \quad g \in \{1, \ldots, G\},$$

46 where $g$ and $i$ index the group and samples respectively. DE model (1.1) assumes that there is a *common*
47 parameter $\beta_0^*$ shared between all groups which models similarities between all samples. And there are
48 *individual* per-group parameters $\beta_g^*$s each characterize the deviation of group $g$, Figure 1.
49      *The setting.* Our goal is to design an estimation procedure which consistently recovers all
50 parameters of DE (1.1) fast and with small number of samples. We specifically focus on the high-
51 dimensional small sample regime where the number of samples $n_g$ for each group is much smaller
52 than the ambient dimensionality, i.e., $\forall g : n_g \ll p$. Similar to all other high-dimensional models, we
53 assume that the parameters $\beta_g$ are structured, i.e., for suitable convex functions $f_g$'s, $f_g(\beta_g)$ is small.
54 For example, when the structure is sparsity, $f_g$s are $L_1$-norms. Further, for the technical analysis and
55 proofs, we focus on the case of linear models, i.e., $\phi(\mathbf{x}, \beta) = \mathbf{x}^T\beta$. The results seamlessly extend to
56 more general non-linear models, e.g., generalized linear models, broad families of semi-parametric and
57 single-index models, non-convex models, etc., using existing results, i.e., how models like LASSO have
58 been extended to these settings [27].

**1.1. Related Work.** In the context of *Multi-Task Learning* (MTL), similar models have been proposed which have the general form of $y_{gi} = \mathbf{x}_{gi}^T(\boldsymbol{\beta}_{1g}^* + \boldsymbol{\beta}_{2g}^*) + w_{gi}$ where $\mathbf{B}_1 = [\boldsymbol{\beta}_{11}, \ldots, \boldsymbol{\beta}_{1G}]$ and $\mathbf{B}_2 = [\boldsymbol{\beta}_{21}, \ldots, \boldsymbol{\beta}_{2G}]$ are two parameter matrices [42]. To capture the relation of tasks, different types of constraints are assumed for parameter matrices. For example, [16] assumes $\mathbf{B}_1$ and $\mathbf{B}_2$ are sparse and low rank respectively. In this parameter matrix decomposition framework for MLT, the most related work to ours is the Dirty Model (DM) proposed in [24] where authors regularize the regression with $\|\mathbf{B}_1\|_{1,\infty}$ and $\|\mathbf{B}_2\|_{1,1}$ where norms are $p, q$-norms on *rows* of matrices, i.e., $\|\cdot\|_{p,q} = \|(\|\cdot\|_q, \ldots, \|\cdot\|_q)\|_p$.

If in our DE model we pick all structure inducing functions $f_g$ to be $l_1$-norm, the resulting model is very similar to the DM where $\|\mathbf{B}_1\|_{1,\infty}$ induces similarity between tasks and $\|\mathbf{B}_2\|_{1,1}$ models their discrepancies. On the other hand, the degree of freedom of DM model is higher than DE because $\|\mathbf{B}_1\|_{1,\infty}$ regularizer enforces shared support of $\boldsymbol{\beta}_{1g}^*$s, i.e., $\text{supp}(\boldsymbol{\beta}_{1i}^*) = \text{supp}(\boldsymbol{\beta}_{1j}^*)$ but allows $\boldsymbol{\beta}_{1i}^* \neq \boldsymbol{\beta}_{1j}^*$ while in DE we have a single common parameter $\boldsymbol{\beta}_0^*$. So one would expect that DE estimators should have smaller sample complexity compared to their DM counterparts and our analysis confirm that our estimator is more data efficient than DM estimator of [24]. Mainly, they require every task $i$ to have large enough samples to learn its own common parameters $\boldsymbol{\beta}_i$ but since DE shares the common parameter it only requires the *total dataset over all tasks* to be sufficiently large.

The linear DE model where $\boldsymbol{\beta}_g$'s are sparse has recently gained attention because of its application in wide range of domains such as personalized medicine [17], sentiment analysis, banking strategy [20], single cell data analysis [31], road safety [30], and disease subtype analysis [17]. More generally, in any high-dimensional problem where the population consists of groups, data enrichment framework has the potential to boost the prediction accuracy and results in a more interpretable set of parameters.

*Motivation.* In spite of the recent surge in applying data enrichment framework to different domains, limited advances have been made in understanding the statistical and computational properties of suitable estimators for the DE model (1.1). In fact, non-asymptotic statistical properties, including sample complexity and statistical rates of convergence, of regularized estimators for the data enriched model is still an open question [20, 30]. To the best of our knowledge, the only theoretical guarantee for data enrichment is provided in [31] where authors prove sparsistency of their proposed method under the stringent irrepresentability condition of the design matrix for recovering *supports* of common and individual parameters. Existing support recovery guarantees [31], sample complexity and $l_2$ consistency results [24] of related MTL models are restricted to sparsity and $l_1$-norm, while our estimator and *norm consistency* analysis work for *any* structure induced by arbitrary convex functions $f_g$. Moreover, no computational results, such as rates of convergence of the estimation procedures exist in the literature.

**1.2. Notation and Preliminaries.** We denote sets by curly $\mathcal{V}$, matrices by bold capital $\mathbf{V}$, random variables by capital $V$, and vectors by small bold $\mathbf{v}$ letters. We take $[G] = \{1, \ldots, G\}$ and $[G_+] = [G] \cup \{0\}$. Throughout the manuscript $c_i$ and $C_i$ denote positive absolute constants.

*Sub-Gaussian random variable and vector.* A random variable $V$ is sub-Gaussian if its moments satisfies $\forall p \geq 1 : (\mathbb{E}|V|^p)^{1/p} \leq K_2\sqrt{p}$. The minimum value of $K_2$ is called the sub-Gaussian norm of $V$, denoted by $\||V\||_{\psi_2}$ [39]. A random vector $\mathbf{v} \in \mathbb{R}^p$ is sub-Gaussian if the one-dimensional marginals $\langle \mathbf{v}, \mathbf{u} \rangle$ are sub-Gaussian random variables for all $\mathbf{u} \in \mathbb{R}^p$. The sub-Gaussian norm of $\mathbf{v}$ is defined [39] as $\||\mathbf{v}\||_{\psi_2} = \sup_{\mathbf{u} \in \mathbb{S}^{p-1}} \||\langle \mathbf{v}, \mathbf{u} \rangle\||_{\psi_2}$. For any set $\mathcal{V} \in \mathbb{R}^p$ the Gaussian width of the set $\mathcal{V}$ is defined as $\omega(\mathcal{V}) = \mathbb{E}_{\mathbf{g}} \left[\sup_{\mathbf{u} \in \mathcal{V}} \langle \mathbf{g}, \mathbf{u} \rangle\right]$ [40], where the expectation is over $\mathbf{g} \sim N(\mathbf{0}, \mathbf{I}_{p \times p})$, a vector of independent zero-mean unit-variance Gaussian.

Given $G$ groups and $n_g$ samples in each as $\{\{\mathbf{x}_{gi}, y_{gi}\}_{i=1}^{n_g}\}_{g=1}^{G}$, we can form the per group design matrix $\mathbf{X}_g \in \mathbb{R}^{n_g \times p}$ and output vector $\mathbf{y}_g \in \mathbb{R}^{n_g}$. The total number of samples is $n = \sum_{g=1}^{G} n_g$ and the data enriched model takes the following vector form:

(1.2) $$\mathbf{y}_g = \mathbf{X}_g(\boldsymbol{\beta}_0^* + \boldsymbol{\beta}_g^*) + \mathbf{w}_g, \quad \forall g \in [G]$$

where each row of $\mathbf{X}_g$ is $\mathbf{x}_{gi}^T$ and $\mathbf{w}_g^T = (w_{g1}, \ldots, w_{gn_g})$ is the noise vector.

**1.3. Our Contributions.** We propose the following Data Enrichment (DE) estimator $\hat{\boldsymbol{\beta}}$ for recovering the structured parameters where the structure is induced by *convex* functions $f_g(\cdot)$:

(1.3) $$\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}}_0^T, \ldots, \hat{\boldsymbol{\beta}}_G^T) \in \operatorname*{argmin}_{\boldsymbol{\beta}_0, \ldots, \boldsymbol{\beta}_G} \frac{1}{n} \sum_{g=1}^{G} \|\mathbf{y}_g - \mathbf{X}_g(\boldsymbol{\beta}_0 + \boldsymbol{\beta}_g)\|_2^2,$$

$$\text{s.t.} \quad \forall g \in [G] \cup \{0\} : f_g(\boldsymbol{\beta}_g) \leq f_g(\boldsymbol{\beta}_g^*).$$

We present several statistical and computational results for the DE estimator (1.3):

- The DE estimator (1.3) succeeds if a geometric condition that we call *Data EnRichment Incoherence Condition* (DERIC) is satisfied, Figure 2b. Compared to other known geometric conditions in the literature such as structural coherence [21] and stable recovery conditions [25], DERIC is a considerably weaker condition, Figure 2a.
- Assuming DERIC holds, we establish a high probability non-asymptotic bound on the weighted sum of parameter-wise estimation error, $\boldsymbol{\delta}_g = \hat{\boldsymbol{\beta}}_g - \boldsymbol{\beta}_g^*$ as:

(1.4) $$\sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 \leq \gamma O\left(\frac{\max_{g \in [G]} \omega(\mathcal{C}_g \cap \mathbb{S}^{p-1})}{\sqrt{n}}\right),$$

where $n_0 \triangleq n$ is the total number of samples, $\gamma \triangleq \max_{g \in [G]} \frac{n}{n_g}$ is the *sample condition number*, and $\mathcal{C}_g$ is the error cone corresponding to $\boldsymbol{\beta}_g^*$ exactly defined in Section **??**. To the best of our knowledge, this is the first statistical estimation guarantee for the data enrichment.
- We also establish the sample complexity of the DE estimator for all parameters as $\forall g \in [G] \cup \{0\}$ : $n_g = O(\omega(\mathcal{C}_g \cap \mathbb{S}^{p-1}))^2$. We emphasize that our result proofs that the recovery of the common parameter $\boldsymbol{\beta}_0$ by DE estimator (1.3) benefits from *all* of the $n$ pooled samples.
- We present an efficient projected block gradient descent algorithm DICER, to solve DE's objective (1.3) which converges geometrically to the statistical error bound of (1.4). To the best of our knowledge, this is the first rigorous computational result for the high-dimensional data-enriched regression.
- We illustrate promising empirical performance of the model on synthetic data as well as on the problem of finding bio-markers associated with drug sensitivity of cell lines from different cancer types, where the support of estimated individual parameters $\text{supp}(\hat{\boldsymbol{\beta}}_g)$ for each cancer type $g$ represents a different set of bio-markers per cancer type.

The rest of this paper is organized as follows: First, we characterize the error set of our estimator and provide a deterministic error bound in Section 2. Then in Section 3, we discuss the restricted eigenvalue condition and calculate the sample complexity required for the recovery of the true parameters by our estimator under DERIC condition. We close the statistical analysis in Section 4 by providing
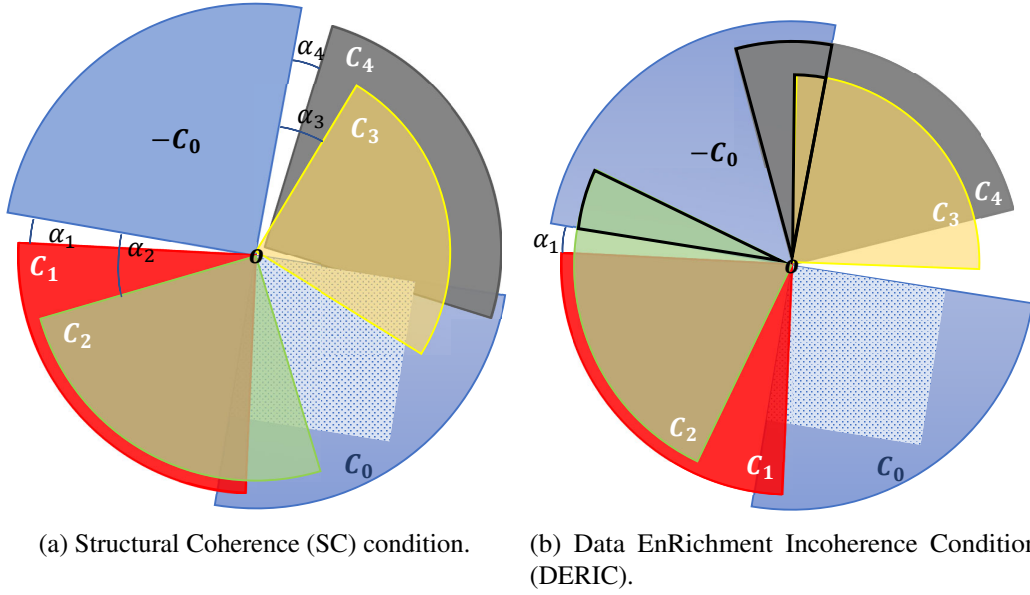
(a) Structural Coherence (SC) condition.

(b) Data EnRichment Incoherence Condition (DERIC).

Figure 2: a) State-of-the-art condition for recovering common and individual parameters in superposition models where $\mathcal{C}_g = \mathrm{Cone}(\mathcal{E}_g)$ are error cones and $\mathcal{E}_g = \left\{ \boldsymbol{\delta}_g | f_g(\boldsymbol{\beta}_g^* + \boldsymbol{\delta}_g) \leq f_g(\boldsymbol{\beta}_g^*) \right\}$ are the error sets for each parameter $\boldsymbol{\beta}_g^* \in [G]$ [21]. b) Our more relaxed recovery condition which allows *arbitrary non-zero fraction* of the error cones of individual parameters intersect with $-\mathcal{C}_0$.

non-asymptotic high probability error bound for parameter recovery. We delineate our geometrically convergent algorithm, DICER in Section 5 and finally supplement our work with synthetic and real experiments in Sections 6 and 7.

**2. The Data Enrichment Estimator.** A compact form of our proposed DE estimator (1.3) is:

(2.1) $$\hat{\boldsymbol{\beta}} \in \operatorname*{argmin}_{\boldsymbol{\beta}} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2, \quad \text{s.t. } \forall g \in [G] \cup \{0\} : f_g(\boldsymbol{\beta}_g) \leq f_g(\boldsymbol{\beta}_g^*),$$

where $\mathbf{y} = (\mathbf{y}_1^T, \dots \mathbf{y}_G^T)^T \in \mathbb{R}^n$, $\boldsymbol{\beta} = (\boldsymbol{\beta}_0^T, \dots, \boldsymbol{\beta}_G^T)^T \in \mathbb{R}^{(G+1)p}$ and

(2.2) $$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 & \mathbf{X}_1 & 0 & \cdots & 0 \\ \mathbf{X}_2 & 0 & \mathbf{X}_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \cdots & \vdots \\ \mathbf{X}_G & 0 & \cdots & \cdots & \mathbf{X}_G \end{pmatrix} \in \mathbb{R}^{n \times (G+1)p} .$$

*Example* 2.1. ($L_1$-**norm**) When all parameters $\boldsymbol{\beta}_g$s are $s_g$-sparse, i.e.,$|\mathrm{supp}(\boldsymbol{\beta}_g^*)| = s_g$ by using $l_1$-norm as the sparsity inducing function, our DE estimator of (2.1) becomes:

(2.3) $$\hat{\boldsymbol{\beta}} \in \operatorname*{argmin}_{\boldsymbol{\beta}} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2, \quad \text{s.t. } \forall g \in [G] \cup \{0\} : \|\boldsymbol{\beta}_g\|_1 \leq \|\boldsymbol{\beta}_g^*\|_1.$$

We call (2.3) *sparse DE* estimator and use it as the running example throughout the paper to illustrate outcomes of our analysis.

Consider the group-wise estimation error $\boldsymbol{\delta}_g = \hat{\boldsymbol{\beta}}_g - \boldsymbol{\beta}_g^*$. Since $\hat{\boldsymbol{\beta}}_g = \boldsymbol{\beta}_g^* + \boldsymbol{\delta}_g$ is a feasible point of (2.1), the error vector $\boldsymbol{\delta}_g$ will belong to the following restricted error set:

$$(2.4) \qquad \mathcal{E}_g = \left\{ \boldsymbol{\delta}_g | f_g(\boldsymbol{\beta}_g^* + \boldsymbol{\delta}_g) \leq f_g(\boldsymbol{\beta}_g^*) \right\}, \quad g \in [G] \cup \{0\}.$$

We denote the cone of the error set as $\mathcal{C}_g \triangleq \text{Cone}(\mathcal{E}_g)$ and the spherical cap corresponding to it as $\mathcal{A}_g \triangleq \mathcal{C}_g \cap \mathbb{S}^{p-1}$. Consider the set $\mathcal{C} = \{\boldsymbol{\delta} = (\boldsymbol{\delta}_0^T, \ldots, \boldsymbol{\delta}_G^T)^T \big| \boldsymbol{\delta}_g \in \mathcal{C}_g\}$, following two subsets of $\mathcal{C}$ play key roles in our analysis:

$$(2.5) \qquad \mathcal{H} \triangleq \left\{ \boldsymbol{\delta} \in \mathcal{C} \big| \sum_{g=0}^{G} \frac{n_g}{n} \|\boldsymbol{\delta}_g\|_2 = 1 \right\}, \quad \bar{\mathcal{H}} \triangleq \left\{ \boldsymbol{\delta} \in \mathcal{C} \big| \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 = 1 \right\}.$$

Starting from the optimality of $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}^* + \boldsymbol{\delta}$ as $\frac{1}{n}\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2 \leq \frac{1}{n}\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*\|_2^2$, we have: $\frac{1}{n}\|\mathbf{X}\boldsymbol{\delta}\|_2^2 \leq \frac{1}{n}2\mathbf{w}^T\mathbf{X}\boldsymbol{\delta}$ where $\mathbf{w} = [\mathbf{w}_1^T, \ldots, \mathbf{w}_G^T]^T \in \mathbb{R}^n$ is the vector of all noises. Using this basic inequality, we can establish the following deterministic error bound.

