# High Dimensional Data Sharing:
# Multi-Task Learning with Theoretical Guarantee [*]

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

**Abstract.** We consider the problem of multi-task learning in high dimension. In particular, we introduce an estimator and investigate its statistical and computational properties for the problem of multiple connected linear regressions known as Data Sharing. The between-tasks connections are captured by a cross-tasks *common parameter* which gets refined by per-task *individual parameters*. Any convex function, e.g., norm, can characterize the structure of both common and individual parameters. We delineate the sample complexity of our estimator and provide high probability non-asymptotic bound for estimation error of all parameters under a geometric condition. We show that the recovery of the common parameter benefits from *all* of the pooled samples. We propose an iterative estimation algorithm with a geometric convergence rate and supplement our theoretical analysis with experiments on synthetic data. Overall, we present a first through statistical and computational analysis of inference in the data sharing model.

**Key words.** multi-task learning, superposition models, high-dimensional statistics, convergence rate analysis

**AMS subject classifications.** 62F10, 62J05, 90C25

**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 [13, 19, 20]. 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 [11, 12, 39]. 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 [3, 30], broad families of semi-parametric and single-index models [9, 36], non-convex models [7, 24], 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 the 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
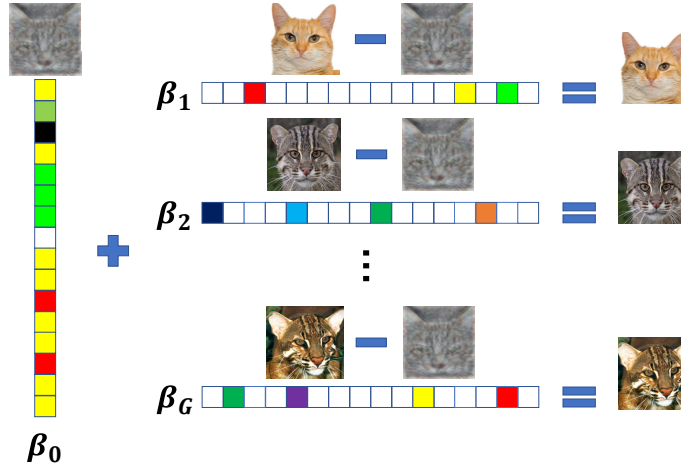
Figure 1: A conceptual illustration of data sharing model for learning representation of cat species. The common parameter $\boldsymbol{\beta}_0$ captures a *generic cat* which consists of shared features among all cats.

35  group which is not the case in several problems, e.g., anti-cancer drug sensitivity prediction [5, 23].

36  The middle ground model for such a scenario is the *superposition* of common and individual
37  parameters $\boldsymbol{\beta}_0^* + \boldsymbol{\beta}_g^*$ which has been of recent interest [22, 27, 43]. Such a collection of *coupled*
38  superposition models is known by multiple names in the statistical machine learning community. It is
39  a form of multi-task learning [25, 45] when we consider regression in each group as a task. It is also
40  called data sharing [21] since information contained in different groups is shared through the common
41  parameter $\boldsymbol{\beta}_0^*$. And finally, it has been called data enrichment [1, 2, 16] because we enrich our data set
42  with pooling multiple samples from different but related sources.

43  Following the successful application of such a modeling scheme in recent years [18, 21, 32, 33],
44  we consider the below *data sharing* (DS) model:

45  (1.1)
$$y_{gi} = \phi(\mathbf{x}_{gi}, (\boldsymbol{\beta}_0^* + \boldsymbol{\beta}_g^*)) + w_{gi}, \quad g \in \{1, \dots, G\},$$

46  where $g$ and $i$ index the group and samples respectively. The DS model (1.1) assumes that there is a
47  *common* parameter $\boldsymbol{\beta}_0^*$ shared between all groups which models similarities between all samples. And
48  there are *individual* per-group parameters $\boldsymbol{\beta}_g^*$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 DS (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 $\boldsymbol{\beta}_g$ are structured, i.e., for suitable convex functions $f_g$'s, $f_g(\boldsymbol{\beta}_g)$ 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., $\phi(\mathbf{x}, \boldsymbol{\beta}) = \mathbf{x}^T \boldsymbol{\beta}$. 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 [26, 29, 34, 35, 44].

**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 [45]. To capture the relation between tasks, different types of constraints are assumed for parameter matrices. For example, [17] 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 Statistical Model (DSM) proposed in [25] where authors regularize the regression with $\|\mathbf{B}_1\|_{1,\infty}$ and $\|\mathbf{B}_2\|_{1,1}$ where norms are $i, j$-norms on *rows*, $\mathbf{b}$, of matrices, i.e., $\|\mathbf{B}\|_{i,j} = \|(\|\mathbf{b}_1\|_j, \ldots, \|\mathbf{b}_p\|_j)\|_i$ and the norms are defined as $\|\mathbf{b}\|_i = \left(\sum_{g=1}^{G} |b_g|^i\right)^{1/i}$ and $\|\mathbf{b}\|_\infty = \max_{g \in G} |b_g|$.

If in our DS model we pick all structure inducing functions $f_g$ to be $l_1$-norm, the resulting model is very similar to the DSM 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 DSM model is higher than DS 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 DS we have a single common parameter $\boldsymbol{\beta}_0^*$. So one would expect that DS estimators should have smaller sample complexity compared to their DSM counterparts and our analysis confirm that our estimator is more data efficient than DSM estimator of [25], Table 1. Mainly, DSM requires every task $i$ to have large enough samples to learn its own common parameters $\boldsymbol{\beta}_i$ but since DS shares the common parameter it only requires the *total dataset over all tasks* to be sufficiently large.

The linear DS 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 [18], sentiment analysis, banking strategy [21], single cell data analysis [33], road safety [32], and disease subtype analysis [18]. More generally, in any high-dimensional problem where the population consists of groups, data sharing 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 sharing framework to different domains, limited advances have been made in understanding the statistical and computational properties of suitable estimators for the DS model (1.1). In fact, non-asymptotic statistical properties, including sample complexity and statistical rates of convergence, of regularized estimators for the data sharing model is still an open question [21, 32]. To the best of our knowledge, the only theoretical guarantee for data sharing is provided in [33] where authors prove sparsistency of their proposed method under the irrepresentability condition of the design matrix for recovering *supports* of common and individual parameters. Existing support recovery guarantees [33], sample complexity and $l_2$ consistency results [25] 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. 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 sharing model takes the following vector form:

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

100    where each row of $\mathbf{X}_g$ is $\mathbf{x}_{gi}^T$ and $\mathbf{w}_g^T = (w_{g1}, \dots, w_{gn_g})$ is the noise vector. It is useful for indexing to

101    consider the common parameter as the zeroth group and define $n_0 \triangleq n$ and $\mathbf{X}_0 \triangleq [\mathbf{X}_1^T, \dots, \mathbf{X}_G^T]^T$.

102        *Sub-Gaussian random variable and vector.* A random variable $V$ is sub-Gaussian if its

103    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

104    norm of $V$, denoted by $\|V\|_{\psi_2}$ [41]. A random vector $\mathbf{v} \in \mathbb{R}^p$ is sub-Gaussian if the one-dimensional

105    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 $\mathbf{v}$

106    is defined [41] 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

107    set $\mathcal{V}$ is defined as $\omega(\mathcal{V}) = \mathbb{E}_{\mathbf{g}} \left[ \sup_{\mathbf{u} \in \mathcal{V}} \langle \mathbf{g}, \mathbf{u} \rangle \right]$ [42], where the expectation is over $\mathbf{g} \sim N(\mathbf{0}, \mathbf{I}_{p \times p})$,

108    a vector of independent zero-mean unit-variance Gaussian. The marginal tail function is defined as

109    $Q_\xi(\mathbf{u}) = \mathbb{P}(|\langle \mathbf{x}, , \mathbf{u} \rangle| > \xi)$ for a fixed vector $\mathbf{u}$, random vector $\mathbf{x}$ and constant $\xi > 0$.

