# Linear Regression with Distributed Learning: A Generalization Error Perspective

Martin Hellkvist, Ayça Özçelikkale, Anders Ahlén

*Abstract*—**Distributed learning provides an attractive framework for scaling the learning task by sharing the computational load over multiple nodes in a network. Here, we investigate the performance of distributed learning for large-scale linear regression where the model parameters, i.e., the unknowns, are distributed over the network. We adopt a statistical learning approach. In contrast to works that focus on the performance on the training data, we focus on the generalization error, i.e., the performance on unseen data. We provide high-probability bounds on the generalization error for both isotropic and correlated Gaussian data as well as sub-gaussian data. These results reveal the dependence of the generalization performance on the partitioning of the model over the network. In particular, our results show that the generalization error of the distributed solution can be substantially higher than that of the centralized solution even when the error on the training data is at the same level for both the centralized and distributed approaches. Our numerical results illustrate the performance with both real-world image data as well as synthetic data.**

*Index Terms*—**Distributed estimation, distributed optimization, supervised learning, generalization error, networked systems.**

## I. INTRODUCTION

Distributed learning provides a framework for sharing the computational burden of large-scale learning tasks over multiple nodes while addressing growing concerns related to security and data privacy [1], [2]. Accordingly, the field of distributed learning is progressing rapidly due to the increasing need and interest from both industry and academia, with applications ranging from edge computing [3], [4] to large-scale machine learning [5]–[7]. In this article, we consider distributed learning from the point of view of generalization error and contribute to the field by highlighting and characterizing potential pitfalls, and providing guidelines for best practice.

In particular, we consider the statistical learning problem where a set of training data $\{(y_i, \boldsymbol{a}_i)\}_{i=1}^n$ from a certain distribution is used to train a model, i.e., estimate parameters in a specified model structure, so that it correctly predicts the output $y_i \in \mathbb{R}$ given the corresponding input $\boldsymbol{a}_i \in \mathbb{R}^{p \times 1}$. The performance of the trained model is often measured by its *training error*, i.e., the error that the model makes over the training data, and more importantly its *generalization error*, i.e., the error that the model makes when estimating

$y$ using $\boldsymbol{a}$ when a new pair $(y, \boldsymbol{a})$ comes from the same distribution as the training data. The generalization error is an inherent part of statistical learning frameworks, where it is innately embedded in the expected error values, and the expectation is taken with respect to the signal model. The generalization error has been thoroughly studied for different centralized approaches, for instance in minimum mean-square error estimation frameworks [8]. Recently, the dependence of the generalization error on the number of model parameters and the training sample size has been investigated for a range of models, such as neural networks and decision trees, and the "double descent" risk curve has been proposed [9]. The generalization error associated with the least-squares estimate under isotropic Gaussian data and Fourier features with partial models [10], [11], the effect of regularization under data correlation [12], [13], as well as sub-gaussian regressor distributions [14], [15] have been presented. These works emphasize trade-offs between model complexity and training sample size in a centralized learning setting, particularly in overparametrized scenarios.

The growing need for distributed learning has lead to the development of several methods, e.g., primal-dual methods [16], [17] and dual decomposition methods [18], [19], where the alternating direction method of multipliers (ADMM) [20] stands out as one of the most extensively studied algorithms. Accordingly, different aspects of distributed learning methods have been explored, including privacy protecting methods [2], time-varying constraints [21], adaptive network architectures [22], and communication efficient methods such as the quantized stochastic gradient descent [23], as well as novel metrics for communication efficiency [24]. Trade-offs between computation and communication [25] has been explored as well. In the case of generalization error, a significant part of the existing work for distributed learning is performed under the mean-square error estimation framework, including Kalman filtering [26], [27], least-mean squares [28], [29] and the affine projection algorithm [30]. A characterization of the generalization error in the case of linear discriminant analysis is presented in [31]. The average behaviour of the generalization error for regression is presented in [32] under isotropic Gaussian data.

Despite this vast interest in distributed learning, this line of work typically assumes that it is the sensor readings that are distributed over the network [29], in contrast to the scenario where the model unknowns are distributed over the network [27], [33]. We address this gap by providing high probability bounds on the generalization error in a distributed linear regression problem under a broad family of training

data distributions. The setting with distributed unknowns is particularly suited to the problems with large number of unknowns [25], such as neural networks [34]. Motivated by the recent results on overparameterization in linear regression [12]–[15], we pay special attention to the overparameterized setting where the number of unknowns governed by each node is larger than the number of observations.

We consider the influential distributed learning algorithm CoCoA [25], developed from its predecessors CoCoA-v1 [35] and CoCoA$^+$ [36]. CoCoA is applicable to convex optimization problems, and allows the nodes to use any local solver of their choice for their local subproblems, enabling the usage of solvers with variable accuracy and a flexible trade-off between computation and communication [25]. In [25], the convergence rate of CoCoA was quantified in terms of the convexity properties of the optimization problem and accuracy of local solvers. In contrast to the work in [25], we focus on the generalization error of CoCoA and the effect of different data partitioning schemes over the nodes.

In this article, we show that the generalization error depends heavily on the partitioning of the model parameters among the nodes. In particular, we have the following main contributions: We provide bounds on the generalization error that hold with high probability for both isotropic Gaussian as well as correlated Gaussian data. Furthermore, for block-correlated and underparameterized local problems with general covariance structure, we generalize these results to sub-gaussian data, which include the Bernoulli and uniform distributions as special cases. For the isotropic Gaussian case, we compare these probabilistic results with analytical results on the average behaviour [32], which we also extend to the setting with noisy training data in this article. We note that the presented results cover a wide set of distributions, compared to the scope of [32], which is limited to the isotropic Gaussian distribution. Our numerical results illustrate the generalization error performance with both synthetic data from these distributions and real-world image data [37].