**Theorem 2.2.** *For the DE estimator* (2.1), *assume there exist* $0 < \kappa \leq \inf_{\mathbf{u}\in\mathcal{H}} \frac{1}{n}\|\mathbf{Xu}\|_2^2$. *Then, for the sample condition number* $\gamma = \max_{g\in[G]} \frac{n}{n_g}$, *the following deterministic upper bounds holds:*

$$\sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 \leq \frac{2\gamma \sup_{\mathbf{u}\in\bar{\mathcal{H}}} \boldsymbol{\omega}^T \mathbf{Xu}}{n\kappa}.$$

*Proof.* We lower bound the LHS and upper bound the RHS of the optimality inequality $\frac{1}{n}\|\mathbf{X}\boldsymbol{\delta}\|_2^2 \leq \frac{1}{n}2\mathbf{w}^T\mathbf{X}\boldsymbol{\delta}$ using the definition of the sets $\mathcal{H}$ and $\bar{\mathcal{H}}$ respectively. Starting with the lower bound using the definition of set $\mathcal{H}$ (2.5) we have:

$$(2.6) \qquad \frac{1}{n}\|\mathbf{X}\boldsymbol{\delta}\|_2^2 \geq \frac{1}{n} \inf_{\mathbf{u}\in\mathcal{H}} \|\mathbf{Xu}\|_2^2 \left( \sum_{g=0}^{G} \frac{n_g}{n} \|\boldsymbol{\delta}_g\|_2 \right)^2$$

$$\geq \kappa \left( \sum_{g=0}^{G} \frac{n_g}{n} \|\boldsymbol{\delta}_g\|_2 \right)^2$$

$$\geq \kappa \left( \min_{g\in[G]} \frac{n_g}{n} \right) \left( \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 \right)^2$$

where $0 < \kappa \leq \frac{1}{n} \inf_{\mathbf{u}\in\mathcal{H}} \|\mathbf{Xu}\|_2^2$ is known as Restricted Eigenvalue (RE) condition. The upper bound factorizes as:

$$(2.7) \qquad \frac{2}{n}\mathbf{w}^T\mathbf{X}\boldsymbol{\delta} \leq \frac{2}{n} \sup_{\mathbf{u}\in\bar{\mathcal{H}}} \mathbf{w}^T\mathbf{Xu} \left( \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 \right), \quad \mathbf{u} \in \mathcal{H}$$

Putting together inequalities (2.6) and (2.7) completes the proof. ∎

*Remark* 2.3. Consider the setting where $n_g = \Theta(\frac{n}{G})$ so that each group has approximately $\frac{1}{G}$ fraction of the samples. Then, $\gamma = \Theta(G)$ and hence

$$\frac{1}{G}\sum_{g=0}^{G}\|\delta_g\|_2 \leq O(G^{1/2})\frac{\sup_{\mathbf{u}\in\bar{\mathcal{H}}}\boldsymbol{\omega}^T\mathbf{X}\mathbf{u}}{n}.$$

**3. Restricted Eigenvalue Condition.** The main assumptions of Theorem 2.2 is known as Restricted Eigenvalue (RE) condition in the literature of high-dimensional statistics [2, 29, 34]: $\inf_{\mathbf{u}\in\mathcal{H}}\frac{1}{n}\|\mathbf{X}\mathbf{u}\|_2^2 \geq \kappa > 0$. The RE condition posits that the minimum eigenvalues of the matrix $\mathbf{X}^T\mathbf{X}$ in directions restricted to $\mathcal{H}$ is strictly positive. In this section, we show that for the design matrix $\mathbf{X}$ defined in (2.2), the RE condition holds with high probability under a suitable geometric condition we call *Data EnRichment Incoherence Condition* (DERIC) and for enough number of samples. We precisely characterize total and per-group sample complexities required for successful parameter recovery. For the analysis, similar to existing work [21, 26, 38], we assume the design matrix to be isotropic sub-Gaussian.[1]

*Definition* 3.1. *We assume* $\mathbf{x}_{gi}$ *are i.i.d. random vectors from a non-degenerate zero-mean, isotropic sub-Gaussian distribution. In other words,* $\mathbb{E}[\mathbf{x}] = 0$, $\mathbb{E}[\mathbf{x}^T\mathbf{x}] = \mathbf{I}_{p\times p}$, *and* $\|\|\mathbf{x}\|\|_{\psi_2} \leq k_x$. *As a consequence,* $\exists \alpha > 0$ *such that* $\forall \mathbf{u} \in \mathbb{S}^{p-1}$ *we have* $\mathbb{E}|\langle\mathbf{x},\mathbf{u}\rangle| \geq \alpha$. *Further, we assume noise* $\mathbf{w}_{gi}$ *are i.i.d. zero-mean, unit-variance sub-Gaussian with* $\|\|\mathbf{w}_{gi}\|\|_{\psi_2} \leq k_w$.

**3.1. Geometric Condition of Recovery.** Unlike standard high-dimensional statistical estimation, for RE condition to be true, parameters of superposition models need to satisfy geometric conditions which limits the interaction of parameters with each other to make sure that recovery is possible. In this section, we elaborate our sufficient geometric condition for recovery and compare it with state-of-the-art condition for recovery of superposition models.

To intuitively illustrate the necessity of such a geometric condition, consider the simplest superposition model i.e., $\boldsymbol{\beta}_0^* + \boldsymbol{\beta}_g^*$. Without any restriction on parameter interactions, any estimates such that $\hat{\boldsymbol{\beta}}_0 + \hat{\boldsymbol{\beta}}_g = \boldsymbol{\beta}_0^* + \boldsymbol{\beta}_g^*$ are valid ones. To avoid such trivial solutions two error cones need to satisfy $\boldsymbol{\delta}_g \neq -\boldsymbol{\delta}_0$. In general, the RE condition of individual superposition models can be established under the so-called Structural Coherence (SC) condition [21, 25] which is the generalization of this idea for superposition of multiple parameters as $\sum_{g=0}^{G}\boldsymbol{\beta}_g^*$.

*Definition* 3.2 (Structural Coherence (SC) [21, 25]). *Consider a superposition model of the form* $y = \mathbf{x}^T\sum_{g=0}^{G}\boldsymbol{\beta}_g^* + w$. *The SC condition requires that*

$$(3.1) \qquad \forall\boldsymbol{\delta}_g \in \mathcal{C}_g, \exists\lambda \quad s.t. \quad \|\sum_{g=0}^{G}\boldsymbol{\delta}_g\|_2 \geq \lambda\sum_{g=0}^{G}\|\boldsymbol{\delta}_g\|_2,$$

*and leads to the corresponding RE condition for the superposition model.*

*Remark* 3.3. Note that the SC conditions is satisfied if none of the individual error cones $\mathcal{C}_g$ intersect with the inverted error cone $-\mathcal{C}_0$, i.e., $\forall g, \alpha_g < 1$ where $\alpha_g = \sup\langle\boldsymbol{\delta}_0/\|\boldsymbol{\delta}_0\|_2, \boldsymbol{\delta}_g/\|\boldsymbol{\delta}_g\|_2\rangle$ [21, 38], Figure 2a.

---

[1]Extension to an-isotropic sub-Gaussian case is straightforward by techniques developed in [2, 35].

207     The SC condition on each individual problem fails to utilize the true coupling structure in the
208 data enriched model, where $\boldsymbol{\beta}_0^*$ is involved in all groups. In fact, below we show, using SC on each
209 individual model leads to radically pessimistic estimates of the sample complexity for $\boldsymbol{\beta}_0^*$ recovery.

210     **Proposition 3.4.** *Assume observations distributed as defined in Definition 3.1 and pair-wise SC*
211 *conditions are satisfied. Consider each superposition model* (1.2) *in isolation; to recover the common*
212 *parameter $\boldsymbol{\beta}_0^*$ requires at least one group $i$ to have $n_i = O(\omega^2(\mathcal{A}_0))$. To recover the rest of individual*
213 *parameters, we need $\forall g \neq i : n_g = O(\omega^2(\mathcal{A}_g))$ samples.*

214 In other words, by separate analysis of superposition estimators at least one problem needs to have
215 sufficient samples for recovering the common parameter $\boldsymbol{\beta}_0$ and therefore the common parameter
216 recovery does not benefit from the pooled $n$ samples. But given the nature of coupling in the data
217 enriched model, we hope to be able to get a better sample complexity specifically for the common
218 parameter $\boldsymbol{\beta}_0$.

219     Here, we introduce DERIC, a considerably weaker geometric condition compared to SC of [21, 25].

220     **Definition 3.5** (Data EnRichment Incoherence Condition (DERIC)). *There exists a non-empty*
221 *set $\mathcal{I} \subseteq [G]_{\backslash}$ of groups where for some scalars $0 < \bar{\rho} \leq 1$ and $\lambda_{\min} > 0$ the following holds:*
222     1. $\sum_{i \in \mathcal{I}} n_i \geq \lceil \bar{\rho} n \rceil$.
223     2. $\forall i \in \mathcal{I}, \forall \boldsymbol{\delta}_i \in \mathcal{C}_i$, and $\boldsymbol{\delta}_0 \in \mathcal{C}_0$: $\|\boldsymbol{\delta}_i + \boldsymbol{\delta}_0\|_2 \geq \lambda_{\min}(\|\boldsymbol{\delta}_0\|_2 + \|\boldsymbol{\delta}_i\|_2)$
224 *Observe that $0 < \lambda_{\min}, \bar{\rho} \leq 1$ by definition.*

225     *Remark 3.6.* Clearly DERIC and SC conditions are satisfied if the error cones $\mathcal{C}_g$ and $\mathcal{C}_0$ does
226 not have a ray in common, i.e., $\sup\langle \boldsymbol{\delta}_0/\|\boldsymbol{\delta}_0\|_2, \boldsymbol{\delta}_g/\|\boldsymbol{\delta}_g\|_2 \rangle < 1$ [38, 21], Figure 2a. In particular, SC
227 requires that none of the individual error cones $\mathcal{C}_g$ intersect with the inverted error cone $-\mathcal{C}_0$. Instead of
228 this stringent geometric condition, DERIC allows $-\mathcal{C}_0$ to intersect with an arbitrarily large fraction of
229 the $\mathcal{C}_g$ cones, Figure 2b. As the number of intersections increases, our bound becomes looser.

230     **3.2. Sample Complexity.** Using DERIC and the small ball method [26], a recent tool from
231 empirical process theory in the following theorem, we get a better sample complexity required for
232 satisfying the RE condition:

233     **Theorem 3.7.** *Let $\mathbf{x}_{gi}s$ be random vectors defined in Definition 3.1. Assume DERIC condition*
234 *of Definition 3.5 holds for error cones $\mathcal{C}_g$s and $\psi_{\mathcal{I}} = \lambda_{\min}\bar{\rho}/3$. Then, for all $\boldsymbol{\delta} \in \mathcal{H}$, when we have*
235 *enough number of samples as $\forall g \in [G]_{\backslash} : n_g \geq m_g = O(k^6 \alpha^{-6} \psi_{\mathcal{I}}^{-2} \omega(\mathcal{A}_g)^2)$, with probability at least*
236 $1 - e^{-n\kappa_{\min}/4}$ *we have:*

$$\inf_{\boldsymbol{\delta} \in \mathcal{H}} \frac{1}{\sqrt{n}} \|\mathbf{X}\boldsymbol{\delta}\|_2 \geq \frac{\kappa_{\min}}{2}$$

238 *where $\kappa_{\min} = \min_{g \in [G]_{\backslash}} C\psi_{\mathcal{I}} \frac{\alpha^3}{k^2} - \frac{2c_g k \omega(\mathcal{A}_g)}{\sqrt{n_g}}$ and $\kappa = \frac{\kappa_{\min}^2}{4}$ is the lower bound of the RE condition.*

239     *Example 3.8.* ($L_1$-**norm**) The Gaussian width of the spherical cap of a $p$-dimensional $s$-sparse
240 vector is $\omega(\mathcal{A}) = \Theta(\sqrt{s \log p})$ [2, 40]. Therefore, the number of samples per group and total required for
241 satisfaction of the RE condition in the sparse DE estimator (2.3) is $\forall g \in [G] : n_g \geq m_g = \Theta(s_g \log p)$.

242     **4. General Error Bound.** In this section, we provide a high probability upper bound for the
243 estimation error of the common and individual parameters. To avoid cluttering the notation, we rename
244 the vector of all noises as $\boldsymbol{\omega}_0 \triangleq \boldsymbol{\omega}$. First, we massage the deterministic upper bound of Theorem 2.2 as

follows:

$$\boldsymbol{\omega}^T \mathbf{X} \boldsymbol{\delta} = \sum_{g=0}^{G} \langle \mathbf{X}_g^T \boldsymbol{\omega}_g, \boldsymbol{\delta}_g \rangle$$

$$= \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 \langle \mathbf{X}_g^T \frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2}, \frac{\boldsymbol{\delta}_g}{\|\boldsymbol{\delta}_g\|_2} \rangle \sqrt{\frac{n}{n_g}} \|\boldsymbol{\omega}_g\|_2$$

Assume $b_g = \langle \mathbf{X}_g^T \frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2}, \frac{\boldsymbol{\delta}_g}{\|\boldsymbol{\delta}_g\|_2} \rangle \sqrt{\frac{n}{n_g}} \|\boldsymbol{\omega}_g\|_2$ and $a_g = \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2$. Then the above term is the inner product of two vectors $\mathbf{a} = (a_0, \dots, a_G)$ and $\mathbf{b} = (b_0, \dots, b_G)$ for which we have:

$$\sup_{\mathbf{a} \in \mathcal{H}} \mathbf{a}^T \mathbf{b} = \sup_{\|\mathbf{a}\|_1 = 1} \mathbf{a}^T \mathbf{b} \leq \|\mathbf{b}\|_\infty = \max_{g \in [G]} b_g,$$

where the inequality holds because of the definition of the dual norm. We can upper bounds $b_g$s with high probability and then by union bound below theorem establishes a high probability upper bound for the deterministic bound of Theorem 2.2, i.e., $\frac{2}{n} \boldsymbol{\omega}^T \mathbf{X} \mathbf{u}$.

**Theorem 4.1.** *Assume $\mathbf{x}_{gi}$ and $\omega_{gi}$ distributed according to Definition 3.1, then with probability at least $1 - \sigma \exp\left(-\min_{g \in [G]} \left[\nu_g n_g - \log(G+1), \frac{\tau^2}{\eta_g^2 k^2}\right]\right)$ we have:*

$$\frac{2}{n} \boldsymbol{\omega}^T \mathbf{X} \boldsymbol{\delta} \leq \sqrt{\frac{8K^2 + 4}{n}} \max_{g \in [G]} \left( \zeta_g k \omega(\mathcal{A}_g) + \epsilon_g \sqrt{\log(G+1)} + \tau \right)$$

*where $\sigma = \max_{g \in [G]} \sigma_g$ and $\tau > 0$.*

The following corollary characterizes the general error bound and results from the direct combination of Theorem 2.2, Theorem 3.7, and Theorem 4.1.

**Corollary 4.2.** *For $\mathbf{x}_{gi}$ and $\omega_{gi}$ described in Definition 3.1 when we have enough number of samples $\forall g \in [G] : n_g > m_g$ which lead to $\kappa > 0$, the following general error bound holds with high probability for estimator (2.1):*

$$(4.1) \qquad \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 \leq C\gamma \frac{k\zeta \max_{g \in [G]} \omega(\mathcal{A}_g) + \epsilon \sqrt{\log(G+1)} + \tau}{\kappa_{\min}^2 \sqrt{n}}$$

*where $C = 8\sqrt{2K^2 + 1}$, $\zeta = \max_{g \in [G]} \zeta_g$, $\epsilon = \max_{g \in [G]} \epsilon_g$, $\gamma = \max_{g \in [G]\setminus} n/n_g$ and $\tau > 0$.*

*Example* 4.3. **($L_1$-norm)** For sparse DE estimator of (2.3), results of Theorem 3.7 and 4.1 translates to the following: For enough number of samples as $\forall g \in [G] : n_g \geq m_g = O(s_g \log p)$, the error bound of (4.1) simplifies to:

$$(4.2) \qquad \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 = O\left( \sqrt{\frac{(\max_{g \in [G]} s_g) \log p}{n}} \right)$$

Therefore, individual errors are bounded as $\|\boldsymbol{\delta}_g\|_2 = O(\sqrt{(\max_{g \in [G]} s_g) \log p / n_g})$ which is slightly worse than $O(\sqrt{s_g \log p / n_g})$, the well-known error bound for recovering an $s_g$-sparse vector from $n_g$ observations using LASSO or similar estimators [2, 12, 8, 13, 4]. Note that $\max_{g \in [G]} s_g$ (instead of $s_g$) is the price we pay to recover the common parameter $\boldsymbol{\beta}_0$.

---

**Algorithm 5.1** DICER

---

1: **input:** $\mathbf{X}, \mathbf{y}$, learning rates $(\mu_0, \dots, \mu_G)$, initialization $\boldsymbol{\beta}^{(1)} = \mathbf{0}$

2: **output:** $\hat{\boldsymbol{\beta}}$

3: **for** t = 1 **to** T **do**

4:     **for** g=1 **to** G **do**

5:       $\boldsymbol{\beta}_g^{(t+1)} = \Pi_{\Omega_{f_g}} \left( \boldsymbol{\beta}_g^{(t)} + \mu_g \mathbf{X}_g^T \left( \mathbf{y}_g - \mathbf{X}_g \left( \boldsymbol{\beta}_0^{(t)} + \boldsymbol{\beta}_g^{(t)} \right) \right) \right)$

6:     **end for**

7:   $\boldsymbol{\beta}_0^{(t+1)} = \Pi_{\Omega_{f_0}} \left( \boldsymbol{\beta}_0^{(t)} + \mu_0 \mathbf{X}_0^T \left( \mathbf{y} - \mathbf{X}_0 \boldsymbol{\beta}_0^{(t)} - \begin{pmatrix} \mathbf{X}_1 \boldsymbol{\beta}_1^{(t)} \\ \vdots \\ \mathbf{X}_G \boldsymbol{\beta}_G^{(t)} \end{pmatrix} \right) \right)$

8: **end for**

---

**5. Estimation Algorithm.** We propose *Data enrIChER* (DICER) a projected block gradient descent algorithm, Algorithm 5.1, where $\Pi_{\Omega_{f_g}}$ is the Euclidean projection onto the set $\Omega_{f_g}(d_g) = \{f_g(\boldsymbol{\beta}) \leq d_g\}$ where $d_g = f_g(\boldsymbol{\beta}_g^*)$ and is dropped to avoid cluttering. In practice, $d_g$ can be determined by cross-validation.

To analysis convergence properties of DICER, we should upper bound the error of each iteration. Let's $\boldsymbol{\delta}^{(t)} = \boldsymbol{\beta}^{(t)} - \boldsymbol{\beta}^*$ be the error of iteration $t$ of DICER, i.e., the distance from the true parameter (not the optimization minimum, $\hat{\boldsymbol{\beta}}$). We show that $\|\boldsymbol{\delta}^{(t)}\|_2$ decreases exponentially fast in $t$ to the statistical error $\|\boldsymbol{\delta}\|_2 = \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2$. We first start with the required definitions for our analysis.

**Definition 5.1.** *We define the following positive constants as functions of step sizes $\mu_g > 0$:*

$$\forall g \in [G] : \rho_g(\mu_g) = \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{B}_g} \mathbf{v}^T \left( \mathbf{I}_g - \mu_g \mathbf{X}_g^T \mathbf{X}_g \right) \mathbf{u},$$

$$\eta_g(\mu_g) = \mu_g \sup_{\mathbf{v} \in \mathcal{B}_g} \mathbf{v}^T \mathbf{X}_g^T \frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2},$$

$$\forall g \in [G]_{\backslash} : \phi_g(\mu_g) = \mu_g \sup_{\mathbf{v} \in \mathcal{B}_g, \mathbf{u} \in \mathcal{B}_0} -\mathbf{v}^T \mathbf{X}_g^T \mathbf{X}_g \mathbf{u},$$

*where $\mathcal{B}_g = \mathcal{C}_g \cap \mathbb{B}^p$ is the intersection of the error cone and the unit ball.*

In the following theorem, we establish a deterministic bound on iteration errors $\|\boldsymbol{\delta}_g^{(t)}\|_2$ which depends on constants defined in Definition 5.1.