110        **1.3. Our Contributions.** We propose the following Data Shared (DS) estimator $\hat{\boldsymbol{\beta}}$ for recovering

111    the structured parameters where the structure is induced by *convex* functions $f_g(\cdot)$:

112    (1.3)        $\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}}_0^T, \dots, \hat{\boldsymbol{\beta}}_G^T) \in \underset{\boldsymbol{\beta}_0, \dots, \boldsymbol{\beta}_G}{\operatorname{argmin}} \frac{1}{n} \sum_{g=1}^{G} \|\mathbf{y}_g - \mathbf{X}_g(\boldsymbol{\beta}_0 + \boldsymbol{\beta}_g)\|_2^2, \text{ s.t. } \forall g \in [G_+] : f_g(\boldsymbol{\beta}_g) \leq f_g(\boldsymbol{\beta}_g^*).$

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

114    • The DS estimator (1.3) succeeds if a geometric condition that we call *DAta SHaring Incoherence*

115      *conditioN* (DASHIN) is satisfied, Figure 2b. Compared to other known geometric conditions in

116      the literature such as structural coherence [22] and stable recovery conditions [27], DASHIN is a

117      considerably weaker condition, Figure 2a.

118    • Assuming DASHIN holds, we establish a high probability non-asymptotic bound on the weighted

119      sum of parameter-wise estimation error, $\boldsymbol{\delta}_g = \hat{\boldsymbol{\beta}}_g - \boldsymbol{\beta}_g^*$ as:

120        (1.4)                    $\sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 \leq \gamma O \left( \frac{\max_{g \in [G]} \omega(\mathcal{C}_g \cap \mathbb{S}^{p-1})}{\sqrt{n}} \right),$

121      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

122      $\mathcal{C}_g$ is the error cone corresponding to $\boldsymbol{\beta}_g^*$ exactly defined in section 2. To the best of our knowledge,

123      this is the first statistical estimation guarantee for the data sharing.

124    • We also establish the sample complexity of the DS estimator for all parameters as $\forall g \in [G_+] : n_g =$

125      $O(\omega(\mathcal{C}_g \cap \mathbb{S}^{p-1}))^2$. We emphasize that our result proofs that the recovery of the common parameter

126      $\boldsymbol{\beta}_0$ by DS estimator (1.3) benefits from *all* of the $n$ pooled samples.

127    • We present an efficient projected block gradient descent algorithm DASHER, to solve DE's objective

128      (1.3) which converges geometrically to the statistical error bound of (1.4). To the best of our

129      knowledge, this is the first rigorous computational result for the high-dimensional data-shared

130      regression.

131        The rest of this paper is organized as follows: First, we characterize the error set of our estimator and

132    provide a deterministic error bound in section 2. Then in section 3, we discuss the restricted eigenvalue

133    condition and calculate the sample complexity required for the recovery of the true parameters by

134    our estimator under DASHIN condition. We close the statistical analysis in section 4 by providing

135    non-asymptotic high probability error bound for parameter recovery. We delineate our geometrically

136    convergent algorithm, DASHER in section 5 and finally supplement our work with experiments on

137    synthetic data in section 6.

(a) Structural Coherence (SC) condition.  (b) DAta SHaring Incoherence conditioN (DASHIN).
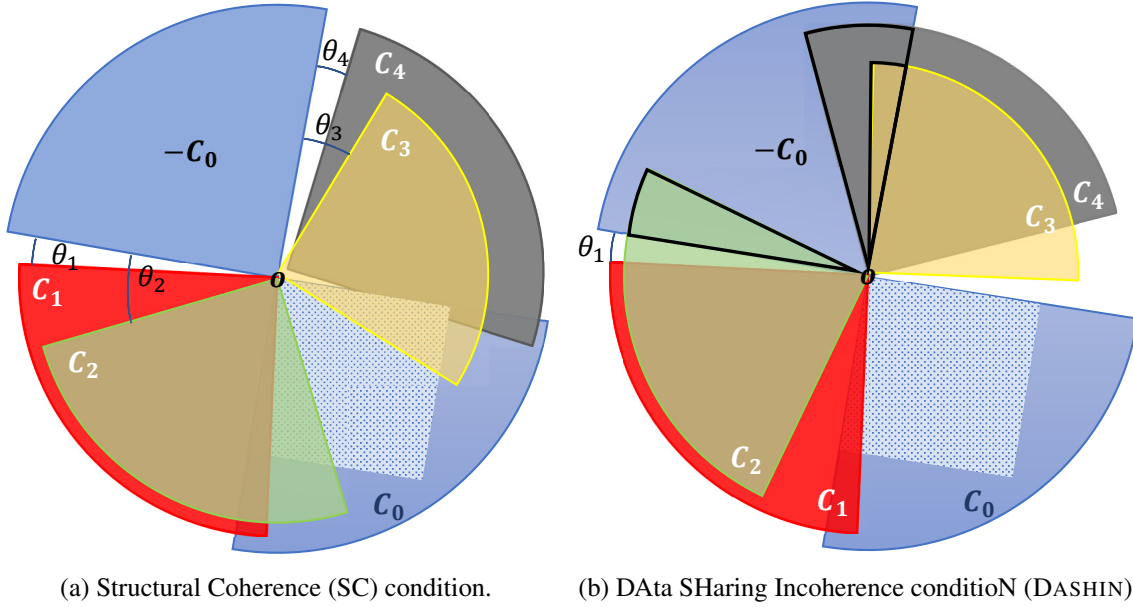
Figure 2: Comparison of geometric recovery condition for superposition models known as Structural Coherence (SC) [22] and our DASHIN recovery condition for data sharing model which is a system of coupled superposition models (1.2). For each parameter $\boldsymbol{\beta}_g^* \in [G]$, $\mathcal{E}_g = \left\{ \boldsymbol{\delta}_g | f_g(\boldsymbol{\beta}_g^* + \boldsymbol{\delta}_g) \leq f_g(\boldsymbol{\beta}_g^*) \right\}$ is the error set and $\mathcal{C}_g = \text{Cone}(\mathcal{E}_g)$ is the error cone. For all $i, j$, SC requires $-\mathcal{C}_i \cap \mathcal{C}_j = \{0\}$. In panel (a) we only show this condition for $i = 0$, i.e., $-\mathcal{C}_0 \cap \mathcal{C}_j = \{0\}$ where all $\theta_j > 0$. DASHINon the other hand only needs one of the $\mathcal{C}_g, g \in [G]$, does not intersect with the inverse of the common parameter error cone $-\mathcal{C}_0$. In panel (b) $-\mathcal{C}_0 \cap \mathcal{C}_1 = \{0\}$ is enough for recovering all parameters.

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

(2.1) $$\hat{\boldsymbol{\beta}} \in \underset{\boldsymbol{\beta}}{\arg\min} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2, \quad \text{s.t. } \forall g \in [G_+] : f_g(\boldsymbol{\beta}_g) \leq f_g(\boldsymbol{\beta}_g^*),$$

where $\mathbf{y} = (\mathbf{y}_1^T, \ldots \mathbf{y}_G^T)^T \in \mathbb{R}^n$, $\boldsymbol{\beta} = (\boldsymbol{\beta}_0^T, \ldots, \boldsymbol{\beta}_G^T)^T \in \mathbb{R}^{(G+1)p}$ 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.,$|\text{supp}(\boldsymbol{\beta}_g^*)| = s_g$ by using $l_1$-norm as the sparsity inducing function, our DS estimator of (2.1) becomes:

(2.3) $$\hat{\boldsymbol{\beta}} \in \underset{\boldsymbol{\beta}}{\arg\min} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2, \quad \text{s.t. } \forall g \in [G_+] : \|\boldsymbol{\beta}_g\|_1 \leq \|\boldsymbol{\beta}_g^*\|_1.$$

We call Example 2.1 *sparse DS* 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:

$$(2.4) \qquad \mathcal{E}_g = \left\{ \boldsymbol{\delta}_g \mid f_g(\boldsymbol{\beta}_g^* + \boldsymbol{\delta}_g) \le f_g(\boldsymbol{\beta}_g^*) \right\}, \quad g \in [G_+].$$