Our results highlight a typically overlooked relationship between the training and generalization error in distributed learning. These findings illustrate that distributed learning schemes can significantly amplify the gap between the training error and the generalization error. More precisely, a distributed solution with a training error that is on the same level as that of the centralized solution is not guaranteed to have a generalization error that is as low as that of the centralized solution.

The rest of the paper is organized as follows: Section II and Section III present the problem formulation and the distributed solution approach, respectively. Section IV provides preliminary results on the generalization error. In Section V, VI and VII, we present the results for the isotropic Gaussian, correlated Gaussian, and the sub-gaussian settings, respectively. The numerical results are presented in Section VIII. We present further discussions of our results in Section IX and conclude the article in Section X.

**Notation:** We denote the Moore-Penrose pseudoinverse and the transpose of a matrix $\boldsymbol{A}$ as $\boldsymbol{A}^+$ and $\boldsymbol{A}^\mathrm{T}$, respectively. The $p \times p$ identity matrix is denoted as $\boldsymbol{I}_p$. The positive semi-
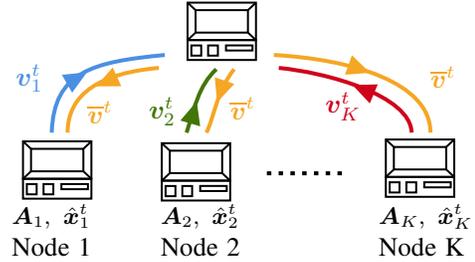


Fig. 1: Distributed learning with CoCoA.

definite (p.s.d.) partial ordering for real symmetric matrices is denoted by $\succeq$. We use $\|\cdot\|$ to denote either the spectral norm or the Euclidean norm, depending on whether the argument is matrix- or vector valued. Throughout the paper, we often partition vectors by blocks of their entries, and matrices either by their blocks of columns or rows. For instance, the column-wise partitioning of a matrix $\boldsymbol{A} \in \mathbb{R}^{n \times p}$ into $K$ blocks is given by $\boldsymbol{A} = [\boldsymbol{A}_1, \cdots, \boldsymbol{A}_K]$, with $\boldsymbol{A}_k \in \mathbb{R}^{n \times p_k}$. The row-wise partitioning of a vector $\boldsymbol{x} \in \mathbb{R}^{p \times 1}$ into $K$ blocks $\boldsymbol{x}_k \in \mathbb{R}^{p_k \times 1}$ is given by $\boldsymbol{x} = [\boldsymbol{x}_1; \cdots; \boldsymbol{x}_K]$, where the semicolon denotes row-wise separation. We use $\sigma_{\max}(\cdot)$, $\sigma_{\min}(\cdot)$ to denote the largest and smallest singular values of a matrix, and $\lambda_{\max}(\cdot)$, $\lambda_{\min}(\cdot)$ to denote the largest and smallest eigenvalues. The notation $(\cdot)_+$ is used as a short-hand for $\max\{0, \cdot\}$. In expressions such as $(\cdot)_+^2$, the $\max$ function takes precedence over the square, i.e., $(\cdot)_+^2 = ((\cdot)_+)^2$.

## II. PROBLEM STATEMENT

We focus on the linear model

$$y_i = \boldsymbol{a}_i^\mathrm{T} \boldsymbol{x} + w_i, \tag{1}$$

where $y_i \in \mathbb{R}$ is the $i^\mathrm{th}$ observation, $\boldsymbol{a}_i \in \mathbb{R}^{p \times 1}$ is the $i^\mathrm{th}$ regressor, $w_i \in \mathbb{R}$ is the corresponding unknown disturbance, and $\boldsymbol{x} \in \mathbb{R}^{p \times 1}$ is the vector of unknown model parameters. We consider the problem of estimating $\boldsymbol{x}$ given $n$ pairs of observations and regressors, i.e., the training dataset $\{(y_i, \boldsymbol{a}_i)\}_{i=1}^n$ by minimizing the following regularized cost function:

$$\min_{\boldsymbol{x} \in \mathbb{R}^{p \times 1}} \frac{1}{2}\|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|^2 + \frac{\lambda}{2}\|\boldsymbol{x}\|^2, \tag{2}$$

where $\boldsymbol{A} \in \mathbb{R}^{n \times p}$ is the regressor matrix whose $i^\mathrm{th}$ row is given by $\boldsymbol{a}_i^\mathrm{T} \in \mathbb{R}^{1 \times p}$, $\boldsymbol{y} \in \mathbb{R}^{n \times 1}$ is the vector of observations $y_i$, and $\lambda \geq 0$ is a regularization parameter.

We consider the setting where the regressors $\boldsymbol{a}_i^\mathrm{T}$ are independent and identically distributed (i.i.d.) zero-mean random vectors with a given distribution $\mathcal{D}(\boldsymbol{\Sigma})$, with the covariance matrix $\boldsymbol{\Sigma} = \mathbb{E}_{\boldsymbol{a}_i}[\boldsymbol{a}_i \boldsymbol{a}_i^\mathrm{T}] \in \mathbb{R}^{p \times p}$. Under this regressor model, we investigate the generalization error of the solution to (2) found by the distributed solver CoCoA [25]. In order to simplify the theoretical analysis, we mainly consider the unregularized and noise-free setting, i.e., with $\lambda = 0$ and $\boldsymbol{w} = 0$. Under these simplifications, we derive bounds illustrating how the generalization error depends on the partitioning of the model over the network. In order to provide background for these results, we consider the more general case with $\boldsymbol{w} \neq 0$ in Sections II – IV together with a discussion on the case with $\lambda > 0$. In Section IV and Section V-A, we provide insights

---

**Algorithm 1:** Implementation of CoCoA [25] for (2).

---

**1 Input**: Data matrix $\boldsymbol{A}$ distributed column-wise according to partitioning $\{p_1, \cdots, p_k\}$. Observations $\boldsymbol{y}$. Regularization parameter $\lambda$, aggregation parameter $\bar{\varphi} \in (0,1]$ and subproblem parameter $\bar{\sigma}$.