**Theorem 5.2.** *For Algorithm 5.1 initialized by $\boldsymbol{\beta}^{(1)} = \mathbf{0}$, we have the following deterministic bound for the error at iteration $t + 1$:*

(5.1)
$$\sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g^{(t+1)}\|_2$$

$$\leq \rho^t \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\beta}_g^*\|_2 + \frac{1 - \rho^t}{1 - \rho} \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \eta_g \|\boldsymbol{\omega}_g\|_2,$$

*where $\rho \triangleq \max \left( \rho_0 + \sum_{g=1}^{G} \sqrt{\frac{n_g}{n}} \phi_g, \max_{g \in [G]} \left[ \rho_g + \sqrt{\frac{n}{n_g}} \frac{\mu_0}{\mu_g} \phi_g \right] \right)$.*

The RHS of (5.1) consists of two terms. If we keep $\rho < 1$, the first term approaches zero fast, and the second term determines the bound. In the following, we show that for specific choices of step sizes $\mu_g$s, the second term can be upper bounded using the analysis of Section 4. More specifically, the first term corresponds to the optimization error which shrinks in every iteration while the second term is constant times the upper bound of the statistical error characterized in Corollary 4.2. Therefore, if we keep $\rho$ below one, the estimation error of DE algorithm geometrically converges to the approximate statistical error bound.

One way for having $\rho < 1$ is to keep all arguments of $\max(\cdots)$ defining $\rho$ strictly below 1. To this end, we first establish high probability upper bound for $\rho_g$, $\eta_g$, and $\phi_g$ (in the Appendix A.6) and then show that with enough number of samples and proper step sizes $\mu_g$, $\rho$ can be kept strictly below one with high probability. The high probability bounds for constants in Definition 5.1 and the deterministic bound of Theorem 5.2 leads to the following theorem which shows that for enough number of samples, of the same order as the statistical sample complexity of Theorem 3.7, we can keep $\rho$ below one and have geometric convergence.

**Theorem 5.3.** *Let $\tau = C\sqrt{\log(G+1)} + b$ for $b > 0$ and $\omega_{0g} = \omega(\mathcal{A}_0) + \omega(\mathcal{A}_g)$. For the step sizes of:*

$$\mu_0 = \frac{\min_{g \in [G]_\setminus} h_g(\tau)^{-2}}{4n}, \forall \in [G]_\setminus : \mu_g = \frac{h_g(\tau)^{-1}}{2\sqrt{nn_g}}$$

*where $h_g(\tau) = \left(1 + c_{0g}\frac{\omega_{0g}+\tau}{\sqrt{n_g}}\right)$ and sample complexities of $\forall g \in [G] : n_g \geq 2c_g^2(2\omega(\mathcal{A}_g) + \tau)^2$, updates of the Algorithm 5.1 obey the following with high probability:*

$$\sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g^{(t+1)}\|_2 \leq r(\tau)^t \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\beta}_g^*\|_2$$

$$+ \frac{(G+1)\sqrt{(2K^2+1)}}{\sqrt{n}(1-r(\tau))} \left(\zeta k \max_{g \in [G]} \omega(\mathcal{A}_g) + \tau\right),$$

*where $r(\tau) < 1$.*

**Corollary 5.4.** *For enough number of samples, iterations of DE algorithm with step sizes $\mu_0 = \Theta(\frac{1}{n})$ and $\mu_g = \Theta(\frac{1}{\sqrt{nn_g}})$ geometrically converges to the following with high probability:*

(5.2)
$$\sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g^\infty\|_2 \leq c\frac{\zeta k \max_{g \in [G]} \omega(\mathcal{A}_g) + C\sqrt{\log(G+1)} + b}{\sqrt{n}(1-r(\tau))}$$

*where $c = (G+1)\sqrt{(2K^2+1)}$.*

It is instructive to compare RHS of (5.2) with that of (4.1): $\kappa_{\min}$ defined in Theorem 3.7 corresponds to $(1 - r(\tau))$ and the extra $G + 1$ factor corresponds to the sample condition number $\gamma = \max_{g \in [G]} \frac{n}{n_g}$. Therefore, Corollary 5.4 shows that DICER converges to a scaled variant of statical error bound determined in Corollary 4.2.

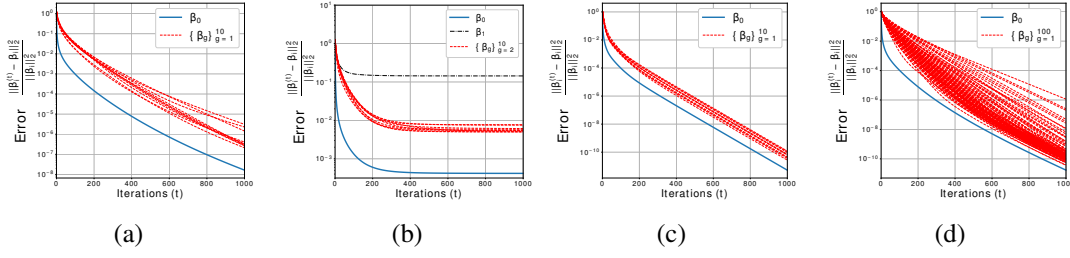(a)                    (b)                    (c)                    (d)

Figure 3: a) Noiseless fast convergence. b) Noise on the first group does not impact other groups as much. c) Increasing sample size improves rate of convergence. d) Our algorithm convergences fast even with a large number of groups $G = 100$.

**6. Synthetic Experiments.** We considered sparsity based simulations with varying $G$ and sparsity levels. In our first set of simulations, we set $p = 100$, $G = 10$ and sparsity of the individual parameters to be $s = 10$. We generated a dense $\boldsymbol{\beta}_0$ with $\|\boldsymbol{\beta}_0\| = p$ and did not impose any constraint. Iterates $\{\boldsymbol{\beta}_g^{(t)}\}_{g=1}^G$ are obtained by projection onto the $\ell_1$ ball $\|\boldsymbol{\beta}_g\|_1$. Nonzero entries of $\boldsymbol{\beta}_g$ are generated with $\mathcal{N}(0,1)$ and nonzero supports are picked uniformly at random. Inspired from our theoretical step size choices, in all experiments, we used simplified learning rates of $\frac{1}{n}$ for $\boldsymbol{\beta}_0$ and $\frac{1}{\sqrt{nn_g}}$ for $\boldsymbol{\beta}_g$, $g \in [G]_{\backslash}$. Observe that, cones of the individual parameters intersect with that of $\boldsymbol{\beta}_0$ hence this setup actually violates DERIC (which requires an arbitrarily small constant fraction of groups to be non-intersecting). Our intuition is that the individual parameters are mostly incoherent with each other and the existence of a nonzero perturbation over $\boldsymbol{\beta}_g$'s that keeps all measurements intact is unlikely. Remarkably, experimental results still show successful learning of all parameters from small amount of samples. We picked $n_g = 60$ for each group. Hence, in total, we have $11p = 1100$ unknowns, $200 = G \times 10 + 100$ degrees of freedom and $G \times 60 = 600$ samples. In all figures, we study the normalized squared error $\frac{\|\boldsymbol{\beta}_g^{(t)} - \boldsymbol{\beta}_g\|_2^2}{\|\boldsymbol{\beta}_g\|_2^2}$ and average 10 independent realization for each curve. Figure 3a shows the estimation performance as a function of iteration number $t$. While each group might behave slightly different, we do observe that all parameters are linear converging to ground truth.

In Figure 3b, we test the noise robustness of our algorithm. We add a $\mathcal{N}(0,1)$ noise to the $n_1 = 60$ measurements of the first group *only*. The other groups are left untouched. While all parameters suffer nonzero estimation error, we observe that, the global parameter $\boldsymbol{\beta}_0$ and noise-free groups $\{\boldsymbol{\beta}_g\}_{g=2}^G$ have substantially less estimation error. This implies that noise in one group mostly affects itself rather than the global estimation. In Figure 3c, we increased the sample size to $n_g = 150$ per group. We observe that, in comparison to Figure 3a, rate of convergence receives a boost from the additional samples as predicted by our theory.

Finally, Figure 3d considers a very high-dimensional problem where $p = 1000$, $G = 100$, individual parameters are 10 sparse, $\boldsymbol{\beta}_0$ is 100 sparse and $n_g = 150$. The total degrees of freedom is 1100, number of unknowns are 101000 and total number of datapoints are $150 \times 100 = 15000$. While individual parameters have substantial variation in terms of convergence rate, at the end of 1000 iteration, all parameters have relative reconstruction error below $10^{-6}$.

**7. Anti-Cancer Drug Sensitivity Prediction.** In this section, we investigate the application of DICER in analyzing the response of cancer tumor cell lines to different doses of various drugs. Each
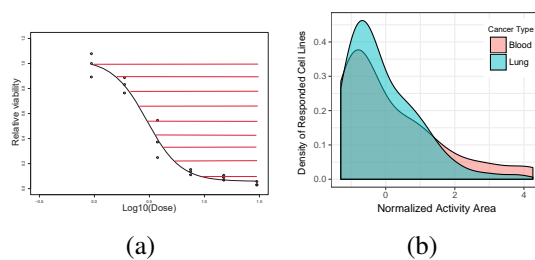
(a)                                          (b)

Figure 4: a) A sample fitted dose-response curve where Activity Area $y_{gi}$ is shaded. b) Distribution of responses to Saracatinib for which some lung and blood cancer cell lines have responded.



(a) Lung and Blood
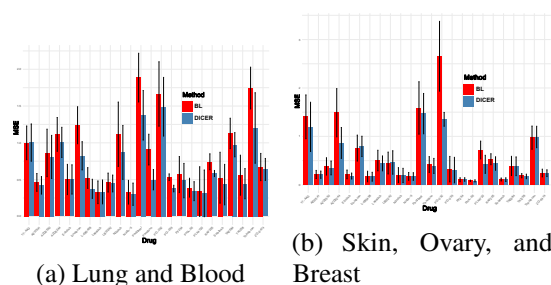
(b) Skin, Ovary, and Breast

Figure 5: Distribution of MSE Comparison of Mean Square Error of DICER and BL in predicting the response to 24 drugs in TWO and THREE experiments. Each bar is the mean of MSE for 5-fold cross-validation.

cancer type (lung, blood, etc.) is a group $g$ in our DE model and the respond of patient $i$ with cancer $g$ to the drug is our output $y_{gi}$. The set of features for each patient $\mathbf{x}_{gi}$ consists of gene expressions, copy number variation, and mutations and $y_{gi}$ is the "activity area" above the dose-response curve, Figure 4a. Given $\mathbf{x}_{gi}$ and a drug, we want accurately predict a patient's response to a drug and identifying genetic predictors of drug sensitivity. We use Cancer Cell Line Encyclopedia (CCLE) [3] which is a compilation $\sim$500 human cancer cell lines for 36 cancer types where their responses to 24 anticancer drugs have been measured.We perform two *experiments* where the number of cancers in each data set are $G = 2$ or 3 and we name them TWO and THREE experiments, respectively. We consider lung and blood[2] for TWO while for THREE we predict the drug sensitivity of skin, breast and ovary cancer cell lines. Beyond these five cancer types, others have less than 50 samples, so we remove them from consideration. Each experiment consists of 24 *problems* each corresponds to a drug. Not all of the 500 cell lines have been treated with all of the drugs. Therefore each problem has a different number of samples $n$ where $n \in [90, 160]$ for TWO and $n \in [70, 100]$ for THREE experiments. We perform a standard preprocessing [3] where we remove features with less than .2 absolute correlation with the response.The features that get removed vary by problem, so the dimension $p$ is reduced from from $> 30,000$ to $p \in [1000, 15000]$.

---

[2]By blood cancer, we mean any cancer originate from haematopoietic and lymphoid tissues.
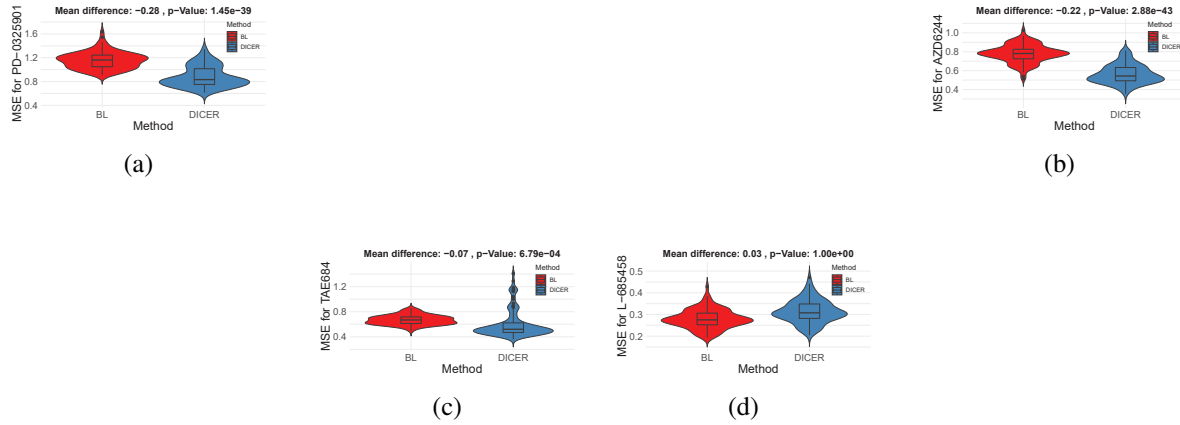
(a)                (b)



(c)          (d)

Figure 6: MSE for 100 bootstrapped dataset of four drugs of experiment TWO. (a),(b) Sample result of large difference between mean of MSEs and small p-values. (c) Smaller mean difference with significant p-value. (d) One of the three cases where DICER is outperformed by the baseline.

**Prediction:** In each TWO and THREE experiments, we predict the drug sensitivity for 24 different drugs using sparse DE estimator (2.3). Since the values of $d_g$ in constraint sets $\Omega_{f_g}(d_g)$ are unknown, we tune them by 5-fold cross-validation and report the mean squared error (MSE) of DICER and a baseline method. Our *baseline* method BL is the LASSO [37] equivalent of DE where we set $\forall g \in [G] \backslash d_g = 0$ and only estimate the common parameter $\beta_0$. Figure 5a and 5b illustrate the performance of DICER and BL for both experiments. Note that DICER outperforms BL in 21 and 18 out of 24 problems in TWO and THREE experiments, respectively.

To ensure that the prediction improvement of DICER over the baseline is statistically significant, we supplement our analysis with the bootstrapped error of both methods for the TWO experiment. For each problem in the TWO experiment, we generate 100 bootstrapped data sets by sampling with replacement as $\{(\mathbf{X}_{\text{TWO}}^{(i)}, \mathbf{y}_{\text{TWO}}^{(i)})\}_{i=1}^{100}$. Then, we fix $d_g$s hyper-parameters to values determined by cross-validation in the last stage and run both methods and compute pairs of MSEs as $\{(\text{MSE}_{\text{DICER}}^{(i)}, \text{MSE}_{\text{BL}}^{(i)})\}_{i=1}^{100}$ for each problem (drug). We perform paired t-test to determine if difference between means of two methods' MSEs is significant. In 21 out of 24 problems DICER's MSE is lesser than BL's with significance level of $\alpha = 0.05$. A representative set of results is demonstrated in Figure 6.

**Interpretation** We select Saracatinib, a drug which shows activity on both lung and blood cancer cell lines, Figure 4b. Then, during the bootstrap experiment on TWO, we record support of the estimated parameters by DICER. We pick the top five most frequently selected genes across 100 bootstrapped runs for further analysis. Now, we have three lists of genes for common, lung , and blood parameters. We perform gene enrichment analysis using ToppGene [15] to see where in functional/disease/drug databases these genes have been observed together with statistical significance. Table 1 summarizes a highlight of our findings which shows lung and blood parameters' supports are capturing a meaningful set of genes as a biomarkers.

| Blood and Lymph | | Lung | |
|---|---|---|---|
| Highlights | p-Val | Highlights | p-Val |
| Viral leukemogenesis | 6.17E-4 | Primary mucoepidermoid carcinoma of lung | 1.85E-4 |
| Primary cutaneous marginal zone B-cell lymphoma | 1.85E-3 | Lung carcinoma cell type unspecified stage IV | 1.85E-4 |
| Burkitt Lymphoma | 5.50E-3 | Primary adenocarcinoma of lung | 1.85E-4 |

Table 1: Highlights of interpretable outcomes of DICER. p-Values are computed by Fisher's exact test [15].

**Appendix A. Proofs of Theorems.** In this Section we present detail proof for each theorem and proposition. To avoid cluttering, during our proofs, we state some needed results as lemmas and provide their proof in the next Section B.

**A.1. Proof of Theorem 2.2.** —

**A.2. Proof of Proposition 3.4.**

*Proof.* Consider only one group for regression in isolation. Note that $\mathbf{y}_g = \mathbf{X}_g(\boldsymbol{\beta}_g^* + \boldsymbol{\beta}_0^*) + \boldsymbol{\omega}_g$ is a superposition model and as shown in [21] the sample complexity required for the RE condition and subsequently recovering $\boldsymbol{\beta}_0^*$ and $\boldsymbol{\beta}_g^*$ is $n_g \geq c(\max_{g \in [G]} \omega(\mathcal{A}_g) + \sqrt{\log 2})^2$. ∎

**A.3. Proof of Theorem 3.7.** Let's simplify the LHS of the RE condition:

$$\frac{1}{\sqrt{n}}\|\mathbf{X}\boldsymbol{\delta}\|_2 = \left(\frac{1}{n}\sum_{g=1}^{G}\sum_{i=1}^{n_g}|\langle\mathbf{x}_{gi}, \boldsymbol{\delta}_0 + \boldsymbol{\delta}_g\rangle|^2\right)^{\frac{1}{2}}$$

$$\geq \frac{1}{n}\sum_{g=1}^{G}\sum_{i=1}^{n_g}|\langle\mathbf{x}_{gi}, \boldsymbol{\delta}_0 + \boldsymbol{\delta}_g\rangle|$$

$$\geq \frac{1}{n}\sum_{g=1}^{G}\xi\|\boldsymbol{\delta}_0 + \boldsymbol{\delta}_g\|_2\sum_{i=1}^{n_g}\mathbb{1}\left(|\langle\mathbf{x}_{gi}, \boldsymbol{\delta}_0 + \boldsymbol{\delta}_g\rangle| \geq \xi\|\boldsymbol{\delta}_0 + \boldsymbol{\delta}_g\|_2\right),$$

where the first inequality is due to Lyapunov's inequality. To avoid cluttering we denote $\boldsymbol{\delta}_{0g} = \boldsymbol{\delta}_0 + \boldsymbol{\delta}_g$ where $\boldsymbol{\delta}_0 \in \mathcal{C}_0$ and $\boldsymbol{\delta}_g \in \mathcal{C}_g$. Now we add and subtract the corresponding per-group marginal tail function, $Q_{\xi_g}(\boldsymbol{\delta}_{0g}) = \mathbb{P}(|\langle\mathbf{x}, , \boldsymbol{\delta}_{0g}\rangle| > \xi_g)$ where $\xi_g > 0$. Let $\xi_g = \|\boldsymbol{\delta}_{0g}\|_2\xi$ then the LHS of the RE condition reduces to:

(A.1)
$$\inf_{\boldsymbol{\delta}\in\mathcal{H}}\frac{1}{\sqrt{n}}\|\mathbf{X}\boldsymbol{\delta}\|_2 \geq \inf_{\boldsymbol{\delta}\in\mathcal{H}}\sum_{g=1}^{G}\frac{n_g}{n}\xi_g Q_{2\xi_g}(\boldsymbol{\delta}_{0g})$$

$$- \sup_{\boldsymbol{\delta}\in\mathcal{H}}\frac{1}{n}\sum_{g=1}^{G}\xi_g\sum_{i=1}^{n_g}\left[Q_{2\xi_g}(\boldsymbol{\delta}_{0g}) - \mathbb{1}(|\langle\mathbf{x}_{gi}, \boldsymbol{\delta}_{0g}\rangle| \geq \xi_g)\right]$$

$$= t_1(\mathbf{X}) - t_2(\mathbf{X})$$

For the ease of exposition we have written the LHS of (A.1) as the difference of two terms, i.e., $t_1(\mathbf{X}) - t_2(\mathbf{X})$ and in the followings we lower bound the first term $t_1$ and upper bound the second term $t_2$.