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

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

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

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

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

*Proof.* We lower bound the LHS and upper bound the RHS of the optimality inequality $\frac{1}{n}\|\mathbf{X}\boldsymbol{\delta}\|_2^2 \le \frac{1}{n} 2\mathbf{w}^T \mathbf{X}\boldsymbol{\delta}$ 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:

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

$$(2.6) \qquad \ge \kappa \left( \min_{g \in [G]} \sqrt{\frac{n_g}{n}} \right)^2 \left( \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 \right)^2 = \kappa \left( \min_{g \in [G]} \frac{n_g}{n} \right) \left( \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 \right)^2$$

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

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

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

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

173     **3. Restricted Eigenvalue Condition.** The main assumption of Theorem 2.2 is known as
174 Restricted Eigenvalue (RE) condition in the literature of high-dimensional statistics [4, 31, 37]:
175 $\inf_{\mathbf{u} \in \mathcal{H}} \frac{1}{n} \|\mathbf{X}\mathbf{u}\|_2^2 \geq \kappa > 0$. The RE condition posits that the minimum eigenvalues of the matrix
176 $\mathbf{X}^T\mathbf{X}$ in directions restricted to $\mathcal{H}$ is strictly positive. In this section, we show that for the design
177 matrix $\mathbf{X}$ defined in (2.2), the RE condition holds with high probability under a suitable geometric
178 condition we call *DAta SHaring Incoherence conditioN* (DASHIN) and for enough number of samples.
179 We precisely characterize total and per-group sample complexities required for successful parameter
180 recovery. For the analysis, similar to existing work [22, 28, 40], we assume the design matrix to be
181 isotropic sub-Gaussian.[1]

182     **Definition 3.1.** *We assume* $\mathbf{x}_{gi}$ *are i.i.d. random vectors from a non-degenerate zero-mean,*
183 *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*
184 *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*
185 *i.i.d. zero-mean, unit-variance sub-Gaussian with* $\|\|\mathbf{w}_{gi}\|\|_{\psi_2} \leq k_w$.

186     **3.1. Geometric Condition for Recovery.** Unlike standard high-dimensional statistical esti-
187 mation, for RE condition to be true, parameters of superposition models need to satisfy geometric
188 conditions which limits the interaction of the error cones of parameters with each other to make sure
189 that recovery is possible. In this section, we elaborate our sufficient geometric condition for recovery
190 and compare it with state-of-the-art condition for recovery of superposition models.
191     To intuitively illustrate the necessity of such a geometric condition, consider the simplest superpo-
192 sition model i.e., $\boldsymbol{\beta}_0^* + \boldsymbol{\beta}_g^*$. Without any restriction on interactions of error cones, any estimates such
193 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
194 $\boldsymbol{\delta}_g \neq -\boldsymbol{\delta}_0$. In general, the RE condition of individual superposition models can be established under
195 the so-called Structural Coherence (SC) condition [22, 27] which is the generalization of this idea for
196 superposition of multiple parameters as $\sum_{g=0}^G \boldsymbol{\beta}_g^*$.

197     **Definition 3.2 (Structural Coherence (SC) [22, 27]).** *Consider a superposition model of*
198 *the form* $y = \mathbf{x}^T \sum_{g=0}^G \boldsymbol{\beta}_g^* + w$. *The SC condition requires that* $\forall \boldsymbol{\delta}_g \in \mathcal{C}_g, \exists \lambda$ *s.t.* $\|\sum_{g=0}^G \boldsymbol{\delta}_g\|_2 \geq$
199 $\lambda \sum_{g=0}^G \|\boldsymbol{\delta}_g\|_2$, *and leads to the RE condition* $\frac{1}{\sqrt{n}}\|\mathbf{X}\sum_{g=1}^G \boldsymbol{\delta}_g\|_2 \geq \kappa \sum_{g=1}^G \|\boldsymbol{\delta}_g\|_2$.

200     *Remark* 3.3. Note that the SC condition is satisfied if none of the individual error cones $\mathcal{C}_g$ intersect
201 with the inverted error cone $-\mathcal{C}_0$ [22, 40], i.e., $\forall g, \theta_g > 0$ in Figure 2a where

$$\cos(\theta_g) = \sup_{\boldsymbol{\delta}_0 \in \mathcal{C}_0, \boldsymbol{\delta}_g \in \mathcal{C}_g} -\langle \boldsymbol{\delta}_0/\|\boldsymbol{\delta}_0\|_2, \boldsymbol{\delta}_g/\|\boldsymbol{\delta}_g\|_2 \rangle.$$

203 Next, we introduce DASHIN, a considerably weaker geometric condition compared to SC which leads
204 to recovery of all parameters in the data sharing model.

205     **Definition 3.4 (DAta SHaring Incoherence conditioN (DASHIN)).** *There exists a non-empty*
206 *set* $\mathcal{I} \subseteq [G]$ *of groups where for some scalars* $0 < \bar{\rho} \leq 1$ *and* $\lambda_{\min} > 0$ *the following holds:*
207     *1.* $\sum_{g \in \mathcal{I}} n_g \geq \lceil \bar{\rho} n \rceil$.
208     *2.* $\forall g \in \mathcal{I}, \forall \boldsymbol{\delta}_g \in \mathcal{C}_g$, *and* $\boldsymbol{\delta}_0 \in \mathcal{C}_0$: $\|\boldsymbol{\delta}_g + \boldsymbol{\delta}_0\|_2 \geq \lambda_{\min}(\|\boldsymbol{\delta}_0\|_2 + \|\boldsymbol{\delta}_g\|_2)$
209 *Observe that* $0 < \lambda_{\min}, \bar{\rho} \leq 1$ *by definition.*

---

[1]Extension to an-isotropic sub-Gaussian case is straightforward by techniques developed in [4, 38].

*Remark* 3.5. DASHIN is a refinement of SC for the specific problem of data sharing, i.e., system of coupled superposition model each with two components. DASHIN holds even if only one of the $\mathcal{C}_g$s does not intersect with $-\mathcal{C}_0$. More specifically, DASHIN holds if $\exists g, \theta_g > 0$ in Figure 2b which allows $-\mathcal{C}_0$ to intersect with an arbitrarily large fraction of the $\mathcal{C}_g$ cones and as the number of intersections increases, our final error bound becomes looser.

**3.2. Sample Complexity.** An alternative to our DS estimator (1.3) may be based on $G$ *isolated* superposition model $\mathbf{y}_g = \mathbf{X}_g(\boldsymbol{\beta}_0^* + \boldsymbol{\beta}_g^*) + \mathbf{w}_g$ each with two components. Now, if SC holds for at least one of the superposition models, i.e., $\exists g, -\mathcal{C}_0 \cap \mathcal{C}_g = \{0\}$, one can recover $\hat{\boldsymbol{\beta}}_0$ and plug it in to the remaining $G - 1$ superposition estimators to estimate the corresponding $\hat{\boldsymbol{\beta}}_g$s. We call such an estimator, *plugin superposition* estimator for which it seems that DASHIN has no advantage over SC. But the disadvantage of plugin superposition estimator is that it fails to utilize the true coupling structure in the data sharing model, where $\boldsymbol{\beta}_0^*$ is involved in all groups. In fact, below we show that the plugin superposition estimator under SC condition leads to a pessimistic sample complexity for $\boldsymbol{\beta}_0^*$ recovery.