**2 Initialize**: $\hat{\boldsymbol{x}}^0 = 0 \in \mathbb{R}^{p \times 1}$, $\boldsymbol{v}_k^0 = 0 \in \mathbb{R}^{p \times 1}$ $\forall k$.

**3 for** $t = 0, 1, \ldots, T$ **do**

**4**    $\bar{\boldsymbol{v}}^t = \frac{1}{K} \sum_{k=1}^{K} \boldsymbol{v}_k^t$

**5**    **for** $k \in \{1, 2, \ldots, K\}$ **do**

**6**      $\boldsymbol{c}_k^t = \lambda \hat{\boldsymbol{x}}_k^t - \boldsymbol{A}_k^{\mathrm{T}}(\boldsymbol{y} - \bar{\boldsymbol{v}}^t)$

**7**      $\Delta \boldsymbol{x}_k^t = -(\bar{\sigma} \boldsymbol{A}_k^{\mathrm{T}} \boldsymbol{A}_k + \lambda \boldsymbol{I}_{p_k})^{+} \boldsymbol{c}_k^t$

**8**      $\hat{\boldsymbol{x}}_k^{t+1} = \hat{\boldsymbol{x}}_k^t + \bar{\varphi} \Delta \boldsymbol{x}_k^t$

**9**      $\boldsymbol{v}_k^{t+1} = \bar{\boldsymbol{v}}^t + \bar{\varphi} K \boldsymbol{A}_k \Delta \boldsymbol{x}_k^t$

---

about why the training noise does not necessarily weaken the dependence of generalization error on partitioning. For the regularized case, i.e., with $\lambda > 0$, and for the case with non-zero noise, we provide numerical results which illustrate how the same heavy dependence on partitioning occurs for $\lambda > 0$ before convergence; and for $\boldsymbol{w} \neq 0$ even after convergence, see Section VIII. In the remainder of this section, we define the generalization error. We provide details for CoCoA in Section III.

Let $\hat{\boldsymbol{x}}$ be an estimate of $\boldsymbol{x}$ found using a given set of training data $\{(y_i, \boldsymbol{a}_i)\}_{i=1}^n$, where $\boldsymbol{a}_i \sim \mathcal{D}(\boldsymbol{\Sigma})$ and $y_i = \boldsymbol{a}_i^{\mathrm{T}} \boldsymbol{x} + w_i$. Let $(y, \boldsymbol{a})$ be a new input-output pair with $\boldsymbol{a} \sim \mathcal{D}(\boldsymbol{\Sigma})$ and $y = \boldsymbol{a}^{\mathrm{T}} \boldsymbol{x} + w$. Then, the generalization error is given by

$$\kappa(\hat{\boldsymbol{x}}) = \mathbb{E}_{\boldsymbol{a}}[(\boldsymbol{a}^{\mathrm{T}} \boldsymbol{x} - \boldsymbol{a}^{\mathrm{T}} \hat{\boldsymbol{x}})^2] \quad (3)$$

$$= (\boldsymbol{x} - \hat{\boldsymbol{x}})^{\mathrm{T}} \mathbb{E}_{\boldsymbol{a}}[\boldsymbol{a}\boldsymbol{a}^{\mathrm{T}}](\boldsymbol{x} - \hat{\boldsymbol{x}}) \quad (4)$$

$$= (\boldsymbol{x} - \hat{\boldsymbol{x}})^{\mathrm{T}} \boldsymbol{\Sigma} (\boldsymbol{x} - \hat{\boldsymbol{x}}), \quad (5)$$

where we have used that $\hat{\boldsymbol{x}}$ is a fixed estimate under the given training data. The notation $\mathbb{E}_{\boldsymbol{a}}[\cdot]$ is used to emphasize that the expectation is over the previously unseen regressor $\boldsymbol{a}$. One may alternatively consider the prediction error in $y$ instead of $\boldsymbol{a}^{\mathrm{T}} \boldsymbol{x}$:

$$\mathbb{E}_{\boldsymbol{a},w}[(y - \boldsymbol{a}^{\mathrm{T}} \hat{\boldsymbol{x}})^2] = \kappa(\hat{\boldsymbol{x}}) + \sigma_w^2, \quad (6)$$

where the noise $w$ in the test data is assumed to be zero-mean with variance $\sigma_w^2$ and statistically independent with the regressor $\boldsymbol{a}$. Since the noise in test data just gives an additive, irreducible term, we focus directly on $\kappa(\hat{\boldsymbol{x}})$ in our technical development. We are interested in the behaviour of the generalization error $\kappa(\hat{\boldsymbol{x}}) \in \mathbb{R}$ with respect to the distribution of the training data, i.e., $\boldsymbol{A}$, and the partitioning of the data over the nodes.