**A.3.1. Lower Bounding the First Term.** Our main result is the following lemma which uses the DERIC condition of the Definition 3.5 and provides a lower bound for the first term $t_1(\mathbf{X})$:

**Lemma A.1.** *Suppose DERIC holds. Let $\psi_{\mathcal{I}} = \frac{\lambda_{\min}\bar{\rho}}{3}$. For any $\boldsymbol{\delta} \in \mathcal{H}$, we have:*

(A.2)
$$\sum_{g=1}^{G}\frac{n_g}{n}\xi_g Q_{2\xi_g}(\boldsymbol{\delta}_{0g}) \geq \psi_{\mathcal{I}}\xi\frac{(\alpha - 2\xi)^2}{4ck^2}\left(\|\boldsymbol{\delta}_0\|_2 + \sum_{g=1}^{n}\frac{n_g}{n}\|\boldsymbol{\delta}_g\|_2\right),$$

*which implies that $t_1(\mathbf{X}) = \inf_{\boldsymbol{\delta}\in\mathcal{H}}\sum_{g=1}^{G}\frac{n_G}{n}\xi_g Q_{2\xi_g}(\boldsymbol{\delta}_{0g})$ satisfies the same RHS bound of (A.2).*

**A.3.2. Upper Bounding the Second Term.** Let's focus on the second term, i.e., $t_2(\mathbf{X})$. First we want to show that the second term satisfies the bounded difference property defined in Section 3.2. of [6]. In other words, by changing each of $\mathbf{x}_{gi}$ the value of $t_2(\mathbf{X})$ at most change by one. First, we rewrite $t_2$ as follows:

$$h\left(\mathbf{x}_{11},\ldots,\mathbf{x}_{jk},\ldots,\mathbf{x}_{Gn_G}\right) = t_2\left(\mathbf{x}_{11},\ldots,\mathbf{x}_{jk},\ldots,\mathbf{x}_{Gn_G}\right) = \sup_{\boldsymbol{\delta}\in\mathcal{H}} g\left(\mathbf{x}_{11},\ldots,\mathbf{x}_{jk},\ldots,\mathbf{x}_{Gn_G}\right)$$

where $g\left(\mathbf{x}_{11},\ldots,\mathbf{x}_{jk},\ldots,\mathbf{x}_{Gn_G}\right) = \sum_{g=1}^{G}\frac{\xi_g}{n}\sum_{i=1}^{n_g}\left[Q_{2\xi_g}(\boldsymbol{\delta}_{0g}) - \mathbb{1}(|\langle\mathbf{x}_{gi},\boldsymbol{\delta}_{0g}\rangle| \geq \xi_g)\right]$. To avoid cluttering let's $\mathcal{X} = \{\mathbf{x}_{11},\ldots,\mathbf{x}_{jk},\ldots,\mathbf{x}_{Gn_G}\}$. We want to show that $t_2$ has the bounded difference property, meaning:

$$\sup_{\mathcal{X},\mathbf{x}'_{jk}} |h\left(\mathbf{x}_{11},\ldots,\mathbf{x}_{jk},\ldots,\mathbf{x}_{Gn_G}\right) - h\left(\mathbf{x}_{11},\ldots,\mathbf{x}'_{jk},\ldots,\mathbf{x}_{Gn_G}\right)| \leq c_i$$

for some constant $c_i$. Note that for bounded functions $f, g : \mathcal{X} \to \mathbb{R}$, we have $|\sup_{\mathcal{X}} f - \sup_{\mathcal{X}} g| \leq \sup_{\mathcal{X}} |f - g|$. Therefore:

$$\sup_{\mathcal{X},\mathbf{x}'_{jk}} |h\left(\mathbf{x}_{11},\ldots,\mathbf{x}_{jk},\ldots,\mathbf{x}_{Gn_G}\right) - h\left(\mathbf{x}_{11},\ldots,\mathbf{x}'_{jk},\ldots,\mathbf{x}_{Gn_G}\right)|$$

$$\leq \sup_{\mathcal{X},\mathbf{x}'_{jk}} \sup_{\boldsymbol{\delta}\in\mathcal{H}} |g\left(\mathbf{x}_{11},\ldots,\mathbf{x}_{jk},\ldots,\mathbf{x}_{Gn_G}\right) - g\left(\mathbf{x}_{11},\ldots,\mathbf{x}'_{jk},\ldots,\mathbf{x}_{Gn_G}\right)|$$

$$\leq \sup_{\mathcal{X},\mathbf{x}'_{jk}} \sup_{\boldsymbol{\delta}\in\mathcal{H}} \sup_{\mathbf{x}_{jk},\mathbf{x}'_{jk}} \frac{\xi_j}{n}\left(\mathbb{1}(|\langle\mathbf{x}'_{jk},\boldsymbol{\delta}_{0j}\rangle| \geq \xi_j) - \mathbb{1}(|\langle\mathbf{x}_{jk},\boldsymbol{\delta}_{0j}\rangle| \geq \xi_j)\right)$$

$$\leq \sup_{\mathcal{X},\mathbf{x}'_{jk}} \sup_{\boldsymbol{\delta}\in\mathcal{H}} \frac{\xi_j}{n}$$

$$= \frac{\xi}{n} \sup_{\boldsymbol{\delta}\in\mathcal{H}} \|\boldsymbol{\delta}_0 + \boldsymbol{\delta}_g\|_2$$

$$= \frac{\xi}{n} \sup_{\boldsymbol{\delta}\in\mathcal{H}} \|\boldsymbol{\delta}_0\|_2 + \|\boldsymbol{\delta}_g\|_2$$

$$(\boldsymbol{\delta}\in\mathcal{H}) = \xi\left(\frac{1}{n} + \frac{1}{n_g}\right)$$

$$\leq \frac{2\xi}{n}$$

Note that for $\boldsymbol{\delta}\in\mathcal{H}$ we have $\|\boldsymbol{\delta}_0\|_2 + \frac{n_g}{n}\|\boldsymbol{\delta}_g\|_2 \leq 1$ which results in $\|\boldsymbol{\delta}_0\|_2 \leq 1$ and $\|\boldsymbol{\delta}_g\|_2 \leq \frac{n}{n_g}$. Now, we can invoke the bounded difference inequality from Theorem 6.2 of [6] which says that with probability at least $1 - e^{-\tau^2/2}$ we have: $t_2(\mathbf{X}) \leq \mathbb{E}t_2(\mathbf{X}) + \frac{\tau}{\sqrt{n}}$.

Having this concentration bound, it is enough to bound the expectation of the second term. Following lemma provides us with the bound on the expectation.

**Lemma A.2.** *For the random vector $\mathbf{x}$ of Definition 3.1, we have the following bound:*

$$\frac{2}{n}\mathbb{E}\sup_{\boldsymbol{\delta}_{[G]}}\sum_{g=1}^{G}\xi_g\sum_{i=1}^{n_g}\left[Q_{2\xi_g}(\boldsymbol{\delta}_{0g}) - \mathbb{1}(|\langle\mathbf{x}_{gi},\boldsymbol{\delta}_{0g}\rangle| \geq \xi_g)\right] \leq \frac{2}{\sqrt{n}}\sum_{g=0}^{G}\sqrt{\frac{n_g}{n}}c_g k\omega(\mathcal{A}_g)\|\boldsymbol{\delta}_g\|_2$$

**A.3.3. Continuing the Proof of Theorem 3.7.** Set $n_0 = n$. Putting back bounds of $t_1(\mathbf{X})$ and $t_2(\mathbf{X})$ together from Lemma A.1 and A.2, with probability at least $1 - e^{-\frac{\tau^2}{2}}$ we have:

$$\inf_{\boldsymbol{\delta} \in \mathcal{H}} \frac{1}{\sqrt{n}} \|\mathbf{X}\boldsymbol{\delta}\|_2 \geq \sum_{g=0}^{G} \frac{n_g}{n} \psi_{\mathcal{I}} \xi \|\boldsymbol{\delta}_g\|_2 \frac{(\alpha - 2\xi)^2}{4ck^2} - \frac{2}{\sqrt{n}} \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} c_g k \omega(\mathcal{A}_g) \|\boldsymbol{\delta}_g\|_2 - \frac{\tau}{\sqrt{n}}$$

$$\left( q = \frac{(\alpha - 2\xi)^2}{4ck^2} \right) = \sum_{g=0}^{G} \frac{n_g}{n} \psi_{\mathcal{I}} \xi \|\boldsymbol{\delta}_g\|_2 q - \frac{2c}{\sqrt{n}} \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} k \omega(\mathcal{A}_g) \|\boldsymbol{\delta}_g\|_2 - \frac{\tau}{\sqrt{n}}$$

$$= n^{-1} \sum_{g=0}^{G} n_g \|\boldsymbol{\delta}_g\|_2 (\psi_{\mathcal{I}} \xi q - 2ck \frac{\omega(\mathcal{A}_g)}{\sqrt{n_g}}) - \frac{\tau}{\sqrt{n}}$$

$$\left( \kappa_g = \psi_{\mathcal{I}} \xi q - \frac{2ck\omega(\mathcal{A}_g)}{\sqrt{n_g}} \right) = \sum_{g=0}^{G} \frac{n_g}{n} \|\boldsymbol{\delta}_g\|_2 \kappa_g - \frac{\tau}{\sqrt{n}}$$

$$\geq \kappa_{\min} \sum_{g=0}^{G} \frac{n_g}{n} \|\boldsymbol{\delta}_g\|_2 - \frac{\tau}{\sqrt{n}}$$

$$(\boldsymbol{\delta} \in \mathcal{H}) = \kappa_{\min} - \frac{\tau}{\sqrt{n}}$$

where $\kappa_{\min} = \mathrm{argmin}_{g \in [G]} \kappa_g$. Note that all $\kappa_g$s should be bounded away from zero. To this end we need the follow sample complexities:

(A.3) $$\forall g \in [G] : \quad \left( \frac{2ck}{\psi_{\mathcal{I}} \xi q} \right)^2 \omega(\mathcal{A}_g)^2 \leq n_g$$

Taking $\xi = \frac{\alpha}{6}$ we can simplify the sample complexities to the followings:

(A.4) $$\forall g \in [G] : \quad \left( \frac{Ck^3}{\psi_{\mathcal{I}} \alpha^3} \right)^2 \omega(\mathcal{A}_g)^2 \leq n_g$$

Finally, to conclude, we take $\tau = \sqrt{n} \kappa_{\min}/2$.

**A.4. Proof of Theorem 4.1.**

*Proof.* From now on, to avoid cluttering the notation assume $\boldsymbol{\omega} = \boldsymbol{\omega}_0$. We massage the equation as follows:

$$\boldsymbol{\omega}^T \mathbf{X} \boldsymbol{\delta} = \sum_{g=0}^{G} \langle \mathbf{X}_g^T \boldsymbol{\omega}_g, \boldsymbol{\delta}_g \rangle = \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 \langle \mathbf{X}_g^T \frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2}, \frac{\boldsymbol{\delta}_g}{\|\boldsymbol{\delta}_g\|_2} \rangle \sqrt{\frac{n}{n_g}} \|\boldsymbol{\omega}_g\|_2$$

Assume $b_g = \langle \mathbf{X}_g^T \frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2}, \frac{\boldsymbol{\delta}_g}{\|\boldsymbol{\delta}_g\|_2} \rangle \sqrt{\frac{n}{n_g}} \|\boldsymbol{\omega}_g\|_2$ and $a_g = \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2$. Then the above term is the inner product of two vectors $\mathbf{a} = (a_0, \ldots, a_G)$ and $\mathbf{b} = (b_0, \ldots, b_G)$ for which we have:

$$\sup_{\mathbf{a} \in \mathcal{H}} \mathbf{a}^T \mathbf{b} = \sup_{\|\mathbf{a}\|_1 = 1} \mathbf{a}^T \mathbf{b}$$

$$\text{(definition of the dual norm)} \leq \|\mathbf{b}\|_\infty$$

$$= \max_{g \in [G]} b_g$$

Now we can go back to the original form:

(A.5)
$$\sup_{\boldsymbol{\delta} \in \mathcal{H}} \boldsymbol{\omega}^T \mathbf{X} \boldsymbol{\delta} \leq \max_{g \in [G]} \langle \mathbf{X}_g^T \frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2}, \frac{\boldsymbol{\delta}_g}{\|\boldsymbol{\delta}_g\|_2} \rangle \sqrt{\frac{n}{n_g}} \|\boldsymbol{\omega}_g\|_2$$

$$\leq \max_{g \in [G]} \sqrt{\frac{n}{n_g}} \|\boldsymbol{\omega}_g\|_2 \sup_{\mathbf{u}_g \in \mathcal{C}_g \cap \mathbb{S}^{p-1}} \langle \mathbf{X}_g^T \frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2}, \mathbf{u}_g \rangle$$

To avoid cluttering we name $h_g(\boldsymbol{\omega}_g, \mathbf{X}_g) = \|\boldsymbol{\omega}_g\|_2 \sup_{\mathbf{u}_g \in \mathcal{A}_g} \langle \mathbf{X}_g^T \frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2}, \mathbf{u}_g \rangle$ and $e_g(\tau) = \sqrt{(2K^2 + 1)n_g} \left( v_g C_g k \omega(\mathcal{A}_g)\right.$
Then from (A.5), we have:

$$\mathbb{P}\left(\frac{2}{n} \sup_{\boldsymbol{\delta} \in \mathcal{H}} \boldsymbol{\omega}^T \mathbf{X} \boldsymbol{\delta} > \frac{2}{n} \max_{g \in [G]} \sqrt{\frac{n}{n_g}} e_g(\tau)\right) \leq \mathbb{P}\left(\frac{2}{n} \max_{g \in [G]} \sqrt{\frac{n}{n_g}} h_g(\boldsymbol{\omega}_g, \mathbf{X}_g) > \frac{2}{n} \max_{g \in [G]} \sqrt{\frac{n}{n_g}} e_g(\tau)\right)$$

To simplify the notation, we drop arguments of $h_g$ for now. From the union bound we have:

$$\mathbb{P}\left(\frac{2}{n} \max_{g \in [G]} \sqrt{\frac{n}{n_g}} h_g > \frac{2}{n} \max_{g \in [G]} \sqrt{\frac{n}{n_g}} e_g(\tau)\right) \leq \sum_{g=0}^{G} \mathbb{P}\left(h_g > \max_{g \in [G]} e_g(\tau)\right)$$

$$\leq \sum_{g=0}^{G} \mathbb{P}\left(h_g > e_g(\tau)\right)$$

$$\leq (G+1) \max_{g \in [G]} \mathbb{P}\left(h_g > e_g(\tau)\right)$$

$$\leq \sigma \exp\left(-\min_{g \in [G]}\left[\nu_g n_g - \log(G+1), \frac{\tau^2}{\eta_g^2 k^2}\right]\right)$$

where $\sigma = \max_{g \in [G]} \sigma_g$ and the last inequality is a result of the following lemma:

**Lemma A.3.** *For $\mathbf{x}_{gi}$ and $\omega_{gi}$ defined in Definition 3.1 and $\tau > 0$, with probability at least*
$1 - \frac{\sigma_g}{(G+1)} \exp\left(-\min\left[\nu_g n_g - \log(G+1), \frac{\tau^2}{\eta_g^2 k^2}\right]\right)$ *we have:*

$$\sqrt{\frac{n}{n_g}} \|\boldsymbol{\omega}_g\|_2 \sup_{\mathbf{u}_g \in \mathcal{A}_g} \langle \mathbf{X}_g^T \frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2}, \mathbf{u}_g \rangle \leq \sqrt{(2K^2 + 1)n} \left(\zeta_g k \omega(\mathcal{A}_g) + \epsilon_g \sqrt{\log(G+1)} + \tau\right),$$

*where $\sigma_g, \eta_g, \zeta_g$ and $\epsilon_g$ are group dependent constants.*

### A.5. Proof of Theorem 5.2.

*Proof.* In the following lemma we establish a recursive relation between errors of consecutive iterations which leads to a bound for the $t$th iteration.