**Proposition 3.6.** *Assume observations distributed as defined in Definition 3.1 and pair-wise SC conditions are satisfied. Consider each superposition model (1.2) in isolation; to recover the common parameter $\boldsymbol{\beta}_0^*$ plugin superposition requires at least one group $i$ to have $n_i = O(\max(\omega^2(\mathcal{A}_0), \omega^2(\mathcal{A}_i)))$. To recover the rest of individual parameters, it needs $\forall g \neq i : n_g = O(\omega^2(\mathcal{A}_g))$ samples.*

In other words, by separate analysis of superposition estimators at least one problem needs to have sufficient samples for recovering the common parameter $\boldsymbol{\beta}_0$ and therefore the common parameter recovery does not benefit from the pooled $n$ samples. But given the nature of coupling in the data sharing model, we hope to be able to get a better sample complexity specifically for the common parameter $\boldsymbol{\beta}_0$. Using DASHIN and the small ball method [28], a tool from empirical process theory in the following theorem, we get a better sample complexity required for satisfying the RE condition:

**Theorem 3.7.** *Let $\mathbf{x}_{gi}$s be random vectors defined in Definition 3.1. Assume DASHIN condition of Definition 3.4 holds for error cones $\mathcal{C}_g$s and $\psi_{\mathcal{I}} = \min\{1/2, \lambda_{\min}\bar{\rho}/3\}$. Then, for all $\boldsymbol{\delta} \in \mathcal{H}$, when we have enough number of samples as $\forall g \in [G_+] : n_g \geq m_g = O(k_x^6 \alpha^{-6} \psi_{\mathcal{I}}^{-2} \omega(\mathcal{A}_g)^2)$, with probability at least $1 - e^{-n\kappa_{\min}/4}$ we have:*

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

*where $\kappa_{\min} = \min_{g \in [G_+]} C\psi_{\mathcal{I}} \frac{\alpha^3}{k_x^2} - \frac{2c_g k_x \omega(\mathcal{A}_g)}{\sqrt{n_g}}$.*

*Remark* 3.8. Note that $\kappa = \frac{\kappa_{\min}^2}{4}$ is the lower bound of the RE condition of Theorem 2.2, i.e., $0 < \kappa \leq \inf_{\mathbf{u} \in \mathcal{H}} \frac{1}{n} \|\mathbf{X}\mathbf{u}\|_2^2$ and is determined by the group with the worst individual RE condition.

*Example* 3.9. ($l_1$-**norm**) The Gaussian width of the spherical cap of a $p$-dimensional $s$-sparse vector is $\omega(\mathcal{A}) = \Theta(\sqrt{s \log p})$ [4, 42]. Therefore, the number of samples per group and total required for satisfaction of the RE condition in the sparse DS estimator Example 2.1 is $\forall g \in [G] : n_g \geq m_g = \Theta(s_g \log p)$. Table 1 compares sample complexities of the sparse DS estimator with three baselines: plugin superposition estimator of Proposition 3.6, G Independent LASSO (GI-LASSO), and Jalali's Dirty Statistical Model (DSM) [25]. Note that GI-LASSO does not recover the common parameter and DSM needs all groups have same number of samples.

| | **GI-LASSO** | **Dirty Stat. Model** | **Plugin Superposition** | **Sparse DS** |
|---|---|---|---|---|
| $m_g$ | $s_{0g} \log p$ | $G \max_{g \in [G]} s_{0g} \log(p)$ | $\exists i \in [G] : \max(s_0, s_i) \log p$ $\forall g \neq i : s_g \log p$ | $s_g \log p$ |

Table 1: Comparison of the order of per group number of samples (sample complexities) of various methods for recovering sparse DS parameters. Let $s_{0g} = |\text{support}(\beta_0^* + \beta_g^*)|$ be the superimposed support where $s_0, s_g \leq \max(s_0, s_g) \leq s_{0g}$.

### 3.3. Proof of Theorem 3.7. Let's simplify the LHS of the RE condition:

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

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

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

$$= \frac{1}{n} \sum_{g=1}^{G} \xi_g \sum_{i=1}^{n_g} \mathbb{1}\left( |\langle \mathbf{x}_{gi}, \boldsymbol{\delta}_{0g} \rangle| \geq \xi_g \right),$$

where to avoid cluttering we denoted $\boldsymbol{\delta}_{0g} = \boldsymbol{\delta}_0 + \boldsymbol{\delta}_g$ and $\xi_g = \xi \|\boldsymbol{\delta}_{0g}\|_2 > 0$. Now we add and subtract the corresponding per-group marginal tail function, $Q_{\xi_g}(\boldsymbol{\delta}_{0g}) = \mathbb{P}(|\langle \mathbf{x}, , \boldsymbol{\delta}_{0g} \rangle| > \xi_g)$ and take inf:

$$\inf_{\boldsymbol{\delta} \in \mathcal{H}} \frac{1}{\sqrt{n}} \|\mathbf{X}\boldsymbol{\delta}\|_2 \geq \inf_{\boldsymbol{\delta} \in \mathcal{H}} \sum_{g=1}^{G} \frac{n_g}{n} \xi_g Q_{2\xi_g}(\boldsymbol{\delta}_{0g}) - \sup_{\boldsymbol{\delta} \in \mathcal{H}} \frac{1}{n} \sum_{g=1}^{G} \xi_g \sum_{i=1}^{n_g} \left[ Q_{2\xi_g}(\boldsymbol{\delta}_{0g}) - \mathbb{1}(|\langle \mathbf{x}_{gi}, \boldsymbol{\delta}_{0g} \rangle| \geq \xi_g) \right]$$

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

For the ease of exposition we consider the LHS of (3.1) 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 $t_2$.

### 3.3.1. Lower Bounding the First Term. Our main result is the following lemma which uses the DASHIN condition of the Definition 3.4 and provides a lower bound for the first term $t_1(\mathbf{X})$:

Lemma 3.10. *Suppose* DASHIN *holds. Let* $\psi_{\mathcal{I}} = \frac{\lambda_{\min}\bar{\rho}}{3}$. *For any* $\boldsymbol{\delta} \in \mathcal{H}$, *we have:*

(3.2) $$\sum_{g=1}^{G} \frac{n_g}{n} \xi_g Q_{2\xi_g}(\boldsymbol{\delta}_{0g}) \geq \psi_{\mathcal{I}} \xi \frac{(\alpha - 2\xi)^2}{4ck_x^2} \sum_{g=0}^{n} \frac{n_g}{n} \|\boldsymbol{\delta}_g\|_2,$$

Lemma 3.10 implies that $t_1(\mathbf{X})$ is lower bounded by the same RHS bound of (3.2).

**3.3.2. Upper Bounding the Second Term.** First we show $t_2(\mathbf{X})$ satisfies the bounded difference property defined in Section 3.2. of [8], i.e., by changing each of $\mathbf{x}_{gi}$ the value of $t_2(\mathbf{X})$ at most change by one. We rewrite $t_2$ as $t_2(\mathbf{X}) = \sup_{\boldsymbol{\delta} \in \mathcal{H}} g_\delta(\mathbf{X})$ where $g_\delta(\mathbf{X})$ is the argument of sup in (3.1). Now we denote the design matrix resulted from replacement of $k$th sample from $j$th group $\mathbf{x}_{jk}$ with another sample $\mathbf{x}'_{jk}$ by $\mathbf{X}'_{jk}$. Then our goal is to show $\forall j \in [G], k \in [n_j], \sup_{\mathbf{X},\mathbf{x}'_{jk}} |t_2(\mathbf{X}) - t_2(\mathbf{X}'_{jk})| \leq c_i$ 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:

$$\sup_{\mathbf{X},\mathbf{x}'_{jk}} |t_2(\mathbf{X}) - t_2(\mathbf{X}'_{jk})| \leq \sup_{\mathbf{X},\mathbf{x}'_{jk}} \sup_{\boldsymbol{\delta} \in \mathcal{H}} \left| g(\mathbf{X}) - g(\mathbf{X}'_{jk}) \right|$$

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

$$\leq \sup_{j} \sup_{\boldsymbol{\delta} \in \mathcal{H}} \frac{\xi_j}{n} = \frac{\xi}{n} \sup_{j} \sup_{\boldsymbol{\delta} \in \mathcal{H}} \|\boldsymbol{\delta}_0 + \boldsymbol{\delta}_j\|_2$$

$$\leq \frac{\xi}{n} \sup_{j} \sup_{\boldsymbol{\delta} \in \mathcal{H}} \|\boldsymbol{\delta}_0\|_2 + \|\boldsymbol{\delta}_j\|_2$$

$$(\boldsymbol{\delta} \in \mathcal{H}) \leq \xi \left( \frac{1}{n} + \frac{1}{n_j} \right) \leq \frac{2\xi}{n}$$

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 $\|\boldsymbol{\delta}_g\|_2 \leq \frac{n}{n_g}$ which justifies the last inequality. Now, we can invoke the bounded difference inequality from Theorem 6.2 of [8] 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 $t_2(\mathbf{X})$ using the following lemma:

**Lemma 3.11.** *For the random vector $\mathbf{x}$ of Definition 3.1, we have the following bound:*

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

**3.3.3. Continuing the Proof of Theorem 3.7.** Define $q \triangleq \frac{(\alpha - 2\xi)^2}{4ck^2}$. Putting back bounds of $t_1(\mathbf{X})$ and $t_2(\mathbf{X})$ together from Lemmas 3.10 and 3.11, with probability at least $1 - e^{-\frac{\tau^2}{2}}$ we have:

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

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

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

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

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

where $\kappa_{\min} = \min_{g\in[G]}\kappa_g$. To conclude the proof, take $\tau = \sqrt{n}\kappa_{\min}/2$.