In the centralized case, a solution to (2) is found as

$$\hat{\boldsymbol{x}}_C = (\boldsymbol{A}^{\mathrm{T}} \boldsymbol{A} + \lambda \boldsymbol{I}_p)^{+} \boldsymbol{A}^{\mathrm{T}} \boldsymbol{y}. \quad (7)$$

In general, with $\lambda = 0$, there can exist multiple solutions to (2). With the Moore-Penrose pseudoinverse, the solution with the minimum Euclidean norm is obtained.

## III. Distributed Solution Approach

We now discuss how to obtain a solution $\hat{\boldsymbol{x}}$ for (2) using the distributed solution approach CoCoA [25], see Figure 1 and Algorithm 1. Here, mutually exclusive subsets of the

$p$ unknown parameters in $\boldsymbol{x}$ and the associated subset of columns in $\boldsymbol{A} \in \mathbb{R}^{n \times p}$ are distributed over $K$ nodes, $K \leq p$. Hence, node $k$ governs the learning of $p_k$ variables, denoted by $\boldsymbol{x}_k \in \mathbb{R}^{p_k \times 1}$, where $\sum_{k=1}^K p_k = p$. We denote the part of $\boldsymbol{A}$ available at node $k$ as $\boldsymbol{A}_k \in \mathbb{R}^{n \times p_k}$. All nodes have access to the vector of observations $\boldsymbol{y} \in \mathbb{R}^{n \times 1}$. Using this partitioning, $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x} + \boldsymbol{w}$ can be expressed as

$$\boldsymbol{y} = [\boldsymbol{A}_1, \cdots, \boldsymbol{A}_K] \begin{bmatrix} \boldsymbol{x}_1 \\ \vdots \\ \boldsymbol{x}_K \end{bmatrix} + \boldsymbol{w} = \sum_{k=1}^K \boldsymbol{A}_k \boldsymbol{x}_k + \boldsymbol{w}, \quad (8)$$

Note that the submatrices $\boldsymbol{A}_k$'s and the observation vector $\boldsymbol{y}$ are fixed over all iterations.

Node $k$ forms an estimate of $\boldsymbol{x}_k$ using $\boldsymbol{y}$, $\boldsymbol{A}_k$ and a centrally computed variable $\bar{\boldsymbol{v}}^t \in \mathbb{R}^{n \times 1}$. Let $\hat{\boldsymbol{x}}_k^t \in \mathbb{R}^{p_k \times 1}$ denote the estimate of $\boldsymbol{x}_k$ at node $k$ and iteration $t$. Accordingly, let $\hat{\boldsymbol{x}}^t = [\hat{\boldsymbol{x}}_1^t; \ldots; \hat{\boldsymbol{x}}_k^t] \in \mathbb{R}^{p \times 1}$ denote the estimate of $\boldsymbol{x}$ at iteration $t$. At iteration $t$, node $k$ receives the centrally computed variable $\bar{\boldsymbol{v}}^t$, which it uses to compute $\Delta \boldsymbol{x}_k^t \in \mathbb{R}^{p_k \times 1}$ (Line 6-7, Alg. 1), i.e., the update for $\hat{\boldsymbol{x}}_k^t$ (Line 8). The node keeps track of its contribution for estimating $\boldsymbol{y}$ by computing the local estimate $\boldsymbol{v}_k^{t+1} \in \mathbb{R}^{n \times 1}$, using $\bar{\boldsymbol{v}}^t$ and $\Delta \boldsymbol{x}_k^t$ (Line 9). Then, the variable $\boldsymbol{v}_k^{t+1}$ is sent to a central node to create $\bar{\boldsymbol{v}}^{t+1}$ (Line 4). The central node then sends $\bar{\boldsymbol{v}}^{t+1}$ to the nodes and the next iteration begins.

We now explain how node $k$ finds the update $\Delta \boldsymbol{x}_k^t$. To find $\Delta \boldsymbol{x}_k^t$, CoCoA solves the following convex minimization problem at each node [25]:

$$\min_{\Delta \boldsymbol{x}_k^t} \frac{1}{K} f(\bar{\boldsymbol{v}}^t) + \nabla_{\bar{\boldsymbol{v}}^t} f(\bar{\boldsymbol{v}}^t)^{\mathrm{T}} \boldsymbol{A}_k \Delta \boldsymbol{x}_k^t \\ + \frac{\bar{\sigma}}{2\tau} \|\boldsymbol{A}_k \Delta \boldsymbol{x}_k^t\|^2 + \frac{\lambda}{2} \|\hat{\boldsymbol{x}}_k^t + \Delta \boldsymbol{x}_k^t\|^2, \quad (9)$$

where $f(\bar{\boldsymbol{v}}^t) = \frac{1}{2} \|\boldsymbol{y} - \bar{\boldsymbol{v}}^t\|^2$ is the first term of the objective function in (2), evaluated at $\bar{\boldsymbol{v}}^t$. Note that $\bar{\boldsymbol{v}}^t = \boldsymbol{A}\hat{\boldsymbol{x}}^t$ by Algorithm 1. The first two terms of (9) comes from the linearization of $f(\cdot)$ around the current value of $\bar{\boldsymbol{v}}^t$, and the third term $\frac{\bar{\sigma}}{2\tau} \|\boldsymbol{A}\Delta \boldsymbol{x}_k^t\|^2$ penalizes large changes in $\bar{\boldsymbol{v}}^t = \boldsymbol{A}\hat{\boldsymbol{x}}^t$. The last term $\frac{\lambda}{2} \|\hat{\boldsymbol{x}}_k^t + \Delta \boldsymbol{x}_k^t\|^2$ corresponds to the local component of the regularization term in (2) evaluated at $\hat{\boldsymbol{x}}_k^t + \Delta \boldsymbol{x}_k^t$.