**Lemma A.4.** *We have the following recursive dependency between the error of $t + 1$th iteration and $t$th iteration of DE:*

$$\|\boldsymbol{\delta}_g^{(t+1)}\|_2 \leq \left(\rho_g(\mu_g)\|\boldsymbol{\delta}_g^{(t)}\|_2 + \xi_g(\mu_g)\|\boldsymbol{\omega}_g\|_2 + \phi_g(\mu_g)\|\boldsymbol{\delta}_0^{(t)}\|_2\right)$$

$$\|\boldsymbol{\delta}_0^{(t+1)}\|_2 \leq \left(\rho_0(\mu_0)\|\boldsymbol{\delta}_0^{(t)}\|_2 + \xi_0(\mu_0)\|\boldsymbol{\omega}_0\|_2 + \mu_0 \sum_{g=1}^{G} \frac{\phi_g(\mu_g)}{\mu_g}\|\boldsymbol{\delta}_g^{(t)}\|_2\right)$$

By recursively applying the result of Lemma A.4, we get the following deterministic bound which depends on constants defined in Definition 5.1:

$$b_{t+1} = \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g^{(t+1)}\|_2 \leq \left( \rho_0 + \sum_{g=1}^{G} \sqrt{\frac{n_g}{n}} \phi_g \right) \|\boldsymbol{\delta}_0^{(t)}\|_2 + \sum_{g=1}^{G} \left( \sqrt{\frac{n_g}{n}} \rho_g + \mu_0 \frac{\phi_g}{\mu_g} \right) \|\boldsymbol{\delta}_g^{(t)}\|_2 + \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \xi_g \|\boldsymbol{\omega}_g\|_2$$

$$(A.6) \qquad \leq \rho \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g^{(t)}\|_2 + \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \xi_g \|\boldsymbol{\omega}_g\|_2$$

where $\rho = \max \left( \rho_0 + \sum_{g=1}^{G} \sqrt{\frac{n_g}{n}} \phi_g, \max_{g \in [G]} \left[ \rho_g + \sqrt{\frac{n}{n_g}} \frac{\mu_0}{\mu_g} \phi_g \right] \right)$. We have:

$$b_{t+1} \leq \rho b_t + \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \xi_g \|\boldsymbol{\omega}_g\|_2$$

$$\leq (\rho)^2 b_{t-1} + (\rho + 1) \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \xi_g \|\boldsymbol{\omega}_g\|_2$$

$$\leq (\rho)^t b_1 + \left( \sum_{i=0}^{t-1} (\rho)^i \right) \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \xi_g \|\boldsymbol{\omega}_g\|_2$$

$$= (\rho)^t \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\beta}_g^1 - \boldsymbol{\beta}_g^*\|_2 + \left( \sum_{i=0}^{t-1} (\rho)^i \right) \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \xi_g \|\boldsymbol{\omega}_g\|_2$$

$$(\boldsymbol{\beta}^1 = 0) \leq (\rho)^t \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\beta}_g^*\|_2 + \frac{1 - (\rho)^t}{1 - \rho} \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \xi_g \|\boldsymbol{\omega}_g\|_2$$

## A.6. Proof of Theorem 5.3.

*Proof.* First we need following two lemmas which are proved separately in the following sections.

**Lemma A.5.** *Consider $a_g \geq 1$, with probability at least $1 - 6 \exp \left( -\gamma_g (\omega(\mathcal{A}_g) + \tau)^2 \right)$ the following upper bound holds:*

$$(A.7) \qquad \rho_g \left( \frac{1}{a_g n_g} \right) \leq \frac{1}{2} \left[ \left( 1 - \frac{1}{a_g} \right) + \sqrt{2} c_g \frac{2\omega(\mathcal{A}_g) + \tau}{a_g \sqrt{n_g}} \right]$$

**Lemma A.6.** *Consider $a_g \geq 1$, with probability at least $1 - 4 \exp \left( -\gamma_g (\omega(\mathcal{A}_g) + \tau)^2 \right)$ the following upper bound holds:*

$$(A.8) \qquad \phi_g \left( \frac{1}{a_g n_g} \right) \leq \frac{1}{a_g} \left( 1 + c_{0g} \frac{\omega(\mathcal{A}_g) + \omega(\mathcal{A}_0) + 2\tau}{\sqrt{n_g}} \right)$$

Note that Lemma A.3 readily provides a high probability upper bound for $\eta_g(1/(a_g n_g))$ as $\sqrt{(2K^2 + 1)} \left( \zeta_g k \omega(\mathcal{A}_g) + \epsilon_g \sqrt{\log} \right.$

Starting from the deterministic form of the bound in Theorem 5.2 and putting in the step sizes as $\mu_g = \frac{1}{n_g a_g}$:

$$\text{(A.9)} \qquad \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g^{(t+1)}\|_2 \le (\rho)^t \sum_{g=0}^{G} \|\boldsymbol{\beta}_g^*\|_2 + \frac{1-(\rho)^t}{1-\rho} \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \eta_g \left(\frac{1}{n_g a_g}\right) \|\boldsymbol{\omega}_g\|_2,$$

where

$$\rho(a_0(\text{A.10})u_G) = \max\left(\rho_0\left(\frac{1}{na_0}\right) + \sum_{g=1}^{G} \sqrt{\frac{n_g}{n}} \phi_g\left(\frac{1}{n_g a_g}\right), \max_{g \in [G]} \rho_g\left(\frac{1}{n_g a_g}\right) + \sqrt{\frac{n}{n_g}} \frac{\mu_0}{\mu_g} \phi_g\left(\frac{1}{n_g a_g}\right)\right)$$

Remember the following two results to upper bound $\rho_g$s and $\phi_g$s from Lemmas A.5 and A.6:

$$\rho_g\left(\frac{1}{a_g n_g}\right) \le \frac{1}{2}\left[\left(1 - \frac{1}{a_g}\right) + \sqrt{2} c_g \frac{2\omega(\mathcal{A}_g) + \tau}{a_g \sqrt{n_g}}\right], \quad \text{w.p.} \quad 1 - 6\exp\left(-\gamma_g(\omega(\mathcal{A}_g) + \tau)^2\right)$$

$$\phi_g\left(\frac{1}{a_g n_g}\right) \le \frac{1}{a_g}\left(1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}}\right), \quad \text{w.p.} \quad 1 - 4\exp\left(-\gamma_g(\omega(\mathcal{A}_g) + \tau)^2\right)$$

First we want to keep $\rho_0 + \sum_{g=1}^{G} \sqrt{\frac{n_g}{n}} \phi_g$ of (A.10) strictly below 1.

$$\rho_0\left(\frac{1}{a_0 n}\right) + \sum_{g=1}^{G} \sqrt{\frac{n_g}{n}} \phi_g\left(\frac{1}{a_g n_g}\right) \le \frac{1}{2}\left[\left(1 - \frac{1}{a_0}\right) + \sqrt{2} c_0 \frac{2\omega_0 + \tau}{a_0 \sqrt{n}}\right]$$

$$+ \frac{1}{2} \sum_{g=1}^{G} \frac{2}{a_g} \sqrt{\frac{n_g}{n}} \left(1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}}\right)$$

Remember that $a_g \ge 1$ was arbitrary. So we pick it as $a_g = 2\sqrt{\frac{n}{n_g}}\left(1 + c_{0g} \frac{\omega_{0g}+\tau}{\sqrt{n_g}}\right)/b_g$ where $b_g \le 2\sqrt{\frac{n}{n_g}}\left(1 + c_{0g} \frac{\omega_{0g}+\tau}{\sqrt{n_g}}\right)$ (because we need $a_g \ge 1$) and the condition becomes:

$$\rho_0\left(\frac{1}{a_0 n}\right) + \sum_{g=1}^{G} \sqrt{\frac{n_g}{n}} \phi_g\left(\frac{1}{a_g n_g}\right) \le \frac{1}{2}\left[\left(1 - \frac{1}{a_0}\right) + \sqrt{2} c_0 \frac{2\omega(\mathcal{A}_0) + \tau}{a_0 \sqrt{n}}\right] + \frac{1}{2} \sum_{g=1}^{G} \frac{n_g}{n} b_g \le 1$$

We want to upper bound the RHS by $1/\theta_f$ which will determine the sample complexity for the shared component:

$$\text{(A.11)} \qquad \sqrt{2} c_0 \frac{2\omega(\mathcal{A}_0) + \tau}{\sqrt{n}} \le a_0\left(1 - \sum_{g=1}^{G} \frac{n_g}{n} b_g\right) + 1$$

Note that any lower bound on the RHS of (A.11) will lead to the correct sample complexity for which the coefficient of $\|\boldsymbol{\delta}_0^{(t)}\|_2$ (determined in (A.10)) will be below one. Since $a_0 \ge 1$ we can ignore the

first term by assuming $\max_{g \in [G]_\backslash} b_g \leq 1$ and the condition becomes:

$$n > 2c_0^2 (2\omega(\mathcal{A}_0) + \tau)^2, \forall g \in [G]_\backslash : a_g = 2b_g^{-1} \sqrt{\frac{n}{n_g}} \left( 1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}} \right),$$

$$a_0 \geq 1, 0 < b_g \leq 2\sqrt{\frac{n}{n_g}} \left( 1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}} \right), \max_{g \in [G]_\backslash} b_g \leq 1,$$

which can be simplified to:

(A.12) $\qquad n > 2c_0^2 (2\omega(\mathcal{A}_0) + \tau)^2, a_0 \geq 1,$

$$\forall g \in [G]_\backslash : a_g = 2b_g^{-1} \sqrt{\frac{n}{n_g}} \left( 1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}} \right), 0 < b_g \leq 1$$

Secondly, we want to bound all of $\rho_g + \mu_0 \sqrt{\frac{n}{n_g} \frac{\phi_g}{\mu_g}}$ terms of (A.10) for $\mu_g = \frac{1}{a_g n_g}$ by 1:

(A.13) $\qquad \rho_g \left( \frac{1}{n_g a_g} \right) + \sqrt{\frac{n}{n_g} \frac{\mu_0}{\mu_g}} \phi_g \left( \frac{1}{n_g a_g} \right) = \rho_g \left( \frac{1}{n_g a_g} \right) + \sqrt{\frac{n_g}{n} \frac{a_g}{a_0}} \phi_g \left( \frac{1}{n_g a_g} \right)$

$$= \frac{1}{2} \left[ \left[ \left( 1 - \frac{1}{a_g} \right) + \sqrt{2} c_g \frac{2\omega_g + \tau}{a_g \sqrt{n_g}} \right] \right.$$

$$\left. + \frac{2}{a_0} \sqrt{\frac{n_g}{n}} \left( 1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}} \right) \right]$$

$$\leq 1$$

The condition becomes:

(A.14) $\qquad \sqrt{2} c_g \frac{2\omega_g + \tau}{\sqrt{n_g}} \leq a_g + 1 - \sqrt{\frac{n_g}{n}} \frac{2a_g}{a_0} \left( 1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}} \right)$

Remember that we chose $a_g = 2b_g^{-1} \sqrt{\frac{n}{n_g}} \left( 1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}} \right)$. We substitute the value of $a_g$ by keeping in mind the constraints for the $b_g$ and the condition reduces to:

(A.15) $\qquad \sqrt{2} c_g \frac{2\omega_g + \tau}{d_g} \leq \sqrt{n_g}, \quad d_g := a_g + 1 - \frac{4}{b_g a_0} \left( 1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}} \right)^2$

for $d_g > 0$. Note that any positive lower bound of the $d_g$ will satisfy the condition in (A.15) and the result is a valid sample complexity. In the following we show that $d_g > 1$. We have $a_0 \geq 1$ condition from (A.12), so we take $a_0 = 4 \max_{g \in [G]_\backslash} \left( 1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}} \right)^2$ and look for a lower bound for $d_g$:

(A.16) $\qquad d_g \geq a_g + 1 - b_g^{-1}$

$$(a_g \text{ from (A.12)}) = 2b_g^{-1} \sqrt{\frac{n}{n_g}} \left( 1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}} \right) + 1 - b_g^{-1}$$

(A.17) $\qquad = 1 + b_g^{-1} \left[ 2\sqrt{\frac{n}{n_g}} \left( 1 + c_{0g} \frac{\omega_{0g} + \tau}{\sqrt{n_g}} \right) - 1 \right]$

The term inside of the last bracket (A.17) is always positive and therefore a lower bound is one, i.e., $d_g \geq 1$. From the condition (A.15) we get the following sample complexity:

$$(A.18) \qquad n_g > 2c_g^2(2\omega_g + \tau)^2$$

Now we need to determine $b_g$ from previous conditions (A.12), knowing that $a_0 = 4\max_{g \in [G]_\backslash} \left(1 + c_{0g}\frac{\omega_{0g}+\tau}{\sqrt{n_g}}\right)^2$
We have $0 < b_g \leq 1$ in (A.12) and we take the largest step by setting $b_g = 1$.

    Here we summarize the setting under which we have the linear convergence:

$$n > 2c_0^2\left(2\omega(\mathcal{A}_0) + \tau\right)^2, \forall g \in [G]_\backslash : n_g \geq 2c_g^2(2\omega(\mathcal{A}_g) + \tau)^2$$

$$(A.19) \qquad a_0 = 4\max_{g \in [G]_\backslash}\left(1 + c_{0g}\frac{\omega_{0g}+\tau}{\sqrt{n_g}}\right)^2, a_g = 2\sqrt{\frac{n}{n_g}}\left(1 + c_{0g}\frac{\omega_{0g}+\tau}{\sqrt{n_g}}\right)$$

$$\mu_0 = \frac{1}{4n} \times \frac{1}{\max_{g \in [G]_\backslash}\left(1 + c_{0g}\frac{\omega_{0g}+\tau}{\sqrt{n_g}}\right)^2}, \mu_g = \frac{1}{2\sqrt{nn_g}}\left(1 + c_{0g}\frac{\omega_{0g}+\tau}{\sqrt{n_g}}\right)^{-1}$$

    Now we rewrite the same analysis using the tail bounds for the coefficients to clarify the probabilities. To simplify the notation, let $r_{g1} = \frac{1}{2}\left[\left(1 - \frac{1}{a_g}\right) + \sqrt{2}c_g\frac{2\omega(\mathcal{A}_g)+\tau}{a_g\sqrt{n_g}}\right]$ and $r_{g2} = \frac{1}{a_g}\left(1 + c_{0g}\frac{\omega_{0g}+\tau}{\sqrt{n_g}}\right)$ and $r_0(\tau) = r_{01} + \sum_{g=1}^G \sqrt{\frac{n_g}{n}}r_{g2}$ and $r_g(\tau) = r_{g1} + \sqrt{\frac{n_g}{n}}\frac{a_g}{a_0}r_{g2}, \forall g \in [G]_\backslash$, and $r(\tau) = \max_{g\in[G]} r_g$. All of which are computed using $a_g$s specified in (A.19). Basically $r$ is an instantiation of an upper bound of the $\rho$ defined in (A.10) using $a_g$s in (A.19).

    We are interested to upper bound the following probability:

$$\mathbb{P}\left(\sum_{g=0}^G \sqrt{\frac{n_g}{n}}\|\boldsymbol{\delta}_g^{(t+1)}\|_2 \geq r(\tau)^t \sum_{g=0}^G \sqrt{\frac{n_g}{n}}\|\boldsymbol{\beta}_g^*\|_2 + \frac{(G+1)\sqrt{(2K^2+1)}}{(1-r(\tau))\sqrt{n}}\left(\zeta k \max_{g\in[G]}\omega(\mathcal{A}_g) + \tau\right)\right)$$

$$\leq \mathbb{P}\left((\rho)^t \sum_{g=0}^G \sqrt{\frac{n_g}{n}}\|\boldsymbol{\beta}_g^*\|_2 + \frac{1-(\rho)^t}{1-\rho}\sum_{g=0}^G \sqrt{\frac{n_g}{n}}\eta_g\left(\frac{1}{n_g a_g}\right)\|\boldsymbol{\omega}_g\|_2\right)$$

$$\geq r(\tau)^t \sum_{g=0}^G \sqrt{\frac{n_g}{n}}\|\boldsymbol{\beta}_g^*\|_2 + \frac{(G+1)\sqrt{(2K^2+1)}}{(1-r(\tau))\sqrt{n}}\left(\zeta k \max_{g\in[G]}\omega(\mathcal{A}_g) + \tau\right)\right)$$

$$\leq \mathbb{P}\left(\rho \geq r(\tau)\right)$$

$$+\mathbb{P}\left(\frac{1}{1-\rho}\sum_{g=0}^G \sqrt{n_g}\eta_g\left(\frac{1}{n_g a_g}\right)\|\boldsymbol{\omega}_g\|_2 \geq \frac{(G+1)\sqrt{(2K^2+1)}}{(1-r(\tau))}\left(\zeta k \max_{g\in[G]}\omega(\mathcal{A}_g) + \tau\right)\right) \qquad (A.20)$$

where the first inequality comes from the deterministic bound of (A.9), We first focus on bounding the

first term $\mathbb{P}\left(\rho \geq r(\tau)\right)$:

$$\mathbb{P}\left(\rho \geq r(\tau)\right)$$

$$= \mathbb{P}\left(\max\left(\rho_0\left(\frac{1}{na_0}\right) + \sum_{g=1}^{G}\sqrt{\frac{n_g}{n}}\phi_g\left(\frac{1}{n_ga_g}\right), \max_{g\in[G]}\rho_g\left(\frac{1}{n_ga_g}\right) + \sqrt{\frac{n}{n_g}}\frac{\mu_0}{\mu_g}\phi_g\left(\frac{1}{n_ga_g}\right)\right) \geq \max_{g\in[G]}r(\tau)\right)$$

$$\leq \mathbb{P}\left(\rho_0\left(\frac{1}{na_0}\right) + \sum_{g=1}^{G}\sqrt{\frac{n_g}{n}}\phi_g\left(\frac{1}{n_ga_g}\right) \geq r_0\right) + \sum_{g=1}^{G}\mathbb{P}\left(\rho_g\left(\frac{1}{n_ga_g}\right) + \sqrt{\frac{n}{n_g}}\frac{\mu_0}{\mu_g}\phi_g\left(\frac{1}{n_ga_g}\right) \geq r_g\right)$$

$$\leq \mathbb{P}\left(\rho_0\left(\frac{1}{na_0}\right) \geq r_{01}\right) + \sum_{g=1}^{G}\mathbb{P}\left(\phi_g\left(\frac{1}{n_ga_g}\right) \geq r_{g2}\right) + \sum_{g=1}^{G}\left[\mathbb{P}\left(\rho_g\left(\frac{1}{n_ga_g}\right) \geq r_{g1}\right) + \mathbb{P}\left(\phi_g\left(\frac{1}{n_ga_g}\right) \geq r_{g2}\right)\right]$$

$$\leq \sum_{g=0}^{G}\mathbb{P}\left(\rho_g\left(\frac{1}{n_ga_g}\right) \geq r_{g1}\right) + 2\sum_{g=1}^{G}\mathbb{P}\left(\phi_g\left(\frac{1}{n_ga_g}\right) \geq r_{g2}\right)$$

$$\leq \sum_{g=0}^{G}6\exp\left(-\gamma_g(\omega(\mathcal{A}_g) + \tau)^2\right) + 2\sum_{g=1}^{G}4\exp\left(-\gamma_g(\omega(\mathcal{A}_g) + \tau)^2\right)$$

$$\leq 6(G+1)\exp\left(-\gamma\min_{g\in[G]}(\omega(\mathcal{A}_g) + \tau)^2\right) + 8G\exp\left(-\gamma\min_{g\in[G]_\setminus}(\omega(\mathcal{A}_g) + \tau)^2\right)$$

$$\leq 14(G+1)\exp\left(\overset{(A.21)}{}-\gamma\min_{g\in[G]}(\omega(\mathcal{A}_g) + \tau)^2\right)$$