To satisfy the RE condition all $\kappa_g$s should be bounded away from zero. To this end we need the following sample complexities $\forall g\in[G_+] : \left(\frac{2c_g k}{\psi_{\mathcal{I}}\xi q}\right)^2\omega(\mathcal{A}_g)^2 \leq n_g$ where by taking $\xi = \frac{\alpha}{6}$ simplifies to: $\forall g\in[G_+] : O\left(k^6\psi_{\mathcal{I}}^{-2}\alpha^{-6}\omega(\mathcal{A}_g)^2\right) \leq n_g$ ∎

**4. General Error Bound.** In this section, we present our main statistical result which is a non-asymptotic high probability upper bound for the estimation error of the common and individual parameters.

**Theorem 4.1.** *For $\mathbf{x}_{gi}$ and $w_{gi}$ 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 for estimator (2.1) with probability at least $1 - \sigma\exp\left(-\min\left[\nu\min_{g\in[G]}n_g - \log(G+1), \tau^2\right]\right)$:*

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

*where $\gamma = \max_{g\in[G]}n/n_g$, $\tau > 0$, and $\sigma, \nu,$ and $C$ are constants.*

**Corollary 4.2.** *From Theorem 4.1 one can immediately entail the error bound for estimation of the common and individual parameters as follows:*

$$\forall g\in[G_+] : \quad \|\boldsymbol{\delta}_g\|_2 = O\left(\gamma\frac{\max_{g\in[G_+]}\omega(\mathcal{A}_g) + \sqrt{\log(G+1)}}{\sqrt{n_g}}\right)$$

*Example 4.3.* For the balanced sample condition number $\gamma = \Theta(G)$ discussed in Remark 2.3 we have the following error bound for all parameters:

$$(4.2) \qquad \forall g\in[G_+] : \quad \|\boldsymbol{\delta}_g\|_2 = O\left(G^{3/2}\frac{\max_{g\in[G_+]}\omega(\mathcal{A}_g) + \sqrt{\log(G+1)}}{\sqrt{n}}\right)$$

where the upper bound of error scales as $\frac{1}{\sqrt{n}}$ for all parameters.

*Example* 4.4. ($l_1$-**norm**) For the sparse DS estimator of Example 2.1, results of Theorems 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 upper bound of error simplifies to:

$$(4.3) \qquad \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2 = O\left(\gamma \sqrt{\frac{(\max_{g \in [G_+]} s_g) \log p}{n}}\right),$$

Therefore, individual errors are bounded as $\|\boldsymbol{\delta}_g\|_2 = O(\gamma\sqrt{(\max_{g \in [G]} s_g) \log p/n_g})$ 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 [4, 6, 10, 14, 15].

**4.1. Proof of Theorem 4.1.** To avoid cluttering the notation, we rename the vector of all noises as $\mathbf{w}_0 \triangleq \mathbf{w}$. First, we massage the deterministic upper bound of Theorem 2.2 as follows:

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

Assume $q_g = \langle \mathbf{X}_g^T \frac{\mathbf{w}_g}{\|\mathbf{w}_g\|_2}, \frac{\boldsymbol{\delta}_g}{\|\boldsymbol{\delta}_g\|_2} \rangle \sqrt{\frac{n}{n_g}} \|\mathbf{w}_g\|_2$ and $p_g = \sqrt{\frac{n_g}{n}} \|\boldsymbol{\delta}_g\|_2$. Then the above term is the inner product of two vectors $\mathbf{p} = (p_0, \ldots, p_G)$ and $\mathbf{q} = (q_0, \ldots, q_G)$ for which we have:

$$\sup_{\mathbf{p} \in \overline{\mathcal{H}}} \mathbf{p}^T \mathbf{q} = \sup_{\|\mathbf{p}\|_1 = 1} \mathbf{p}^T \mathbf{q} \leq \|\mathbf{q}\|_\infty = \max_{g \in [G_+]} q_g,$$

where the inequality holds because of the definition of the dual norm. Now we can go back to the original form:

$$(4.4) \qquad \sup_{\boldsymbol{\delta} \in \mathcal{H}} \mathbf{w}^T \mathbf{X} \boldsymbol{\delta} \leq \max_{g \in [G]} \langle \mathbf{X}_g^T \frac{\mathbf{w}_g}{\|\mathbf{w}_g\|_2}, \frac{\boldsymbol{\delta}_g}{\|\boldsymbol{\delta}_g\|_2} \rangle \sqrt{\frac{n}{n_g}} \|\mathbf{w}_g\|_2$$

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

To avoid cluttering we define a random quantity $h_g(\mathbf{w}_g, \mathbf{X}_g) \triangleq \|\mathbf{w}_g\|_2 \sup_{\mathbf{u}_g \in \mathcal{A}_g} \langle \mathbf{X}_g^T \frac{\mathbf{w}_g}{\|\mathbf{w}_g\|_2}, \mathbf{u}_g \rangle$ and a corresponding constant $e_g(\tau) \triangleq c_g \sqrt{(2k_w^2 + 1)k_x^2 n_g} \left(\omega(\mathcal{A}_g) + \sqrt{\log(G+1)} + \tau\right)$. Then from (4.4), we have:

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

$$(\text{Union Bound}) \leq \sum_{g=0}^{G} \mathbb{P}\left(\sqrt{\frac{n}{n_g}} h_g(\mathbf{w}_g, \mathbf{X}_g) > \max_{g \in [G]} \sqrt{\frac{n}{n_g}} e_g(\tau)\right)$$

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

$$\leq (G+1) \max_{g \in [G_+]} \mathbb{P}\left(h_g(\mathbf{w}_g, \mathbf{X}_g) > e_g(\tau)\right)$$

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

---

**Algorithm 5.1** DASHER

---

1: **input: X, y**, learning rates $(\mu_0, \ldots, \mu_G)$, 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**

---

where the last inequality is the result of the following lemma:

**Lemma 4.5.** *For $\mathbf{x}_{gi}$ and $\omega_{gi}$ defined in* Definition *3.1 and $\tau > 0$, with probability at least* $1 - \frac{\sigma_g}{(G+1)} \exp\left( - \min\left[ \nu n_g - \log(G+1), \tau^2 \right] \right)$ *we have:*

$$\|\mathbf{w}_g\|_2 \sup_{\mathbf{u}_g \in \mathcal{A}_g} \langle \mathbf{X}_g^T \frac{\mathbf{w}_g}{\|\mathbf{w}_g\|_2}, \mathbf{u}_g \rangle \le c_g \sqrt{(2k_w^2 + 1)k_x^2 n_g} \left( \omega(\mathcal{A}_g) + \sqrt{\log(G+1)} + \tau \right),$$

*where $\sigma_g, \nu$ and $c_g$ are constants.*

The proof completes by replacing $\max_{g \in [G]} \sqrt{\frac{n}{n_g}} e_g(\tau)$ as the upper bound of $\sup_{\boldsymbol{\delta} \in \mathcal{H}} \mathbf{w}^T \mathbf{X} \boldsymbol{\delta}$ and $\kappa_{\min}^2 / 4$ as the lower bound of $\kappa$ (from Theorem 3.7) both into the bound of Theorem 2.2 . ∎

**5. Estimation Algorithm.** We propose *DAta SHarER* (DASHER) a projected block gradient descent algorithm, Algorithm 5.1, where $\Pi_{\Omega_{f_g}}$ is the Euclidean projection onto the set $\Omega_{f_g}(d_g) = \{ f_g(\boldsymbol{\beta}) \le d_g \}$ where $d_g = f_g(\boldsymbol{\beta}_g^*)$ and is dropped to avoid cluttering.