The smoothness parameter for $f(\cdot)$ is $\tau = 1$ [25]. Only keeping the terms that depend on $\Delta \boldsymbol{x}_k^t$ reveals that (9) can be equivalently solved by

$$\min_{\Delta \boldsymbol{x}_k^t} (\Delta \boldsymbol{x}_k^t)^{\mathrm{T}} (\frac{\bar{\sigma}}{2} \boldsymbol{A}_k^{\mathrm{T}} \boldsymbol{A}_k + \frac{\lambda}{2} \boldsymbol{I}_{p_k}) \Delta \boldsymbol{x}_k^t \\ + (\lambda \hat{\boldsymbol{x}}_k^t - \boldsymbol{A}_k^{\mathrm{T}} (\boldsymbol{y} - \bar{\boldsymbol{v}}^t))^{\mathrm{T}} \Delta \boldsymbol{x}_k^t. \quad (10)$$

Taking the derivative with respect to $\Delta \boldsymbol{x}_k^t$ and setting it to zero, we obtain

$$(\bar{\sigma} \boldsymbol{A}_k^{\mathrm{T}} \boldsymbol{A}_k + \lambda \boldsymbol{I}_{p_k}) \Delta \boldsymbol{x}_k^t = -(\lambda \hat{\boldsymbol{x}}_k^t - \boldsymbol{A}_k^{\mathrm{T}} (\boldsymbol{y} - \bar{\boldsymbol{v}}^t)). \quad (11)$$

With $\lambda = 0$, the existence of a matrix inverse is not guaranteed, in general. Hence, the local solvers use Moore-Penrose pseudoinverse to solve (11) to obtain

$$\Delta \boldsymbol{x}_k^t = -(\bar{\sigma} \boldsymbol{A}_k^{\mathrm{T}} \boldsymbol{A}_k + \lambda \boldsymbol{I}_{p_k})^{+} (\lambda \hat{\boldsymbol{x}}_k^t - \boldsymbol{A}_k^{\mathrm{T}} (\boldsymbol{y} - \bar{\boldsymbol{v}}^t)). \quad (12)$$

The resulting algorithm for estimating $\boldsymbol{x}$ iteratively is presented in Algorithm 1.

Similar to dual decomposition methods [18], [19] and in particular ADMM [20], CoCoA encourages a consensus over nodes by utilizing Lagrangian duality. Although it is inherently

4

connected to ADMM [20], CoCoA utilizes a more simple update for $\hat{x}$ and allows approximate proximal steps, see for instance [25, Eqn. (12), Eqn. (15)].

## IV. GENERALIZATION ERROR WITH CoCoA

We are interested in the behaviour of the generalization error in (5) with respect to different partitioning schemes $\{p_1, \cdots, p_K\}$, as well as different distributions of the data. For the rest of the article, we consider the case with $\lambda = 0$, except for the numerical experiments in Section VIII. We set $\bar{\sigma} = \bar{\varphi}K$, as it is considered a safe choice in terms of convergence [25, Sec. 3.1], and provide additional experiments with other values in Section VIII-G. As shown in Lemma 1 of [32], the iterations of Algorithm 1 can be expressed as

$$\hat{x}^{t+1} = B\hat{x}^t + \tfrac{1}{K}\bar{A}y, \qquad (13)$$

where $B \in \mathbb{R}^{p\times p}$ and $\bar{A} \in \mathbb{R}^{p\times n}$ are given by

$$B = \left(I_p - \tfrac{1}{K}\bar{A}A\right), \qquad (14)$$

and

$$\bar{A} = \begin{bmatrix} A_1^+ \\ A_2^+ \\ \vdots \\ A_K^+ \end{bmatrix}. \qquad (15)$$

Note that, under $\bar{\sigma} = \bar{\varphi}K$, $\bar{\varphi}$ and $\bar{\sigma}$ enters into (13) as $\frac{\bar{\varphi}}{\bar{\sigma}} = \frac{1}{K}$, hence the solutions $\hat{x}^t$ are independent of the particular values of $\bar{\varphi}$ and $\bar{\sigma}$.

Using (13)-(14) and $y = Ax + w$, the error vector $\tilde{x}^t = x - \hat{x}^t$ can be expressed as

$$\tilde{x}^t = x - B\hat{x}^{t-1} - \tfrac{1}{K}\bar{A}y = B\tilde{x}^{t-1} - \tfrac{1}{K}\bar{A}w \qquad (16)$$
$$= B^2\tilde{x}^{t-2} - (B + I_p)\tfrac{1}{K}\bar{A}w = \cdots \qquad (17)$$
$$= B^t x - \tfrac{1}{K}\sum_{i=0}^{t-1}B^i\bar{A}w \triangleq B^t x - R_t w, \qquad (18)$$

where we have defined $R_t = \tfrac{1}{K}\sum_{i=0}^{t-1}B^i\bar{A}$, and used the fact that the algorithm is initialized using $\hat{x}^0 = 0$, thus $\tilde{x}^0 = x - \hat{x}^0 = x$. Hence, the evolution of the error vector is governed by the matrix $B$. Using (5) and (18), we can bound the generalization error $\kappa(\hat{x}^t)$ as

$$\kappa(\hat{x}^t) = \|\Sigma^{1/2}(B^t x - R_t w)\|^2 \qquad (19)$$
$$\leq 2\|\Sigma\|(\|B\|^{2t}\|x\|^2 + \|R_t\|^2\|w\|^2), \qquad (20)$$

where $\|\cdot\|$ denotes the spectral norm for matrices, and the Euclidean norm for vectors. Here, we have used properties of the matrix/vector norms and the fact that $\|\Sigma^{1/2}\|^2 = \|\Sigma\|$ for $\Sigma \succeq 0$.