Now we focus on bounding the second term:

$$\mathbb{P}\left(\frac{1}{1-\rho}\sum_{g=0}^{G}\sqrt{n_g}\eta_g\left(\frac{1}{n_ga_g}\right)\|\boldsymbol{\omega}_g\|_2 \geq \frac{(G+1)\sqrt{(2K^2+1)}}{(1-r(\tau))}\left(\zeta k\max_{g\in[G]}\omega(\mathcal{A}_g) + \tau\right)\right)$$

$$\leq \mathbb{P}\left(\frac{1}{1-\rho}\sum_{g=0}^{G}\sqrt{n_g}\eta_g\left(\frac{1}{n_ga_g}\right)\|\boldsymbol{\omega}_g\|_2 \geq \frac{1}{(1-r(\tau))}\sum_{g=0}^{G}\sqrt{(2K^2+1)}\left(\zeta_g k\omega(\mathcal{A}_g) + \tau\right)\right)$$

$$\leq \mathbb{P}\left(\sum_{g=0}^{G}\sqrt{n_g}\eta_g\left(\frac{1}{n_ga_g}\right)\|\boldsymbol{\omega}_g\|_2 \geq \sum_{g=0}^{G}\sqrt{(2K^2+1)}\left(\zeta_g k\omega(\mathcal{A}_g) + \tau\right)\right) + \mathbb{P}\left(\rho \geq r(\tau)\right)$$

$$\overset{(A.22)}{\leq}\sum_{g=0}^{G}\mathbb{P}\left(\sqrt{n_g}\eta_g\left(\frac{1}{n_ga_g}\right)\|\boldsymbol{\omega}_g\|_2 \geq \sqrt{(2K^2+1)}\left(\zeta_g k\omega(\mathcal{A}_g) + \tau\right)\right) + \mathbb{P}\left(\rho \geq r(\tau)\right)$$

Focusing on the summand of the first term, remember from Definition 5.1 that $\eta_g(\mu_g) = \frac{1}{a_gn_g}\sup_{\mathbf{v}\in\mathcal{B}_g}\mathbf{v}^T\mathbf{X}_g^T\frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2}$, $g \in$ ▮
$[G]$ and $a_g \geq 1$:

$$\mathbb{P}\left(\overset{(A.23)}{\|\boldsymbol{\omega}_g\|_2}\sup_{\mathbf{v}\in\mathcal{B}_g}\mathbf{v}^T\mathbf{X}_g^T\frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2} \geq a_g\sqrt{(2K^2+1)n_g}\left(\zeta_g k\omega(\mathcal{A}_g) + \tau\right)\right) \leq \sigma_g\exp\left(-\min\left[\nu_g n_g, \frac{\tau^2}{\eta_g^2 k^2}\right]\right) ▮$$

where we used the intermediate form of Lemma A.3 for $\tau > 0$. Putting all of the bounds (A.21), (A.22), and (A.23) back into the (A.20):

$$\sigma_g(G+1)\exp\left(-\min_{g\in[G]}\left(\min\left[\nu_g n_g, \frac{\tau^2}{\eta_g^2 k^2}\right]\right)\right) + 28(G+1)\exp\left(-\gamma\min_{g\in[G]}(\omega(\mathcal{A}_g)+\tau)^2\right)$$

$$\leq \upsilon\exp\left[\min_{g\in[G]}\left(-\min\left[\nu_g n_g - \log G, \gamma(\omega(\mathcal{A}_g)+t)^2, \frac{t^2}{\eta_g^2 k^2}\right]\right)\right]$$

where $\upsilon = \max(28, \sigma)$ and $\gamma = \min_{g\in[G]}\gamma_g$ and $\tau = t + \max(\epsilon, \gamma^{-1/2})\sqrt{\log(G+1)}$ where $\epsilon = k\max_{g\in[G]}\eta_g$. Note that $\tau = t + C\sqrt{\log(G+1)}$ increases the sample complexities to the followings:

$$n > 2c_0^2\left(2\omega(\mathcal{A}_0) + C\sqrt{\log(G+1)} + t\right)^2, \forall g\in[G]_\backslash : n_g \geq 2c_g^2(2\omega(\mathcal{A}_g) + C\sqrt{\log(G+1)} + t)^2$$

and it also affects step sizes as follows:

$$\mu_0 = \frac{1}{4n}\times\min_{g\in[G]_\backslash}\left(1 + c_{0g}\frac{\omega_{0g} + C\sqrt{\log(G+1)} + t}{\sqrt{n_g}}\right)^{-2}, \mu_g = \frac{1}{2\sqrt{nn_g}}\left(1 + c_{0g}\frac{\omega_{0g} + C\sqrt{\log(G+1)} + t}{\sqrt{n_g}}\right)^{-1}$$

## Appendix B. Proofs of Lemmas.

Here, we present proofs of each lemma used during the proofs of theorems in Section A.

### B.1. Proof of Lemma A.1.

*Proof.* LHS of (A.2) is the weighted summation of $\xi_g Q_{2\xi_g}(\boldsymbol{\delta}_{0g}) = \|\boldsymbol{\delta}_{0g}\|_2\xi\mathbb{P}(|\langle\mathbf{x}, ,\boldsymbol{\delta}_{0g}/\|\boldsymbol{\delta}_{0g}\|_2\rangle| > \blacksquare$ $2\xi) = \|\boldsymbol{\delta}_{0g}\|_2\xi Q_{2\xi}(\mathbf{u})$ where $\xi > 0$ and $\mathbf{u} = \boldsymbol{\delta}_{0g}/\|\boldsymbol{\delta}_{0g}\|_2$ is a unit length vector. So we can rewrite the LHS of (A.2) as:

$$\sum_{g=1}^{G}\frac{n_g}{n}\xi_g Q_{2\xi_g}(\boldsymbol{\delta}_{0g}) = \sum_{g=1}^{G}\frac{n_g}{n}\|\boldsymbol{\delta}_0 + \boldsymbol{\delta}_g\|_2\xi Q_{2\xi}(\mathbf{u})$$

With this observation, the lower bound of the Lemma A.1 is a direct consequence of the following two results:

**Lemma B.1.** *Let* $\mathbf{u}$ *be any unit length vector and suppose* $\mathbf{x}$ *obeys Definiton 3.1. Then for any* $\mathbf{u}$, *we have*

(B.1) $$Q_{2\xi}(\mathbf{u}) \geq \frac{(\alpha - 2\xi)^2}{4ck^2}.$$

**Lemma B.2.** *Suppose Definition 3.5 holds. Then, we have:*

(B.2) $$\sum_{i=1}^{G}n_i\|\boldsymbol{\delta}_0 + \boldsymbol{\delta}_i\|_2 \geq \frac{\bar{\rho}\lambda_{\min}}{3}\left(Gn\|\boldsymbol{\delta}_0\|_2 + \sum_{i=1}^{G}n_i\|\boldsymbol{\delta}_i\|_2\right), \quad \forall i\in[G] : \boldsymbol{\delta}_i \in \mathcal{C}_i.$$

### B.2. Proof of Lemma A.2.

*Proof.* Consider the following soft indicator function which we use in our derivation:

$$
\psi_a(s) = \begin{cases} 0, & |s| \le a \\ (|s| - a)/a, & a \le |s| \le 2a \\ 1, & 2a < |s| \end{cases}
$$

Now:

$$
\mathbb{E} \sup_{\boldsymbol{\delta}_{[G]}} \sum_{g=1}^{G} \xi_g \sum_{i=1}^{n_g} \left[ Q_{2\xi_g}(\boldsymbol{\delta}_{0g}) - \mathbb{1}(|\langle \mathbf{x}_{gi}, \boldsymbol{\delta}_{0g} \rangle| \ge \xi_g) \right]
$$

$$
= \mathbb{E} \sup_{\boldsymbol{\delta}_{[G]}} \sum_{g=1}^{G} \xi_g \sum_{i=1}^{n_g} \left[ \mathbb{E}\mathbb{1}(|\langle \mathbf{x}_{gi}, \boldsymbol{\delta}_{0g} \rangle| \ge 2\xi_g) - \mathbb{1}(|\langle \mathbf{x}_{gi}, \boldsymbol{\delta}_{0g} \rangle| \ge \xi_g) \right]
$$

$$
\le \mathbb{E} \sup_{\boldsymbol{\delta}_{[G]}} \sum_{g=1}^{G} \xi_g \sum_{i=1}^{n_g} \left[ \mathbb{E}\psi_{\xi_g}(\langle \mathbf{x}, \boldsymbol{\delta}_{0g} \rangle) - \psi_{\xi_g}(\langle \mathbf{x}_{gi}, \boldsymbol{\delta}_{0g} \rangle) \right]
$$

$$
\le 2\mathbb{E} \sup_{\boldsymbol{\delta}_{[G]}} \sum_{g=1}^{G} \xi_g \sum_{i=1}^{n_g} \epsilon_{gi} \psi_{\xi_g}(\langle \mathbf{x}_{gi}, \boldsymbol{\delta}_{0g} \rangle)
$$

$$
\le 2\mathbb{E} \sup_{\boldsymbol{\delta}_{[G]}} \sum_{g=1}^{G} \sum_{i=1}^{n_g} \epsilon_{gi} \langle \mathbf{x}_{gi}, \boldsymbol{\delta}_{0g} \rangle
$$

where $\epsilon_{gi}$ are iid copies of Rademacher random variable which are independent of every other random

variables and themselves. Now we add back $\frac{1}{n}$ and expand $\boldsymbol{\delta}_{0g} = \boldsymbol{\delta}_0 + \boldsymbol{\delta}_g$:

$$\frac{2}{n}\mathbb{E}\sup_{\boldsymbol{\delta}_{[G]}\in\mathcal{C}_{[G]}}\sum_{g=1}^{G}\sum_{i=1}^{n_g}\epsilon_{gi}\langle\mathbf{x}_{gi},\boldsymbol{\delta}_{0g}\rangle = \frac{2}{n}\mathbb{E}\sup_{\boldsymbol{\delta}_0\in\mathcal{C}_0}\sum_{i=1}^{n}\epsilon_i\langle\mathbf{x}_i,\boldsymbol{\delta}_0\rangle + \frac{2}{n}\mathbb{E}\sup_{\boldsymbol{\delta}_{[G]\backslash}\in\mathcal{C}_{[G]\backslash}}\sum_{g=1}^{G}\sum_{i=1}^{n_g}\epsilon_{gi}\langle\mathbf{x}_{gi},\boldsymbol{\delta}_g\rangle$$

$$= \frac{2}{\sqrt{n}}\mathbb{E}\sup_{\boldsymbol{\delta}_0\in\mathcal{C}_0}\sum_{i=1}^{n}\langle\frac{1}{\sqrt{n}}\epsilon_i\mathbf{x}_i,\boldsymbol{\delta}_0\rangle + \frac{2}{\sqrt{n}}\mathbb{E}\sup_{\boldsymbol{\delta}_{[G]\backslash}\in\mathcal{C}_{[G]\backslash}}\sum_{g=1}^{G}\sqrt{\frac{n_g}{n}}\sum_{i=1}^{n_g}\langle\frac{1}{\sqrt{n_g}}\epsilon_{gi}\mathbf{x}_{gi},\boldsymbol{\delta}_g\rangle$$

$$(n_0 := n, \epsilon_{0i} := \epsilon_0, \mathbf{x}_{0i} := \mathbf{x}_i) = \frac{2}{\sqrt{n}}\mathbb{E}\sup_{\boldsymbol{\delta}_{[G]}\in\mathcal{C}_{[G]}}\sum_{g=0}^{G}\sqrt{\frac{n_g}{n}}\sum_{i=1}^{n_g}\langle\frac{1}{\sqrt{n_g}}\epsilon_{gi}\mathbf{x}_{gi},\boldsymbol{\delta}_g\rangle$$

$$(\mathbf{h}_g := \frac{1}{\sqrt{n_g}}\sum_{i=1}^{n_g}\epsilon_{gi}\mathbf{x}_{gi}) = \frac{2}{\sqrt{n}}\mathbb{E}\sup_{\boldsymbol{\delta}_{[G]}\in\mathcal{C}_{[G]}}\sum_{g=0}^{G}\sqrt{\frac{n_g}{n}}\langle\mathbf{h}_g,\boldsymbol{\delta}_g\rangle$$

$$(\mathcal{A}_g\in\mathcal{C}_g\cap\mathbb{S}^{p-1}) \le \frac{2}{\sqrt{n}}\mathbb{E}\sup_{\boldsymbol{\delta}_{[G]}\in\mathcal{A}_{[G]}}\sum_{g=0}^{G}\sqrt{\frac{n_g}{n}}\langle\mathbf{h}_g,\boldsymbol{\delta}_g\rangle\|\boldsymbol{\delta}_g\|_2$$

$$\le \frac{2}{\sqrt{n}}\sum_{g=0}^{G}\sqrt{\frac{n_g}{n}}\mathbb{E}_{\mathbf{h}_g}\sup_{\boldsymbol{\delta}_g\in\mathcal{A}_g}\langle\mathbf{h}_g,\boldsymbol{\delta}_g\rangle\|\boldsymbol{\delta}_g\|_2$$

$$\le \frac{2}{\sqrt{n}}\sum_{g=0}^{G}\sqrt{\frac{n_g}{n}}c_g k\omega(\mathcal{A}_g)\|\boldsymbol{\delta}_g\|_2$$

Note that the $\mathbf{h}_{gi}$ is a sub-Gaussian random vector which let us bound the $\mathbb{E}\sup$ using the Gaussian width [38] in the last step. ∎

### B.3. Proof of Lemma A.3.

*Proof.* To avoid cluttering let $h_g(\boldsymbol{\omega}_g,\mathbf{X}_g) = \sqrt{\frac{n}{n_g}}\|\boldsymbol{\omega}_g\|_2\sup_{\mathbf{u}_g\in\mathcal{A}_g}\langle\mathbf{X}_g^T\frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2},\mathbf{u}_g\rangle$, $e_g = \zeta_g k\omega(\mathcal{A}_g)+$
$\epsilon_g\sqrt{\log G} + \tau$, where $s_g = \sqrt{\frac{n}{n_g}}\sqrt{(2K^2+1)n_g}$.

$$\mathbb{P}(h_g(\boldsymbol{\omega}_g,\mathbf{X}_g) > e_g s_g) = \mathbb{P}\left(h_g(\boldsymbol{\omega}_g,\mathbf{X}_g) > e_g s_g\Big|\sqrt{\frac{n}{n_g}}\|\boldsymbol{\omega}_g\|_2 > s_g\right)\mathbb{P}\left(\sqrt{\frac{n}{n_g}}\|\boldsymbol{\omega}_g\|_2 > s_g\right)$$

$$+ \mathbb{P}\left(h_g(\boldsymbol{\omega}_g,\mathbf{X}_g) > e_g s_g\Big|\sqrt{\frac{n}{n_g}}\|\boldsymbol{\omega}_g\|_2 < s_g\right)\mathbb{P}\left(\sqrt{\frac{n}{n_g}}\|\boldsymbol{\omega}_g\|_2 < s_g\right)$$

$$\le \mathbb{P}\left(\sqrt{\frac{n}{n_g}}\|\boldsymbol{\omega}_g\|_2 > s_g\right) + \mathbb{P}\left(h_g(\boldsymbol{\omega}_g,\mathbf{X}_g) > e_g s_g\Big|\sqrt{\frac{n}{n_g}}\|\boldsymbol{\omega}_g\|_2 < s_g\right)$$

$$\le \mathbb{P}\left(\|\boldsymbol{\omega}_g\|_2 > \sqrt{(2K^2+1)n_g}\right) + \mathbb{P}\left(\sup_{\mathbf{u}_g\in\mathcal{C}_g\cap\mathbb{S}^{p-1}}\langle\mathbf{X}_g^T\frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2},\mathbf{u}_g\rangle > e_g\right)$$

$$\le \mathbb{P}\left(\|\boldsymbol{\omega}_g\|_2 > \sqrt{(2K^2+1)n_g}\right) + \sup_{\mathbf{v}\in\mathbb{S}^{p-1}}\mathbb{P}\left(\sup_{\mathbf{u}_g\in\mathcal{C}_g\cap\mathbb{S}^{p-1}}\langle\mathbf{X}_g^T\mathbf{v},\mathbf{u}_g\rangle > e_g\right)$$

Let's focus on the first term. Since $\boldsymbol{\omega}_g$ consists of i.i.d. centered unit-variance sub-Gaussian elements with $\|\|\omega_{gi}\|\|_{\psi_2} < K$, $\omega_{gi}^2$ is sub-exponential with $\|\|\omega_{gi}\|\|_{\psi_1} < 2K^2$. Let's apply the Bernstein's inequality

to $\|\boldsymbol{\omega}_g\|_2^2 = \sum_{i=1}^{n_g} \omega_{gi}^2$:

$$\mathbb{P}\left(\left|\|\boldsymbol{\omega}_g\|_2^2 - \mathbb{E}\|\boldsymbol{\omega}_g\|_2^2\right| > \tau\right) \leq 2\exp\left(-\nu_g \min\left[\frac{\tau^2}{4K^4 n_g}, \frac{\tau}{2K^2}\right]\right)$$

We also know that $\mathbb{E}\|\omega_g\|_2^2 \leq n_g$ [2] which gives us:

$$\mathbb{P}\left(\|\boldsymbol{\omega}_g\|_2 > \sqrt{n_g + \tau}\right) \leq 2\exp\left(-\nu_g \min\left[\frac{\tau^2}{4K^4 n_g}, \frac{\tau}{2K^2}\right]\right)$$

Finally, we set $\tau = 2K^2 n_g$:

$$\mathbb{P}\left(\|\boldsymbol{\omega}_g\|_2 > \sqrt{(2K^2+1)n_g}\right) \leq 2\exp\left(-\nu_g n_g\right) = \frac{2}{(G+1)}\exp\left(-\nu_g n_g + \log(G+1)\right)$$

Now we upper bound the second term of (B.3). Given any fixed $\mathbf{v} \in \mathbb{S}^{p-1}$, $\mathbf{X}_g \mathbf{v}$ is a sub-Gaussian random vector with $\left\|\left|\mathbf{X}_g^T \mathbf{v}\right|\right\|_{\psi_2} \leq C_g k$ [2]. From Theorem 9 of [2] for any $\mathbf{v} \in \mathbb{S}^{p-1}$ we have:

$$\mathbb{P}\left(\sup_{\mathbf{u}_g \in \mathcal{A}_g} \langle \mathbf{X}_g^T \mathbf{v}, \mathbf{u}_g \rangle > v_g C_g k \omega(\mathcal{A}_g) + t\right) \leq \pi_g \exp\left(-\left(\frac{t}{\theta_g C_g k \phi_g}\right)^2\right)$$

where $\phi_g = \sup_{\mathbf{u}_g \in \mathcal{A}_g} \|\mathbf{u}_g\|_2$ and in our problem $\phi_g = 1$. We now substitute $t = \tau + \epsilon_g \sqrt{\log(G+1)}$ where $\epsilon_g = \theta_g C_g k$.

$$\mathbb{P}\left(\sup_{\mathbf{u}_g \in \mathcal{A}_g} \langle \mathbf{X}_g^T \mathbf{v}, \mathbf{u}_g \rangle > v_g C_g k \omega(\mathcal{A}_g) + \epsilon_g \sqrt{\log(G+1)} + \tau\right) \leq \pi_g \exp\left(-\left(\frac{\tau + \epsilon_g \sqrt{\log(G+1)}}{\epsilon_g}\right)^2\right)$$

$$\leq \pi_g \exp\left(-\log G - \left(\frac{\tau}{\theta_g C_g k}\right)^2\right)$$

$$\leq \frac{\pi_g}{(G+1)} \exp\left(-\left(\frac{\tau}{\theta_g C_g k}\right)^2\right)$$