To analysis convergence properties of DASHER, 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 DASHER, 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 $\mu_g > 0$:*

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

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

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

*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 where to simplify the notation we drop $\mu_g$ arguments.

**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 $t + 1$:*

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

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

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

**Lemma 5.3.** *We have the following recursive dependency between the error of $t + 1$th iteration and $t$th iteration of DASHER:*

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

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

By recursively applying results of Lemma 5.3, we get the following deterministic bound which depends on constants defined in Definition 5.1:

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

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

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:

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

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

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

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

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

377      The RHS of (5.2) consists of two terms. If we keep $\rho < 1$, the first term approaches zero fast,
378 and the second term determines the bound. In the following, we show that for specific choices of step
379 sizes $\mu_g$s we can keep $\rho < 1$ with high probability and the second term can be upper bounded using
380 the analysis of section 4. More specifically, the first term corresponds to the optimization error which
381 shrinks in every iteration while the second term is of the same order of the upper bound of the statistical
382 error characterized in Theorem 4.1.

383      One way for having $\rho < 1$ is to keep all arguments of $\max(\cdots)$ defining $\rho$ strictly below 1. To
384 this end, we first establish high probability upper bound for $\rho_g$, $\eta_g$, and $\phi_g$ (in the subsection SM1.2)
385 and then show that with enough number of samples and proper step sizes $\mu_g$, $\rho$ can be kept strictly
386 below one with high probability. The high probability bounds for constants in Definition 5.1 and
387 the deterministic bound of Theorem 5.2 leads to the following theorem which shows that for enough
388 number of samples, of the same order as the statistical sample complexity of Theorem 3.7, we can keep
389 $\rho$ below one and have geometric convergence.

390      **Theorem 5.4.** *Let* $\tau = \sqrt{\log(G+1)}/\zeta + \epsilon$ *for* $\epsilon, \zeta > 0$. *For the step sizes of:*

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

392 *where* $h_g(\tau) = \left(1 + c_{0g}\frac{\omega(\mathcal{A}_g)+\omega(\mathcal{A}_0)+2\tau}{\sqrt{n_g}}\right)$ *and sample complexities of* $\forall g \in [G_+] : n_g \geq C_g(\omega(\mathcal{A}_g) +$
393 $\tau)^2$, *with probability at least* $1 - \sigma \exp(-\min(\nu \min_{g\in[G]} n_g - \log(G+1), \zeta\epsilon^2))$ *updates of Algo-*
394 *rithm 5.1 obey the following:*

$$\sum_{g=0}^{G} \sqrt{\frac{n_g}{n}}\|\boldsymbol{\delta}_g^{(t+1)}\|_2 \leq r(\tau)^t \sum_{g=0}^{G} \sqrt{\frac{n_g}{n}}\|\boldsymbol{\beta}_g^*\|_2 + \frac{C(G+1)\sqrt{(2k_w^2+1)k_x^2}}{\sqrt{n}(1-r(\tau))}\left(\max_{g\in[G_+]}\omega(\mathcal{A}_g) + \tau\right)$$

396 *where* $r(\tau) < 1$ *is a constant depending on* $\tau$ *defined in* (SM1.1) *and* $\upsilon, \zeta$, *and* $\sigma$ *are constants.*

397      **Corollary 5.5.** *For enough number of samples, iterations of DS algorithm with step sizes* $\mu_0 =$
398 $\Theta(\frac{1}{n})$ *and* $\mu_g = \Theta(\frac{1}{\sqrt{nn_g}})$ *geometrically converges to the following with high probability:*

399 (5.3)
$$\sum_{g=0}^{G} \sqrt{\frac{n_g}{n}}\|\boldsymbol{\delta}_g^{\infty}\|_2 \leq c\frac{\max_{g\in[G_+]}\omega(\mathcal{A}_g) + \sqrt{\log(G+1)}/\zeta + \theta}{\sqrt{n}(1-r(\tau))}$$

400 *where* $c = C(G+1)\sqrt{(2k_w^2+1)k_x^2}$.

401 It is instructive to compare RHS of (5.3) with that of Theorem 4.1: $\kappa_{\min}$ defined in Theorem Theorem 3.7
402 corresponds to $(1 - r(\tau))$ and the extra $G + 1$ factor corresponds to the sample condition number
403 $\gamma = \max_{g\in[G]} \frac{n}{n_g}$. Therefore, Corollary 5.5 shows that with the number of samples in the order
404 of sample complexity determined in Theorem 3.7 DASHER converges to the statistical error bound
405 determined in Theorem 4.1.

406     **5.1. Proof Sketch of Theorem 5.4.** We want to determine $r(\tau) < 1$ such that $\rho < r(\tau)$ with
407 high probability. Here, we provide a proof sketch using the below probabilistic bounds on constants of
408 Definition 5.1 while ignoring detailed computation of subsequent probabilities in finding $r(\tau)$. The full
409 probabilistic proof is provided in subsection SM1.2. First we need the following lemma to upper bound
410 constants of Definition 5.1:

**Lemma 5.6.** *Consider $a_g \geq 1$ the following upper bounds hold:*

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

$$\eta_g \left( \frac{1}{a_g n_g} \right) \leq \frac{c_g k_x (\omega_g + \tau)}{a_g n_g}, \quad \text{w.p. at least} \quad 1 - \pi_g \exp\left( -\tau^2 \right)$$

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

*where $\omega_g = \omega(\mathcal{A}_g)$ and $\omega_{0g} = \omega(\mathcal{A}_g) + \omega(\mathcal{A}_0)$.*

To keep $\rho < 1$ in the deterministic bound of Theorem 5.2 with the step sizes $\mu_g = \frac{1}{n_g a_g}$ we need to find the number of samples which satisfy the following conditions:

- Condition 1: $\rho_0 (\mu_0) + \sum_{g=1}^{G} \sqrt{\frac{n_g}{n}} \phi_g (\mu_g) < 1$
- Condition 2: $\forall g \in [G] : \rho_g (\mu_g) + \sqrt{\frac{n}{n_g} \frac{\mu_0}{\mu_g}} \phi_g (\mu_g) < 1$

where according to the step sizes determine in the Theorem $a_0 \triangleq (4n \max_{g \in [G]} (1 + c_{0g} \frac{\omega_{0g} + 2\tau}{\sqrt{n_g}})^2)^{-1}$ and $a_g \triangleq (2\sqrt{n/n_g}(1 + c_{0g} \frac{\omega_{0g} + 2\tau}{\sqrt{n_g}}))^{-1}$. Condition 1 requires $\rho_0 + \sum_{g=1}^{G} \sqrt{\frac{n_g}{n}} \phi_g$ to be strictly below 1 which is equivalent to:

$$\rho_0 (\mu_0) + \sum_{g=1}^{G} \sqrt{\frac{n_g}{n}} \phi_g (\mu_g) \leq \frac{1}{2} \left[ \left( 1 - \frac{1}{a_0} \right) + \sqrt{2} c_0 \frac{2\omega_0 + \tau}{a_0 \sqrt{n}} \right] + \frac{1}{2} \sum_{g=1}^{G} \frac{2}{a_g} \sqrt{\frac{n_g}{n}} \left( 1 + c_{0g} \frac{\omega_{0g} + 2\tau}{\sqrt{n_g}} \right)$$

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

$$= \frac{1}{2} \left[ \left( 2 - \frac{1}{a_0} \right) + \sqrt{2} c_0 \frac{2\omega_0 + \tau}{a_0 \sqrt{n}} \right] < 1$$

So Condition 1 reduces to $n > 8 c_0^2 (\omega(\mathcal{A}_0) + \tau)^2$.