Considering the noise-free setting, i.e., $w = 0$, we obtain the bound

$$\kappa(\hat{x}^t) = \|\Sigma^{1/2}\tilde{x}^t\|^2 = \|\Sigma^{1/2}B^t x\|^2 \qquad (21)$$
$$\leq \|\Sigma\|\|B\|^{2t}\|x\|^2. \qquad (22)$$

In the upcoming sections, we investigate the generalization error in terms of the behaviour of $\|B\|$ for different statistical models for $A$. As our results in the coming sections illustrate, the spectral properties of $B = I_p - \tfrac{1}{K}\bar{A}A$ can heavily depend on the partitioning parameters $p_k$, $k = 1, \ldots, K$. In particular, if any $p_k$ is close to $n$, then $\|B\|$ cannot be bounded with high probability, which is directly reflected in the generalization

error $\kappa(\hat{x}^t)$. Both the noisy and the noise-free setting are studied in the numerical results of Section VIII. The presented results illustrate how the error can be extremely large also with noisy training data, hence the insights gained from our analytical study of the bound in (22) are also relevant for the noisy setting of (19).

We now compare generalization error and the training error. For $w = 0$, the training error associated with $\hat{x}^t$, i.e., the error in reconstructing $\{a_i^T x\}_{i=1}^n$, can be expressed as follows

$$\tfrac{1}{n}\sum_{i=1}^n(a_i^T x - a_i^T\hat{x}^t)^2 = \tfrac{1}{n}\|A\tilde{x}^t\|^2 = \tfrac{1}{n}\|AB^t x\|^2. \qquad (23)$$

Note that for the training error in (23), $\tilde{x}^t = B^t x$ is multiplied with the current realization $A$. On the other hand, for the generalization error in (19), there is a multiplication with $\Sigma^{1/2}$ due to averaging over realizations of the regressor matrix. This distinction can lead to a significant gap between the training and generalization error. We illustrate later in this section, see (27), that the training error is exactly zero under certain partitioning schemes whereas the generalization error can be large. The numerical results in Section VIII further verify that this observation.

By [25, Thm. 2] and strong duality of (2), the solution produced by Algorithm 1 is optimal for the optimization problem in (2). This is realized by making use of the concept "bounded support modification" in [25, eqn. (18)]. We note that with $\lambda = 0$ and $\sigma_w^2 = 0$, an optimal solution gives exactly zero training error (23). On the the other hand, there can exist multiple solutions with zero training error, but with vastly varying generalization error. Hence, a characterization of the generalization error of the distributed algorithm is needed, which is the focus of this paper.

Motivated by the recent results on overparameterization in linear regression [12]–[15] and the success of massively overparameterized models [38, Table 1], we pay special attention to the overparameterized setting of $p_k \geq n$, $\forall k$. Although our results generally hold for all possible partitioning schemes, we obtain some particularly interesting results for this overparameterized setting. The following lemma shows that the governing matrix $B$ is a projection under $p_k \geq n$, $\forall k$, which will be instrumental for the upcoming results:

**Lemma 1.** *Let $B$ be defined as in (14), and the rows of $A$ be drawn i.i.d. from $\mathcal{N}(0, \Sigma)$, with $\Sigma \in \mathbb{R}^{p\times p}$ being positive definite. If all $A_k$ are broad, i.e., $p_k \geq n$, $\forall k$, then with probability one, we have $B^2 = B$.*

Proof: See Section XI-B.

Note that Lemma 1 together with (13) shows that CoCoA, i.e., Algorithm 1, converges in one iteration if the number of unknowns at each node is larger or equal to the number of training samples, i.e., $\hat{x}^t = \hat{x}^1$ for $t \geq 1$ if $p_k \geq n, \forall k$. Hence we have the following corollary for the generalization error:

**Corollary 1.** *Consider the setting of Lemma 1, and $t \geq 1$. Then $\hat{x}^t = \hat{x}^1$, and*

$$\kappa(\hat{x}^t) = \kappa(\hat{x}^1) = \|\Sigma^{1/2}(Bx - \tfrac{1}{K}\bar{A}w)\|^2 \qquad (24)$$

$$\leq 2\|\Sigma\|\left(\|B\|^2\|x\|^2 + \tfrac{1}{K^2}\|\bar{A}\|^2\|w\|^2\right), \qquad (25)$$

*which for $\boldsymbol{w} = 0$ can be tightened as*

$$\kappa(\hat{\boldsymbol{x}}^t) \le \|\boldsymbol{\Sigma}\| \|\boldsymbol{B}\|^2 \|\boldsymbol{x}\|^2. \tag{26}$$

Proof: Using Lemma 1, we have that $\boldsymbol{B}^t = \boldsymbol{B}$ and $\boldsymbol{R}_t = \frac{1}{K}\bar{\boldsymbol{A}}$, hence $\hat{\boldsymbol{x}}^t = \hat{\boldsymbol{x}}^1$. Combining these equations with (19) gives the expression in (24). Similarly as with (20) and (22), (24) can be upper bounded by (25), and by (26) if $\boldsymbol{w} = 0$. $\square$