Now we put back results to the original inequality (B.3):

$$\mathbb{P}\left(h_g(\boldsymbol{\omega}_g, \mathbf{X}_g) > \sqrt{\frac{n}{n_g}}\sqrt{(2K^2+1)n_g} \times \left(v_g C_g k \omega(\mathcal{A}_g) + \epsilon_g \sqrt{\log(G+1)} + \tau\right)\right)$$

$$\leq \frac{\sigma_g}{(G+1)} \exp\left(-\min\left[\nu_g n_g - \log(G+1), \frac{\tau^2}{\theta_g^2 C_g^2 k^2}\right]\right)$$

$$\leq \frac{\sigma_g}{(G+1)} \exp\left(-\min\left[\nu_g n_g - \log(G+1), \frac{\tau^2}{\eta_g^2 k^2}\right]\right)$$

where $\sigma_g = \pi_g + 2$, $\zeta_g = v_g C_g$, $\eta_g = \theta_g C_g$. ∎

### B.4. Proof of Lemma A.4.

*Proof.* We upper bound the individual error $\|\boldsymbol{\delta}_g^{(t+1)}\|_2$ and the common one $\|\boldsymbol{\delta}_0^{(t+1)}\|_2$ in the followings:

$$\|\boldsymbol{\delta}_g^{(t+1)}\|_2 = \|\boldsymbol{\beta}_g^{(t+1)} - \boldsymbol{\beta}_g^*\|_2$$

$$= \left\| \Pi_{\Omega_{f_g}}\left( \boldsymbol{\beta}_g^{(t)} + \mu_g \mathbf{X}_g^T\left( \mathbf{y}_g - \mathbf{X}_g(\boldsymbol{\beta}_0^{(t)} + \boldsymbol{\beta}_g^{(t)}) \right) \right) - \boldsymbol{\beta}_g^* \right\|_2$$

$$(\text{Lemma 6.3 of [32]}) = \left\| \Pi_{\Omega_{f_g} - \{\boldsymbol{\beta}_g^*\}}\left( \boldsymbol{\beta}_g^{(t)} + \mu_g \mathbf{X}_g^T\left( \mathbf{y}_g - \mathbf{X}_g(\boldsymbol{\beta}_0^{(t)} + \boldsymbol{\beta}_g^{(t)}) \right) - \boldsymbol{\beta}_g^* \right) \right\|_2$$

$$= \left\| \Pi_{\mathcal{E}_g}\left( \boldsymbol{\delta}_g^{(t)} + \mu_g \mathbf{X}_g^T\left( \mathbf{y}_g - \mathbf{X}_g(\boldsymbol{\beta}_0^{(t)} + \boldsymbol{\beta}_g^{(t)}) - \mathbf{X}_g(\boldsymbol{\beta}_0^* + \boldsymbol{\beta}_g^*) + \mathbf{X}_g(\boldsymbol{\beta}_0^* + \boldsymbol{\beta}_g^*) \right) \right) \right\|_2$$

$$= \left\| \Pi_{\mathcal{E}_g}\left( \boldsymbol{\delta}_g^{(t)} + \mu_g \mathbf{X}_g^T\left( \boldsymbol{\omega}_g - \mathbf{X}_g(\boldsymbol{\delta}_0^{(t)} + \boldsymbol{\delta}_g^{(t)}) \right) \right) \right\|_2$$

$$(\text{Lemma 6.4 of [32]}) \leq \left\| \Pi_{\mathcal{C}_g}\left( \boldsymbol{\delta}_g^{(t)} + \mu_g \mathbf{X}_g^T\left( \boldsymbol{\omega}_g - \mathbf{X}_g(\boldsymbol{\delta}_0^{(t)} + \boldsymbol{\delta}_g^{(t)}) \right) \right) \right\|_2$$

$$(\text{Lemma 6.2 of [32]}) \leq \sup_{\mathbf{v} \in \mathcal{C}_g \cap \mathbb{B}^p} \mathbf{v}^T\left( \boldsymbol{\delta}_g^{(t)} + \mu_g \mathbf{X}_g^T\left( \boldsymbol{\omega}_g - \mathbf{X}_g(\boldsymbol{\delta}_0^{(t)} + \boldsymbol{\delta}_g^{(t)}) \right) \right)$$

$$(\mathcal{B}_g = \mathcal{C}_g \cap \mathbb{B}^p) = \sup_{\mathbf{v} \in \mathcal{B}_g} \mathbf{v}^T\left( \boldsymbol{\delta}_g^{(t)} + \mu_g \mathbf{X}_g^T\left( \boldsymbol{\omega}_g - \mathbf{X}_g(\boldsymbol{\delta}_0^{(t)} + \boldsymbol{\delta}_g^{(t)}) \right) \right)$$

$$\leq \sup_{\mathbf{v} \in \mathcal{B}_g} \mathbf{v}^T\left( \mathbf{I}_g - \mu_g \mathbf{X}_g^T \mathbf{X}_g \right)\boldsymbol{\delta}_g^{(t)} + \mu_g \sup_{\mathbf{v} \in \mathcal{B}_g} \mathbf{v}^T \mathbf{X}_g^T \boldsymbol{\omega}_g + \mu_g \sup_{\mathbf{v} \in \mathcal{B}_g} -\mathbf{v}^T \mathbf{X}_g^T \mathbf{X}_g \boldsymbol{\delta}_0^{(t)}$$

$$\leq \left\| \boldsymbol{\delta}_g^{(t)} \right\|_2 \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{B}_g} \mathbf{v}^T\left( \mathbf{I}_g - \mu_g \mathbf{X}_g^T \mathbf{X}_g \right)\mathbf{u} + \mu_g \|\boldsymbol{\omega}_g\|_2 \sup_{\mathbf{v} \in \mathcal{B}_g} \mathbf{v}^T \mathbf{X}_g^T \frac{\boldsymbol{\omega}_g}{\|\boldsymbol{\omega}_g\|_2}$$

$$+ \mu_g \|\boldsymbol{\delta}_0^{(t)}\|_2 \sup_{\mathbf{v} \in \mathcal{B}_g, \mathbf{u} \in \mathcal{B}_0} -\mathbf{v}^T \mathbf{X}_g^T \mathbf{X}_g \mathbf{u}$$

$$= \rho_g(\mu_g)\|\boldsymbol{\delta}_g^{(t)}\|_2 + \xi_g(\mu_g)\|\boldsymbol{\omega}_g\|_2 + \phi_g(\mu_g)\|\boldsymbol{\delta}_0^{(t)}\|_2$$

So the final bound becomes:

$$(\text{B.4}) \qquad \|\boldsymbol{\delta}_g^{(t+1)}\|_2 \leq \rho_g(\mu_g)\|\boldsymbol{\delta}_g^{(t)}\|_2 + \xi_g(\mu_g)\|\boldsymbol{\omega}_g\|_2 + \phi_g(\mu_g)\|\boldsymbol{\delta}_0^{(t)}\|_2$$

∎

Now we upper bound the error of common parameter. Remember common parameter's update:

681 $$\boldsymbol{\beta}_0^{(t+1)} = \Pi_{\Omega_{f_0}} \left( \boldsymbol{\beta}_0^{(t)} + \mu_0 \mathbf{X}_0^T \begin{pmatrix} (\mathbf{y}_1 - \mathbf{X}_1(\boldsymbol{\beta}_0^{(t)} + \boldsymbol{\beta}_1^{(t)})) \\ \vdots \\ (\mathbf{y}_G - \mathbf{X}_G(\boldsymbol{\beta}_0^{(t)} + \boldsymbol{\beta}_G^{(t)})) \end{pmatrix} \right).$$

682 $$\|\boldsymbol{\delta}_0^{(t+1)}\|_2 = \|\boldsymbol{\beta}_0^{(t+1)} - \boldsymbol{\beta}_0^*\|_2$$

683

684 $$= \left\| \Pi_{\Omega_{f_0}} \left( \boldsymbol{\beta}_0^{(t)} + \mu_0 \sum_{g=1}^{G} \mathbf{X}_g^T \left( \mathbf{y}_g - \mathbf{X}_g(\boldsymbol{\beta}_0^{(t)} + \boldsymbol{\beta}_g^{(t)}) \right) \right) - \boldsymbol{\beta}_0^* \right\|_2$$

685 (Lemma 6.3 of [32]) $$= \left\| \Pi_{\Omega_{f_0} - \{\boldsymbol{\beta}_0^*\}} \left( \boldsymbol{\beta}_0^{(t)} + \mu_0 \sum_{g=1}^{G} \mathbf{X}_g^T \left( \mathbf{y}_g - \mathbf{X}_g(\boldsymbol{\beta}_0^{(t)} + \boldsymbol{\beta}_g^{(t)}) \right) \right) - \boldsymbol{\beta}_0^* \right\|_2$$

686 $$= \left\| \Pi_{\mathcal{E}_0} \left( \boldsymbol{\delta}_0^{(t)} + \mu_0 \sum_{g=1}^{G} \mathbf{X}_g^T \left( \mathbf{y}_g - \mathbf{X}_g(\boldsymbol{\beta}_0^{(t)} + \boldsymbol{\beta}_g^{(t)}) \right) \right) \right\|_2$$

687 (Lemma 6.4 of [32]) $$\leq \left\| \Pi_{\mathcal{C}_0} \left( \boldsymbol{\delta}_0^{(t)} + \mu_0 \sum_{g=1}^{G} \mathbf{X}_g^T \left( \boldsymbol{\omega}_g - \mathbf{X}_g(\boldsymbol{\delta}_0^{(t)} + \boldsymbol{\delta}_g^{(t)}) \right) \right) \right\|_2$$

688 (Lemma 6.2 of [32]) $$\leq \sup_{\mathbf{v} \in \mathcal{B}_0} \mathbf{v}^T \left( \boldsymbol{\delta}_0^{(t)} + \mu_0 \sum_{g=1}^{G} \mathbf{X}_g^T \left( \boldsymbol{\omega}_g - \mathbf{X}_g(\boldsymbol{\delta}_0^{(t)} + \boldsymbol{\delta}_g^{(t)}) \right) \right)$$

689 $$\leq \sup_{\mathbf{v} \in \mathcal{B}_0} \mathbf{v}^T \left( \mathbf{I} - \mu_0 \sum_{g=1}^{G} \mathbf{X}_g^T \mathbf{X}_g \right) \boldsymbol{\delta}_0^{(t)} + \mu_0 \sup_{\mathbf{v} \in \mathcal{B}_0} \mathbf{v}^T \sum_{g=1}^{G} \mathbf{X}_g^T \boldsymbol{\omega}_g$$

690 $$+ \mu_0 \sup_{\mathbf{v} \in \mathcal{B}_0} -\mathbf{v}^T \sum_{g=1}^{G} \mathbf{X}_g^T \mathbf{X}_g \boldsymbol{\delta}_g^{(t)}$$

691 $$\leq \|\boldsymbol{\delta}_0^{(t)}\|_2 \sup_{\mathbf{u},\mathbf{v} \in \mathcal{B}_0} \mathbf{v}^T \left( \mathbf{I} - \mu_0 \mathbf{X}_0^T \mathbf{X}_0 \right) \mathbf{u} + \mu_0 \sup_{\mathbf{v} \in \mathcal{B}_0} \mathbf{v}^T \mathbf{X}_0^T \frac{\boldsymbol{\omega}_0}{\|\boldsymbol{\omega}_0\|_2} \|\boldsymbol{\omega}_0\|_2$$

692 $$+ \mu_0 \sum_{g=1}^{G} \sup_{\mathbf{v}_g \in \mathcal{B}_0, \mathbf{u}_g \in \mathcal{B}_g} -\mathbf{v}_g^T \mathbf{X}_g^T \mathbf{X}_g \mathbf{u}_g \|\boldsymbol{\delta}_g^{(t)}\|_2$$

693 (B.5) $$\leq \rho_0(\mu_0) \|\boldsymbol{\delta}_0^{(t)}\|_2 + \xi_0(\mu_0) \|\boldsymbol{\omega}_0\|_2 + \mu_0 \sum_{g=1}^{G} \frac{\phi_g(\mu_g)}{\mu_g} \|\boldsymbol{\delta}_g^{(t)}\|_2$$

694

695 To avoid cluttering we drop $\mu_g$ as the arguments. Putting together (B.4) and (B.5) inequalities we
696 reach to the followings:

697 $$\|\boldsymbol{\delta}_g^{(t+1)}\|_2 \leq \rho_g \|\boldsymbol{\delta}_g^{(t)}\|_2 + \xi_g \|\boldsymbol{\omega}_g\|_2 + \phi_g \|\boldsymbol{\delta}_0^{(t)}\|_2$$

698 $$\|\boldsymbol{\delta}_0^{(t+1)}\|_2 \leq \rho_0 \|\boldsymbol{\delta}_0^{(t)}\|_2 + \xi_0 \|\boldsymbol{\omega}_0\|_2 + \mu_0 \sum_{g=1}^{G} \frac{\phi_g}{\mu_g} \|\boldsymbol{\delta}_g^{(t)}\|_2$$

**B.5. Proof of Lemma A.5.** We will need the following lemma in our proof. It establishes the RE condition for individual isotropic sub-Gaussian designs and provides us with the essential tool for proving high probability bounds.

**Lemma B.3 (Theorem 11 of [2]).** *For all $g \in [G]$, for the matrix $\mathbf{X}_g \in \mathbb{R}^{n_g \times p}$ with independent isotropic sub-Gaussian rows, i.e., $\|\|\mathbf{x}_{gi}\|\|_{\psi_2} \leq k$ and $\mathbb{E}[\mathbf{x}_{gi}\mathbf{x}_{gi}^T] = \mathbf{I}$, the following result holds with probability at least $1 - 2\exp\left(-\gamma_g(\omega(\mathcal{A}_g) + \tau)^2\right)$ for $\tau > 0$:*

$$\forall \mathbf{u}_g \in \mathcal{C}_g : n_g\left(1 - c_g\frac{\omega(\mathcal{A}_g) + \tau}{\sqrt{n_g}}\right)\|\mathbf{u}_g\|_2^2 \leq \|\mathbf{X}_g\mathbf{u}_g\|_2^2 \leq n_g\left(1 + c_g\frac{\omega(\mathcal{A}_g) + \tau}{\sqrt{n_g}}\right)\|\mathbf{u}_g\|_2^2$$

*where $c_g > 0$ is constant.*

The statement of Lemma B.3 characterizes the distortion in the Euclidean distance between points $\mathbf{u}_g \in \mathcal{C}_g$ when the matrix $\mathbf{X}_g/n_g$ is applied to them and states that any sub-Gaussian design matrix is approximately isometry, with high probability:

$$(1 - \alpha)\|\mathbf{u}_g\|_2^2 \leq \frac{1}{n_g}\|\mathbf{X}_g\mathbf{u}_g\|_2^2 \leq (1 + \alpha)\|\mathbf{u}_g\|_2^2$$

where $\alpha = c_g\frac{\omega(\mathcal{A}_g)}{\sqrt{n_g}}$.

Now the proof for Lemma A.5:

*Proof.* First we upper bound each of the coefficients $\forall g \in [G]$:

$$\rho_g(\mu_g) = \sup_{\mathbf{u},\mathbf{v}\in\mathcal{B}_g} \mathbf{v}^T\left(\mathbf{I}_g - \mu_g\mathbf{X}_g^T\mathbf{X}_g\right)\mathbf{u}$$

We upper bound the argument of the sup as follows:

$$\mathbf{v}^T\left(\mathbf{I}_g - \mu_g\mathbf{X}_g^T\mathbf{X}_g\right)\mathbf{u} = \frac{1}{4}\left[(\mathbf{u}+\mathbf{v})^T(\mathbf{I} - \mu_g\mathbf{X}_g^T\mathbf{X}_g)(\mathbf{u}+\mathbf{v}) - (\mathbf{u}-\mathbf{v})^T(\mathbf{I} - \mu_g\mathbf{X}_g^T\mathbf{X}_g)(\mathbf{u}-\mathbf{v})\right]$$

$$= \frac{1}{4}\left[\|\mathbf{u}+\mathbf{v}\|_2^2 - \mu_g\|\mathbf{X}_g(\mathbf{u}+\mathbf{v})\|_2^2 - \|\mathbf{u}-\mathbf{v}\|_2^2 + \mu_g\|\mathbf{X}_g(\mathbf{u}-\mathbf{v})\|_2^2\right]$$

$$(\text{Lemma B.3}) \leq \frac{1}{4}\left[\left(1 - \mu_g n_g\left(1 - c_g\frac{2\omega(\mathcal{A}_g) + \tau}{\sqrt{n_g}}\right)\right)\|\mathbf{u}+\mathbf{v}\|_2\right.$$

$$\left. - \left(1 - \mu_g n_g\left(1 + c_g\frac{2\omega(\mathcal{A}_g) + \tau}{\sqrt{n_g}}\right)\right)\|\mathbf{u}-\mathbf{v}\|_2\right]$$

$$\left(\mu_g = \frac{1}{a_g n_g}\right) \leq \frac{1}{4}\left[\left(1 - \frac{1}{a_g}\right)(\|\mathbf{u}+\mathbf{v}\|_2 - \|\mathbf{u}-\mathbf{v}\|_2) + c_g\frac{2\omega(\mathcal{A}_g) + \tau}{a_g\sqrt{n_g}}(\|\mathbf{u}+\mathbf{v}\|_2 + \|\mathbf{u}-\mathbf{v}\|_2)\right]$$

$$\leq \frac{1}{4}\left[\left(1 - \frac{1}{a_g}\right)2\|\mathbf{v}\|_2 + c_g\frac{2\omega(\mathcal{A}_g) + \tau}{a_g\sqrt{n_g}}2\sqrt{2}\right]$$

where the last line follows from the triangle inequality and the fact that $\|\mathbf{u}+\mathbf{v}\|_2 + \|\mathbf{u}-\mathbf{v}\|_2 \leq 2\sqrt{2}$ which itself follows from $\|\mathbf{u}+\mathbf{v}\|_2^2 + \|\mathbf{u}-\mathbf{v}\|_2^2 \leq 4$. Note that we applied the Lemma B.3 for bigger

726 sets of $\mathcal{A}_g + \mathcal{A}_g$ and $\mathcal{A}_g - \mathcal{A}_g$ where Gaussian width of both of them are upper bounded by $2\omega(\mathcal{A}_g)$.
727 The above holds with high probability (computed below). Now we set :

728 (B.6)
$$\mathbf{v}^T\left(\mathbf{I}_g - \frac{1}{a_g n_g}\mathbf{X}_g^T\mathbf{X}_g\right)\mathbf{u} \leq \frac{1}{2}\left[\left(1 - \frac{1}{a_g}\right) + \sqrt{2}c_g\frac{2\omega(\mathcal{A}_g) + \tau}{a_g\sqrt{n_g}}\right]$$

729 To keep the upper bound of $\rho_g$ in (B.6) below any arbitrary $\frac{1}{b} < 1$ we need $n_g = O(b^2(\omega(\mathcal{A}_g)+\tau)^2)$
730 samples.