Secondly in Condition 2, we want to bound all of $\rho_g + \mu_0 \sqrt{\frac{n}{n_g} \frac{\phi_g}{\mu_g}}$ terms for $\mu_g = \frac{1}{a_g n_g}$ by 1:

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

$$= \frac{1}{2} \left[ \left[ \left( 1 - \frac{1}{a_g} \right) + \sqrt{2} c_g \frac{2\omega_g + \tau}{a_g \sqrt{n_g}} \right] + \frac{2}{a_0} \sqrt{\frac{n_g}{n}} \left( 1 + c_{0g} \frac{\omega_{0g} + 2\tau}{\sqrt{n_g}} \right) \right]$$

$$\leq 1$$

Condition 2 becomes:

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

$$(\text{Substitute } a_g) = 1 + a_g - \frac{4}{a_0} \left( 1 + c_{0g} \frac{\omega_{0g} + 2\tau}{\sqrt{n_g}} \right)^2$$

$$(\text{Substitute } a_0) \leq 1 + a_g$$

(a) $n_g = 60$, $\mathbf{w} = 0$  (b) $n_g = 60$, $\mathbf{w}_1 \neq 0$  (c) $n_g = 150$, $\mathbf{w} = 0$  (d) $n_g = 150$, $\mathbf{w} = 0$
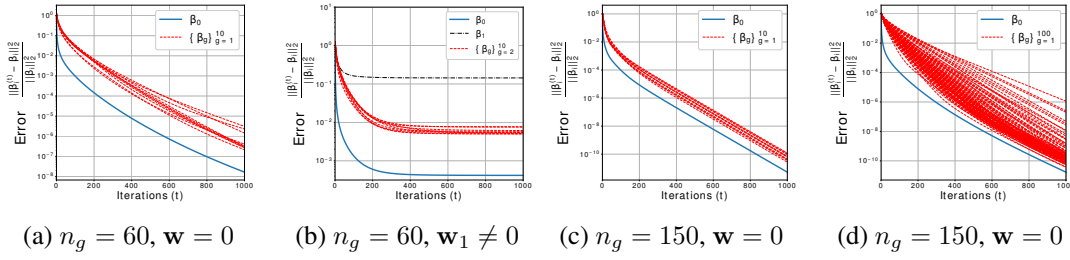
Figure 3: In (a), (b), and (c) experiments $p = 100$, $G = 10$, $\forall g \in [G] : s_g = 10$, and $s_0 = p$. For (d) $p = 1000$, $G = 100$, $\forall g \in [G] : s_g = 10$, and $s_0 = 100$. (a) Noiseless fast convergence. (b) Noise on the first group does not impact other groups as much. (c) Increasing sample size improves rate of convergence. (d) DASHER convergences fast even with a large number of groups $G = 100$.

435  So the sample complexity should be $\sqrt{n_g} > \frac{\sqrt{2}c_g(2\omega_g+2\tau)}{1+a_g}$ and since $a_g > 1$, the final per group sample
436  complexity should be $n_g > 8c_g(\omega(\mathcal{A}_g) + \tau)^2$. ■

437  **6. Experiments on Synthetic Data.** We considered sparsity based simulations with varying
438  $G$ and sparsity levels. In our first set of simulations, we set $p = 100$, $G = 10$ and sparsity of the
439  individual parameters to be $s = 10$. We generated a dense $\boldsymbol{\beta}_0$ with $\|\boldsymbol{\beta}_0\| = p$ and did not impose any
440  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 $\boldsymbol{\beta}_g$
441  are generated with $\mathcal{N}(0, 1)$ and nonzero supports are picked uniformly at random. Inspired from our
442  theoretical step size choices, in all experiments, we used simplified learning rates of $\frac{1}{n}$ for $\boldsymbol{\beta}_0$ and $\frac{1}{\sqrt{nn_g}}$
443  for $\boldsymbol{\beta}_g$, $g \in [G]$. Observe that, cones of the individual parameters intersect with that of $\boldsymbol{\beta}_0$ hence this
444  setup actually violates DASHIN (which requires an arbitrarily small constant fraction of groups to be
445  non-intersecting). Our intuition is that the individual parameters are mostly incoherent with each other
446  and the existence of a nonzero perturbation over $\boldsymbol{\beta}_g$'s that keeps all measurements intact is unlikely.
447  Remarkably, experimental results still show successful learning of all parameters from small amount
448  of samples. We picked $n_g = 60$ for each group. Hence, in total, we have $11p = 1100$ unknowns,
449  $200 = G \times 10 + 100$ degrees of freedom and $G \times 60 = 600$ samples. In all figures, we study the
450  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
451  shows the estimation performance as a function of iteration number $t$. While each group might behave
452  slightly different, we do observe that all parameters are linear converging to ground truth.
453  In Figure 3b, we test the noise robustness of our algorithm. We add a $\mathcal{N}(0, 1)$ noise to the $n_1 = 60$
454  measurements of the first group *only*. The other groups are left untouched. While all parameters suffer
455  nonzero estimation error, we observe that, the global parameter $\boldsymbol{\beta}_0$ and noise-free groups $\{\boldsymbol{\beta}_g\}_{g=2}^G$ have
456  substantially less estimation error. This implies that noise in one group mostly affects itself rather than
457  the global estimation. In Figure 3c, we increased the sample size to $n_g = 150$ per group. We observe
458  that, in comparison to Figure 3a, rate of convergence receives a boost from the additional samples as
459  predicted by our theory.
460  Finally, Figure 3d considers a very high-dimensional problem where $p = 1000$, $G = 100$, individual
461  parameters are 10 sparse, $\boldsymbol{\beta}_0$ is 100 sparse and $n_g = 150$. The total degrees of freedom is 1100, number
462  of unknowns are 101000 and total number of datapoints are $150 \times 100 = 15000$. While individual

463  parameters have substantial variation in terms of convergence rate, at the end of 1000 iteration, all
464  parameters have relative reconstruction error below $10^{-6}$.