Note that Lemma 1 provides an interesting observation for the training error in (23). Using that $\boldsymbol{B}^t = \boldsymbol{B}$ for $t \ge 1$, we can simplify the training error in (23):

$$\frac{1}{n} \|\boldsymbol{A}\boldsymbol{B}\boldsymbol{x}\|^2 = \frac{1}{n} \|\boldsymbol{A}(\boldsymbol{I}_p - \frac{1}{K}\bar{\boldsymbol{A}}\boldsymbol{A})\boldsymbol{x}\|^2 \tag{27}$$

$$= \frac{1}{n} \|(\boldsymbol{A} - \frac{1}{K}\boldsymbol{A}\bar{\boldsymbol{A}}\boldsymbol{A})\boldsymbol{x}\|^2 = 0, \tag{28}$$

where we have used that $p_k \ge n$ to apply $\boldsymbol{A}\bar{\boldsymbol{A}} = K\boldsymbol{I}_n$, using Property (e) from Section XI-A of the Appendix. Our results in Theorem 1 – 4 of the coming sections illustrate that $\|\boldsymbol{B}\|$ can be unboundedly large if $p_k$ is too close to $n$. Hence the generalization error in (26) can be unbounded even though the training error is zero.

We conclude this section by motivating our study of the unregularized case. With $\lambda > 0$, convergence to a solution with arbitrarily small optimality gap is guaranteed with a sufficient number of iterations $T$ [25, Thm. 3]. Hence in our problem setting, for $T$ large enough we have that $\hat{\boldsymbol{x}}^T \to (\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A} + \lambda \boldsymbol{I}_p)^{-1}\boldsymbol{A}^{\mathrm{T}}\boldsymbol{y}$, i.e., convergence to the centralized least-squares (LS) solution in (7). While [25, Thm. 3] shows that smaller $\lambda$ will require larger $T$, these results do not show how the partitioning affects the generalization error. We address this gap by first studying the setting with $\lambda = 0$ analytically, and then showing implications of these results for the setting with $\lambda > 0$ before convergence through numerical results.

## V. ISOTROPIC GAUSSIAN REGRESSORS

In this section we present our analysis on the generalization error associated with Algorithm 1 under isotropic Gaussian regressors, i.e., the entries of $\boldsymbol{A}$ are i.i.d. with $\mathcal{N}(0, 1)$, or equivalently, the rows of $\boldsymbol{A}$ are i.i.d. with $\boldsymbol{a}_i \sim \mathcal{N}(0, \boldsymbol{I}_p)$. In Section VI and Section VII, we extend our results to the correlated Gaussian and the sub-gaussian settings. Focusing first on the isotropic Gaussian distributions allows us to give more precise results than the case for more general distributions, see the discussions after Remark 3 and the discussions at the end of Section VII for details.

**Lemma 2.** (Tracy-Widom fluctuations [39]). *For a matrix $\boldsymbol{M} \in \mathbb{R}^{n \times p}$ with i.i.d. $\mathcal{N}(0, 1)$ distributed entries, the following bound holds with probability at least $\rho = 1 - 2e^{-q^2/2}$, $q \ge 0$*

$$\sqrt{r_{\max}} - \sqrt{r_{\min}} - q \le \sigma_{\min}(\boldsymbol{M})$$
$$\le \sigma_{\max}(\boldsymbol{M}) \le \sqrt{r_{\max}} + \sqrt{r_{\min}} + q, \tag{29}$$

*where $r_{\min} = \min\{n, p\}$ and $r_{\max} = \max\{n, p\}$.*

This result quantifies the deviations of the extreme singular values of a standard Gaussian random matrix from their respective expectations. Our main result in this section uses Lemma 2 to provide high-probability bounds for the spectral norm of $\|\boldsymbol{B}\|$:

**Theorem 1.** *Let $\boldsymbol{A} \in \mathbb{R}^{n \times p}$ be a Gaussian random matrix with i.i.d. $\mathcal{N}(0, 1)$ distributed entries, and let $\boldsymbol{B}$ be defined*

*as in* (14). *Let us define $r_{\min,k} = \min\{p_k, n\}$ and $r_{\max,k} = \max\{p_k, n\}$, and let*

$$\beta_G = 1 + \frac{1}{K}\sqrt{K + \sum_{k=1}^{K}\sum_{\substack{i=1 \\ i \ne k}}^{K} \bar{\gamma}_{k,i}}, \tag{30}$$

*where*

$$\bar{\gamma}_{k,i} = \frac{(\sqrt{r_{\max,k}} + \sqrt{r_{\min,k}} + q_k)^2}{(\sqrt{r_{\max,i}} - \sqrt{r_{\min,i}} - q_i)^2}, \tag{31}$$

*and $0 \le q_i < \sqrt{r_{\max,i}} - \sqrt{r_{\min,i}}$, $i = 1, ..., K$. Then, the following bound*

$$\|\boldsymbol{B}\| \le \beta_G, \tag{32}$$

*holds with probability at least $\rho_G = \prod_{k=1}^{K} (1 - 2e^{-q_k^2/2})_+$.*

Proof: See Section XI-C. Note that the probabilistic result in Theorem 1 is with respect to the distribution of the training data in $\boldsymbol{A}$, whereas the expectation taken in $\kappa(\hat{\boldsymbol{x}}^t)$, see (4), is with respect to the unseen data, i.e., the test regressors.

Theorem 1 provides key insights about $\boldsymbol{B}$, which governs the iterations of COCOA, see (13). The matrix $\boldsymbol{B}$ represents the contribution of the local solutions $\hat{\boldsymbol{x}}_k$ from each node, as well as the interactions of these solutions through the shared variable $\bar{\boldsymbol{v}}$. In one node at a given iteration, the setting can be interpreted as a regression problem with a partial model with missing features, i.e., that some of the features in the full model are ignored during regression. The main technical challenge in Theorem 1 is then to capture the combined contribution of many nodes to the error, each of which sees its own partial model. In (32), the bound $\beta_G$ on the spectral norm of $\boldsymbol{B}$ captures this joint behaviour, with dependency on the dimensions of the local regressor matrices $\boldsymbol{A}_k$.