731 Now we rewrite the same analysis using the tail bounds for the coefficients to clarify the probabilities.
732 Let's set $\mu_g = \frac{1}{a_g n_g}$, $d_g := \frac{1}{2}\left(1 - \frac{1}{a_g}\right) + \sqrt{2}c_g\frac{\omega(\mathcal{A}_g)+\tau/2}{a_g\sqrt{n_g}}$ and name the bad events of $\|\mathbf{X}_g(\mathbf{u}+\mathbf{v})\|_2^2 <$
733 $n_g\left(1 - c_g\frac{2\omega(\mathcal{A}_g)+\tau}{\sqrt{n_g}}\right)$ and $\|\mathbf{X}_g(\mathbf{u}-\mathbf{v})\|_2^2 > n_g\left(1 + c_g\frac{2\omega(\mathcal{A}_g)+\tau}{\sqrt{n_g}}\right)$ as $\mathcal{E}_1$ and $\mathcal{E}_2$ respectively:

734
$$\mathbb{P}(\rho_g \geq d_g) \leq \mathbb{P}(\rho_g \geq d_g | \neg\mathcal{E}_1, \neg\mathcal{E}_2) + 2\mathbb{P}(\mathcal{E}_1) + \mathbb{P}(\mathcal{E}_2)$$

735
$$\text{Lemma B.3} \leq 0 + 6\exp\left(-\gamma_g(\omega(\mathcal{A}_g) + \tau)^2\right)$$

736 which concludes the proof. ∎

### B.6. Proof of Lemma A.6.

738 *Proof.* The following holds for any $\mathbf{u}$ and $\mathbf{v}$ because of $\|\mathbf{X}_g(\mathbf{u}+\mathbf{v})\|_2^2 \geq 0$:

739 (B.7)
$$-\mathbf{v}^T\mathbf{X}_g^T\mathbf{X}_g\mathbf{u} \leq \frac{1}{2}\left(\|\mathbf{X}_g\mathbf{u}\|_2^2 + \|\mathbf{X}_g\mathbf{v}\|_2^2\right)$$

740 Now we can bound $\phi_g$ as follows:

741 (B.8)
$$\phi_g(\mu_g) = \mu_g\sup_{\mathbf{v}\in\mathcal{B}_g,\mathbf{u}\in\mathcal{B}_0} -\mathbf{v}^T\mathbf{X}_g^T\mathbf{X}_g\mathbf{u} \leq \frac{\mu_g}{2}\left(\sup_{\mathbf{u}\in\mathcal{B}_0}\|\mathbf{X}_g\mathbf{u}\|_2^2 + \sup_{\mathbf{v}\in\mathcal{B}_g}\|\mathbf{X}_g\mathbf{v}\|_2^2\right)$$

■

742 So we have:

743 (B.9)
$$\phi_g\left(\frac{1}{a_g n_g}\right) \leq \frac{1}{2a_g}\left(\frac{1}{n_g}\sup_{\mathbf{u}\in\mathcal{B}_0}\|\mathbf{X}_g\mathbf{u}\|_2^2 + \frac{1}{n_g}\sup_{\mathbf{v}\in\mathcal{B}_g}\|\mathbf{X}_g\mathbf{v}\|_2^2\right)$$

744
$$\text{(Lemma B.3)} \leq \frac{1}{a_g}\left(1 + c_{0g}\frac{\omega(\mathcal{A}_g) + \omega(\mathcal{A}_0) + 2\tau}{2\sqrt{n_g}}\right)$$

745
$$(\omega_{0g} = \max(\omega(\mathcal{A}_0), \omega(\mathcal{A}_g)) \leq \frac{1}{a_g}\left(1 + c_{0g}\frac{\omega_{0g} + \tau}{\sqrt{n_g}}\right)$$

746 where $c_{0g} = \max(c_0, c_g)$.

747 To compute the exact probabilities lets define $s_g := \frac{1}{a_g}\left(1 + c_{0g}\frac{\omega(\mathcal{A}_g)+\omega(\mathcal{A}_0)+2\tau}{2\sqrt{n_g}}\right)$ and name the
748 bad events of $\frac{1}{n_g}\sup_{\mathbf{u}\in\mathcal{B}_0}\|\mathbf{X}_g\mathbf{u}\|_2^2 > 1 + c_0\frac{\omega(\mathcal{A}_0)+\tau}{\sqrt{n_g}}$ and $\frac{1}{n_g}\sup_{\mathbf{v}\in\mathcal{B}_g}\|\mathbf{X}_g\mathbf{v}\|_2^2 > 1 + c_g\frac{\omega(\mathcal{A}_g)+\tau}{\sqrt{n_g}}$ as $\mathcal{E}_1$
749 and $\mathcal{E}_2$ respectively.

750 (B.10)
$$\mathbb{P}(\phi_g > s_g) \leq \mathbb{P}(\phi_g > s_g | \neg\mathcal{E}_1)\mathbb{P}(\neg\mathcal{E}_1) + \mathbb{P}(\mathcal{E}_1)$$

751
$$\leq \mathbb{P}(\mathcal{E}_2) + \mathbb{P}(\mathcal{E}_1)$$

752
$$\leq 4\exp\left(-\gamma_g(\omega(\mathcal{A}_g) + \tau)^2\right)$$

### B.7. Proof of Lemma B.1.

*Proof.* To obtain lower bound, we use the Paley–Zygmund inequality for the zero-mean, non-degenerate ($0 < \alpha \leq \mathbb{E}|\langle \mathbf{x}, \mathbf{u} \rangle|, \mathbf{u} \in \mathbb{S}^{p-1}$) sub-Gaussian random vector $\mathbf{x}$ with $\|\|\mathbf{x}\|\|_{\psi_2} \leq k$ [38].

$$Q_{2\xi}(\mathbf{u}) \geq \frac{(\alpha - 2\xi)^2}{4ck^2}.$$

### B.8. Proof of Lemma B.2.

*Proof.* We split $[G]_{\setminus} - \mathcal{I}$ into two groups $\mathcal{J}, \mathcal{K}$. $\mathcal{J}$ consists of $\boldsymbol{\delta}_i$'s with $\|\boldsymbol{\delta}_i\|_2 \geq 2\|\boldsymbol{\delta}_0\|_2$ and $\mathcal{K} = [G]_{\setminus} - \mathcal{I} - \mathcal{J}$. We use the bounds

(B.11)
$$\|\boldsymbol{\delta}_0 + \boldsymbol{\delta}_i\|_2 \geq \begin{cases} \lambda_{\min}(\|\boldsymbol{\delta}_i\|_2 + \|\boldsymbol{\delta}_0\|_2) & \text{if } i \in \mathcal{I} \\ \|\boldsymbol{\delta}_i\|_2/2 & \text{if } i \in \mathcal{J} \\ 0 & \text{if } i \in \mathcal{K} \end{cases}$$

This implies

$$\sum_{i=1}^{G} n_i\|\boldsymbol{\delta}_0 + \boldsymbol{\delta}_i\|_2 \geq \sum_{i \in \mathcal{J}} \frac{n_i}{2}\|\boldsymbol{\delta}_i\|_2 + \lambda_{\min}\sum_{i \in \mathcal{I}} n_i(\|\boldsymbol{\delta}_i\|_2 + \|\boldsymbol{\delta}_0\|_2). \qquad \blacksquare$$

Let $S_{\mathcal{S}} = \sum_{i \in \mathcal{S}} n_i\|\boldsymbol{\delta}_i\|_2$ for $\mathcal{S} = \mathcal{I}, \mathcal{J}, \mathcal{K}$. We know that over $\mathcal{K}$, $\|\boldsymbol{\delta}_i\|_2 \leq 2\|\boldsymbol{\delta}_0\|_2$ which implies $S_{\mathcal{K}} = \sum_{i \in \mathcal{K}} n_i\|\boldsymbol{\delta}_i\|_2 \leq 2\sum_{i \in \mathcal{K}} n_i\|\boldsymbol{\delta}_0\|_2 \leq 2n\|\boldsymbol{\delta}_0\|_2$. Set $\psi_{\mathcal{I}} = \min\{1/2, \lambda_{\min}\bar{\rho}/3\} = \lambda_{\min}\bar{\rho}/3$. Using $1/2 \geq \psi_{\mathcal{I}}$, we write:

$$\sum_{i=1}^{G} n_i\|\boldsymbol{\delta}_0 + \boldsymbol{\delta}_i\|_2 \geq \psi_{\mathcal{I}}S_{\mathcal{J}} + \lambda_{\min}\sum_{i \in \mathcal{I}} n_i(\|\boldsymbol{\delta}_i\|_2 + \|\boldsymbol{\delta}_0\|_2)$$

$$(S_{\mathcal{K}} \leq 2n\|\boldsymbol{\delta}_0\|_2) \geq \psi_{\mathcal{I}}S_{\mathcal{J}} + \psi_{\mathcal{I}}S_{\mathcal{K}} - 2\psi_{\mathcal{I}}n\|\boldsymbol{\delta}_0\|_2 + \left(\sum_{i \in \mathcal{I}} n_i\right)\lambda_{\min}\|\boldsymbol{\delta}_0\|_2 + \lambda_{\min}S_{\mathcal{I}}$$

$$(\lambda_{\min} \geq \psi_{\mathcal{I}}) \geq \psi_{\mathcal{I}}(S_{\mathcal{I}} + S_{\mathcal{J}} + S_{\mathcal{K}}) + \left(\left(\sum_{i \in \mathcal{I}} n_i\right)\lambda_{\min} - 2\psi_{\mathcal{I}}n\right)\|\boldsymbol{\delta}_0\|_2.$$

Now, observe that, assumption of the Definition 3.5, $\sum_{i \in \mathcal{I}} n_i \geq \bar{\rho}n$ implies:

$$\left(\sum_{i \in \mathcal{I}} n_i\right)\lambda_{\min} - 2\psi_{\mathcal{I}}n \geq (\bar{\rho}\lambda_{\min} - 2\psi_{\mathcal{I}})n \geq \psi_{\mathcal{I}}n.$$

Combining all, we obtain:

$$\sum_{i=1}^{G} n_i\|\boldsymbol{\delta}_0 + \boldsymbol{\delta}_i\|_2 \geq \psi_{\mathcal{I}}(S_{\mathcal{I}} + S_{\mathcal{J}} + S_{\mathcal{K}} + \|\boldsymbol{\delta}_0\|_2) = \psi_{\mathcal{I}}(n\|\boldsymbol{\delta}_0\|_2 + \sum_{i=1}^{G} n_i\|\boldsymbol{\delta}_i\|_2).$$

## REFERENCES

[1] F. BACH, R. JENATTON, J. MAIRAL, G. OBOZINSKI, ET AL., *Optimization with sparsity-inducing penalties*, Foundations and Trends® in Machine Learning, 4 (2012), pp. 1–106.

[2] A. BANERJEE, S. CHEN, F. FAZAYELI, AND V. SIVAKUMAR, *Estimation with Norm Regularization*, in Advances in Neural Information Processing Systems, 2014, pp. 1556–1564.

[3] J. BARRETINA, G. CAPONIGRO, N. STRANSKY, K. VENKATESAN, A. A. MARGOLIN, S. KIM, C. J. WILSON, J. LEHÁR, G. V. KRYUKOV, D. SONKIN, ET AL., *The cancer cell line encyclopedia enables predictive modelling of anticancer drug sensitivity*, Nature, 483 (2012), p. 603.

[4] P. J. BICKEL, Y. RITOV, A. B. TSYBAKOV, ET AL., *Simultaneous analysis of lasso and dantzig selector*, The Annals of Statistics, 37 (2009), pp. 1705–1732.

[5] T. BLUMENSATH AND M. E. DAVIES, *Iterative hard thresholding for compressed sensing*, Applied and computational harmonic analysis, 27 (2009), pp. 265–274.

[6] S. BOUCHERON, G. LUGOSI, AND P. MASSART, *Concentration Inequalities: A Nonasymptotic Theory of Independence*, Oxford University Press, 2013.

[7] P. T. BOUFOUNOS AND R. G. BARANIUK, *1-bit compressive sensing*, in Information Sciences and Systems, 2008. CISS 2008. 42nd Annual Conference on, IEEE, 2008, pp. 16–21.

[8] E. CANDES, T. TAO, ET AL., *The dantzig selector: Statistical estimation when p is much larger than n*, The Annals of Statistics, 35 (2007), pp. 2313–2351.

[9] E. J. CANDÈS AND B. RECHT, *Exact matrix completion via convex optimization*, Foundations of Computational mathematics, 9 (2009), p. 717.

[10] E. J. CANDÈS, J. ROMBERG, AND T. TAO, *Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information*, IEEE Transactions on information theory, 52 (2006), pp. 489–509.

[11] E. J. CANDÈS AND T. TAO, *The power of convex relaxation: Near-optimal matrix completion*, IEEE Transactions on Information Theory, 56 (2010), pp. 2053–2080.

[12] V. CHANDRASEKARAN, B. RECHT, P. A. PARRILO, AND A. S. WILLSKY, *The convex geometry of linear inverse problems*, Foundations of Computational Mathematics, 12 (2012), pp. 805–849.

[13] S. CHATTERJEE, S. CHEN, AND A. BANERJEE, *Generalized dantzig selector: Application to the k-support norm*, in Advances in Neural Information Processing Systems, 2014, pp. 1934–1942.

[14] A. CHEN, A. B. OWEN, AND M. SHI, *Data enriched linear regression*, Electronic journal of statistics, 9 (2015), pp. 1078–1112.

[15] J. CHEN, E. E. BARDES, B. J. ARONOW, AND A. G. JEGGA, *Toppgene suite for gene list enrichment analysis and candidate gene prioritization*, Nucleic acids research, 37 (2009), pp. W305–W311.

[16] J. CHEN, J. LIU, AND J. YE, *Learning incoherent sparse and Low-Rank patterns from multiple tasks*, ACM transactions on knowledge discovery from data, 5 (2012), p. 22.

[17] F. DONDELINGER AND S. MUKHERJEE, *High-dimensional regression over disease subgroups*, arXiv preprint arXiv:1611.00953, (2016).

[18] D. L. DONOHO, *Compressed sensing*, IEEE Transactions on information theory, 52 (2006), pp. 1289–1306.

[19] J. FRIEDMAN, T. HASTIE, AND R. TIBSHIRANI, *Sparse inverse covariance estimation with the graphical lasso*, Biostatistics, 9 (2008), pp. 432–441.

[20] S. M. GROSS AND R. TIBSHIRANI, *Data shared lasso: A novel tool to discover uplift*, Computational Statistics & Data Analysis, 101 (2016), pp. 226–235.

[21] Q. GU AND A. BANERJEE, *High dimensional structured superposition models*, in Advances In Neural Information Processing Systems, 2016, pp. 3684–3692.

[22] F. IORIO, T. A. KNIJNENBURG, D. J. VIS, G. R. BIGNELL, M. P. MENDEN, M. SCHUBERT, N. ABEN, E. GONCALVES, S. BARTHORPE, H. LIGHTFOOT, ET AL., *A landscape of pharmacogenomic interactions in cancer*, Cell, 166 (2016), pp. 740–754.

[23] P. JAIN, P. NETRAPALLI, AND S. SANGHAVI, *Low-rank matrix completion using alternating minimization*, in Proceedings of the forty-fifth annual ACM symposium on Theory of computing, ACM, 2013, pp. 665–674.

[24] A. JALALI, P. RAVIKUMAR, S. SANGHAVI, AND C. RUAN, *A Dirty Model for Multi-task Learning*, in Advances in Neural Information Processing Systems, 2010, pp. 964–972.

[25] M. B. MCCOY AND J. A. TROPP, *The achievable performance of convex demixing*, arXiv preprint arXiv:1309.7478, (2013).

[26] S. MENDELSON, *Learning Without Concentration*, in Journal of the ACM (JACM), To appear, 2014.

827  [27]  S. NEGAHBAN AND M. J. WAINWRIGHT, *Restricted strong convexity and weighted matrix completion: Optimal*
828        *bounds with noise*, Journal of Machine Learning Research, 13 (2012), pp. 1665–1697.
829  [28]  S. NEGAHBAN, B. YU, M. J. WAINWRIGHT, AND P. K. RAVIKUMAR, *A unified framework for high-dimensional*
830        *analysis of m-estimators with decomposable regularizers*, in Advances in Neural Information Processing Systems,
831        2009, pp. 1348–1356.
832  [29]  S. N. NEGAHBAN, P. RAVIKUMAR, M. J. WAINWRIGHT, AND B. YU, *A Unified Framework for High-Dimensional*
833        *Analysis of $M$-Estimators with Decomposable Regularizers*, Statistical Science, 27 (2012), pp. 538–557.
834  [30]  E. OLLIER AND V. VIALLON, *Joint estimation of k related regression models with simple l_1-norm penalties*, arXiv
835        preprint arXiv:1411.1594, (2014).
836  [31]  E. OLLIER AND V. VIALLON, *Regression modeling on stratified data with the lasso*, arXiv preprint arXiv:1508.05476,
837        (2015).
838  [32]  S. OYMAK, B. RECHT, AND M. SOLTANOLKOTABI, *Sharp time–data tradeoffs for linear inverse problems*, arXiv
839        preprint arXiv:1507.04793, (2015).
840  [33]  Y. PLAN, R. VERSHYNIN, AND E. YUDOVINA, *High-dimensional estimation with geometric constraints*, Information
841        and Inference: A Journal of the IMA, 6 (2017), pp. 1–40.
842  [34]  G. RASKUTTI, M. J. WAINWRIGHT, AND B. YU, *Restricted eigenvalue properties for correlated gaussian designs*,
843        Journal of Machine Learning Research, 11 (2010), pp. 2241–2259.
844  [35]  M. RUDELSON AND S. ZHOU, *Reconstruction from anisotropic random measurements*, IEEE Transactions on
845        Information Theory, 59 (2013), pp. 3434–3447.
846  [36]  R. TIBSHIRANI, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society. Series B
847        (Methodological), (1996), pp. 267–288.
848  [37]  R. TIBSHIRANI, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society. Series B
849        (Methodological), (1996), pp. 267–288.
850  [38]  J. A. TROPP, *Convex recovery of a structured signal from independent random linear measurements*, in Sampling
851        Theory - a Renaissance, To appear, may 2015, https://arxiv.org/abs/1405.1102.
852  [39]  R. VERSHYNIN, *Introduction to the non-asymptotic analysis of random matrices*, in Compressed Sensing, Cambridge
853        University Press, Cambridge, 2012, pp. 210–268.
854  [40]  R. VERSHYNIN, *High-dimensional probability: An introduction with applications in data science*, vol. 47, Cambridge
855        University Press, 2018.
856  [41]  E. YANG AND P. RAVIKUMAR, *Dirty statistical models*, in Advances in Neural Information Processing Systems, 2013,
857        pp. 611–619.
858  [42]  Y. ZHANG AND Q. YANG, *A survey on Multi-Task learning*, (2017), https://arxiv.org/abs/1707.08114.