465                                              **REFERENCES**

466  [1]  A. ASIAEE, S. OYMAK, K. R. COOMBES, AND A. BANERJEE, *High dimensional data enrichment: Interpretable,*
467         *fast, and data-efficient*, arXiv preprint arXiv:1806.04047, (2018).
468  [2]  A. ASIAEE, S. OYMAK, K. R. COOMBES, AND A. BANERJEE, *Data enrichment: Multi-task learning in high*
469         *dimension with theoretical guarantees*, in Adaptive and Multitask Learning Workshop at ICML, 2019.
470  [3]  F. BACH, R. JENATTON, J. MAIRAL, G. OBOZINSKI, ET AL., *Optimization with sparsity-inducing penalties*,
471         Foundations and Trends® in Machine Learning, 4 (2012), pp. 1–106.
472  [4]  A. BANERJEE, S. CHEN, F. FAZAYELI, AND V. SIVAKUMAR, *Estimation with Norm Regularization*, in Advances in
473         Neural Information Processing Systems, 2014, pp. 1556–1564.
474  [5]  J. BARRETINA, G. CAPONIGRO, N. STRANSKY, K. VENKATESAN, A. A. MARGOLIN, S. KIM, C. J. WILSON,
475         J. LEHÁR, G. V. KRYUKOV, D. SONKIN, ET AL., *The cancer cell line encyclopedia enables predictive modelling*
476         *of anticancer drug sensitivity*, Nature, 483 (2012), p. 603.
477  [6]  P. J. BICKEL, Y. RITOV, A. B. TSYBAKOV, ET AL., *Simultaneous analysis of lasso and dantzig selector*, The Annals
478         of Statistics, 37 (2009), pp. 1705–1732.
479  [7]  T. BLUMENSATH AND M. E. DAVIES, *Iterative hard thresholding for compressed sensing*, Applied and computational
480         harmonic analysis, 27 (2009), pp. 265–274.
481  [8]  S. BOUCHERON, G. LUGOSI, AND P. MASSART, *Concentration Inequalities: A Nonasymptotic Theory of Indepen-*
482         *dence*, Oxford University Press, 2013.
483  [9]  P. T. BOUFOUNOS AND R. G. BARANIUK, *1-bit compressive sensing*, in Information Sciences and Systems, 2008.
484         CISS 2008. 42nd Annual Conference on, IEEE, 2008, pp. 16–21.
485  [10] E. CANDES, T. TAO, ET AL., *The dantzig selector: Statistical estimation when p is much larger than n*, The Annals of
486         Statistics, 35 (2007), pp. 2313–2351.
487  [11] E. J. CANDÈS AND B. RECHT, *Exact matrix completion via convex optimization*, Foundations of Computational
488         mathematics, 9 (2009), p. 717.
489  [12] E. J. CANDÈS, J. ROMBERG, AND T. TAO, *Robust uncertainty principles: Exact signal reconstruction from highly*
490         *incomplete frequency information*, IEEE Transactions on information theory, 52 (2006), pp. 489–509.
491  [13] E. J. CANDÈS AND T. TAO, *The power of convex relaxation: Near-optimal matrix completion*, IEEE Transactions on
492         Information Theory, 56 (2010), pp. 2053–2080.
493  [14] V. CHANDRASEKARAN, B. RECHT, P. A. PARRILO, AND A. S. WILLSKY, *The convex geometry of linear inverse*
494         *problems*, Foundations of Computational Mathematics, 12 (2012), pp. 805–849.
495  [15] S. CHATTERJEE, S. CHEN, AND A. BANERJEE, *Generalized dantzig selector: Application to the k-support norm*, in
496         Advances in Neural Information Processing Systems, 2014, pp. 1934–1942.
497  [16] A. CHEN, A. B. OWEN, AND M. SHI, *Data enriched linear regression*, Electronic journal of statistics, 9 (2015),
498         pp. 1078–1112.
499  [17] J. CHEN, J. LIU, AND J. YE, *Learning incoherent sparse and Low-Rank patterns from multiple tasks*, ACM transactions
500         on knowledge discovery from data, 5 (2012), p. 22.
501  [18] F. DONDELINGER, S. MUKHERJEE, AND ALZHEIMER'S DISEASE NEUROIMAGING INITIATIVE, *The joint lasso:*
502         *high-dimensional regression for group structured data*, Biostatistics, (2018).
503  [19] D. L. DONOHO, *Compressed sensing*, IEEE Transactions on information theory, 52 (2006), pp. 1289–1306.
504  [20] J. FRIEDMAN, T. HASTIE, AND R. TIBSHIRANI, *Sparse inverse covariance estimation with the graphical lasso*,
505         Biostatistics, 9 (2008), pp. 432–441.
506  [21] S. M. GROSS AND R. TIBSHIRANI, *Data shared lasso: A novel tool to discover uplift*, Computational Statistics &
507         Data Analysis, 101 (2016), pp. 226–235.
508  [22] Q. GU AND A. BANERJEE, *High dimensional structured superposition models*, in Advances In Neural Information
509         Processing Systems, 2016, pp. 3684–3692.
510  [23] F. IORIO, T. A. KNIJNENBURG, D. J. VIS, G. R. BIGNELL, M. P. MENDEN, M. SCHUBERT, N. ABEN,
511         E. GONCALVES, S. BARTHORPE, H. LIGHTFOOT, ET AL., *A landscape of pharmacogenomic interactions*
512         *in cancer*, Cell, 166 (2016), pp. 740–754.
513  [24] 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.

[25] 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.

[26] S. KAKADE, O. SHAMIR, K. SINDHARAN, AND A. TEWARI, *Learning exponential families in High-Dimensions: Strong convexity and sparsity*, in Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, Y. W. Teh and M. Titterington, eds., vol. 9 of Proceedings of Machine Learning Research, Chia Laguna Resort, Sardinia, Italy, 2010, PMLR, pp. 381–388.

[27] M. B. MCCOY AND J. A. TROPP, *The achievable performance of convex demixing*, (2013), https://arxiv.org/abs/1309.7478.

[28] S. MENDELSON, *Learning without concentration*, Journal of the ACM, 62 (2015), pp. 21:1–21:25.

[29] S. NEGAHBAN AND M. J. WAINWRIGHT, *Restricted strong convexity and weighted matrix completion: Optimal bounds with noise*, Journal of Machine Learning Research, 13 (2012), pp. 1665–1697.

[30] S. NEGAHBAN, B. YU, M. J. WAINWRIGHT, AND P. K. RAVIKUMAR, *A unified framework for high-dimensional analysis of m-estimators with decomposable regularizers*, in Advances in Neural Information Processing Systems, 2009, pp. 1348–1356.

[31] S. N. NEGAHBAN, P. RAVIKUMAR, M. J. WAINWRIGHT, AND B. YU, *A Unified Framework for High-Dimensional Analysis of $M$-Estimators with Decomposable Regularizers*, Statistical Science, 27 (2012), pp. 538–557.

[32] E. OLLIER AND V. VIALLON, *Joint estimation of $K$ related regression models with simple $L_1$-norm penalties*, (2014), https://arxiv.org/abs/1411.1594.

[33] E. OLLIER AND V. VIALLON, *Regression modeling on stratified data with the lasso*, (2015), https://arxiv.org/abs/1508.05476.

[34] Y. PLAN AND R. VERSHYNIN, *Robust 1-bit compressed sensing and sparse logistic regression: A convex programming approach*, IEEE transactions on information theory / Professional Technical Group on Information Theory, 59 (2013), pp. 482–494.

[35] Y. PLAN AND R. VERSHYNIN, *The generalized lasso with Non-Linear observations*, IEEE transactions on information theory / Professional Technical Group on Information Theory, 62 (2016), pp. 1528–1537.

[36] Y. PLAN, R. VERSHYNIN, AND E. YUDOVINA, *High-dimensional estimation with geometric constraints*, Information and Inference: A Journal of the IMA, 6 (2017), pp. 1–40.

[37] G. RASKUTTI, M. J. WAINWRIGHT, AND B. YU, *Restricted eigenvalue properties for correlated gaussian designs*, Journal of Machine Learning Research, 11 (2010), pp. 2241–2259.

[38] M. RUDELSON AND S. ZHOU, *Reconstruction from anisotropic random measurements*, IEEE Transactions on Information Theory, 59 (2013), pp. 3434–3447.

[39] R. TIBSHIRANI, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society. Series B (Methodological), (1996), pp. 267–288.

[40] J. A. TROPP, *Convex recovery of a structured signal from independent random linear measurements*, in Sampling Theory, a Renaissance, Springer, 2015, pp. 67–101.

[41] R. VERSHYNIN, *Introduction to the non-asymptotic analysis of random matrices*, in Compressed Sensing, Cambridge University Press, Cambridge, 2012, pp. 210–268.

[42] R. VERSHYNIN, *High-dimensional probability: An introduction with applications in data science*, vol. 47, Cambridge University Press, 2018.

[43] E. YANG AND P. K. RAVIKUMAR, *Dirty statistical models*, in Advances in Neural Information Processing Systems 26, C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, eds., Curran Associates, Inc., 2013, pp. 611–619.

[44] Z. YANG, Z. WANG, H. LIU, Y. ELDAR, AND T. ZHANG, *Sparse nonlinear regression: Parameter estimation under nonconvexity*, in Proceedings of The 33rd International Conference on Machine Learning, M. F. Balcan and K. Q. Weinberger, eds., vol. 48 of Proceedings of Machine Learning Research, New York, New York, USA, 2016, PMLR, pp. 2472–2481.

[45] Y. ZHANG AND Q. YANG, *A survey on Multi-Task learning*, (2017), https://arxiv.org/abs/1707.08114.