We now connect the presented bound on $\|\boldsymbol{B}\|$ to the generalization error $\kappa(\hat{\boldsymbol{x}}^t)$, illustrating how the overall performance of the solution can depend on the partitioning scheme. Together with (22), Theorem 1 provides an upper bound on the generalization error:

$$\kappa(\hat{\boldsymbol{x}}^t) \le \beta_G^{2t} \|\boldsymbol{x}\|^2, \tag{33}$$

with success probability, i.e., probability of the upper bound holding, at least $\rho_G$. Ideally, the bound in (32) would be small while the success probability $\rho_G$ is large, meaning that the generalization error is small with a high probability. To have a high success probability $\rho_G$, the variables $q_i$ needs to be large for $i = 1, \ldots, K$. On the other hand, as any $q_i$ approaches its upper bound, i.e., $q_i \to \sqrt{r_{\max,i}} - \sqrt{r_{\min,i}}$, the corresponding denominator in $\beta_G$ goes to zero, and the upper bound becomes larger. Thus, we need $\sqrt{r_{\max,i}} - \sqrt{r_{\min,i}}$ to be sufficiently large for $i = 1, \ldots, K$, so that all $q_i$ can be chosen to guarantee a sufficiently large $\rho_G$, without compromising the level of the upper bound. Note that for $r_{\max,i} \approx r_{\min,i}$, Lemma 2 is also typically uninformative. Hence, for a fixed success probability $\rho_G$, the bound on $\|\boldsymbol{B}\|$ grows as $r_{\max,i}$ and $r_{\min,i}$ get closer. Although results that provide a more accurate picture of the behaviour of the minimum singular value of $\boldsymbol{A}_i$ for $r_{\max,i} \approx r_{\min,i}$ [39, Thm. 3.3] exist, a similar line of argument in terms of the effect of $r_{\max,i} \approx r_{\min,i}$ also holds there.

**Remark 1.** *Theorem 1 shows that the bound $\beta_G$ guaran-*

*tees smaller values on the generalization error with higher probability when the dimensions of the submatrices $n$ and $p_k$ are apart from each other, compared to when they are close. This sufficiency result suggests that it may be the case that the generalization error gets large values with $n \approx p_k$.*

Our results on the average generalization error in Section V-A and also the numerical results in Section VIII show that this is indeed the case. Theorem 1 is consistent with the results for the isotropic Gaussian regressors for the centralized setting in [9]–[11], where the relationship between the number of unknowns and the number of observations determines whether a low generalization error is attainable. In particular, average behaviour of the generalized inverse of Wishart matrices [10], [11], as well as other high probability results [11] play a central role leading to "double descent" curves [11]. Similarly in Theorem 1, the spectral properties of the Wishart matrices $A_k A_k^T$, particularly their singular values' closeness to zero, are of central importance. An important difference between the existing literature and our work in this section (as well as our work in the cases of correlated Gaussian and sub-gaussian distributions in the subsequent sections) is the fact that our results focus on the distributed setting and we explain how the trade-offs between the number of observations and the unknowns studied in the centralized case have important implications for distributed learning.

For the setting where all local models at the nodes are overparameterized, we present the following tighter bound that holds for all $t$ on the generalization error:

**Corollary 2.** Let $p_k \geq n, \forall k$. The generalization error $\kappa(\hat{x}^t)$ is bounded by $\kappa(\hat{x}^t) \leq \beta_G^2 \|x\|^2$ with probability $\rho_G$.
Proof: The result follows directly from Theorem 1 and Corollary 1. □

We now provide the following alternative bound for $\|B\|$ for the case with $n \geq p_k, \forall k$:

**Lemma 3.** *Consider the setting of Theorem 1 under $n \geq p_k, \forall k$. Let*
$$\bar{\beta}_G = \sqrt{\frac{(K-1)^2}{K} + \frac{1}{K^2}\sum_{k=1}^K \sum_{\substack{i=1\\i\neq k}}^K \bar{\gamma}_{k,i}}, \quad (34)$$

*with $\bar{\gamma}_{k,i}$ as defined in Theorem 1. Then, the following bound*
$$\|B\| \leq \bar{\beta}_G, \quad (35)$$

*holds with probability at least $\bar{\rho}_G = \prod_{k=1}^K (1 - 2e^{-q_k^2/2})_+$.*
Proof: See Section XI-E. This result provides an alternative to (33): $\kappa(\hat{x}^t) \leq \bar{\beta}_G^{2t}\|x\|^2$, which holds with probability at least $\bar{\rho}_G$. In Section V-B we compare the high-probability results of this section, in particular Corollary 2 and Lemma 3, to the average generalization error.

*A. The Average Generalization Error*

We now compare our results with the average generalization error
$$\mathbb{E}_{A,w}[\kappa(\hat{x})] = \mathbb{E}_{A,w}[\|x - \hat{x}\|^2], \quad (36)$$
where the expectation is over the regressor matrix $A$ in the training data and the training noise $w$. We consider the following extension of [32, Thm. 1] to the case with training noise:

**Lemma 4.** *Let $A \in \mathbb{R}^{n\times p}$ be a random matrix with i.i.d. $\mathcal{N}(0,1)$ distributed entries. Let $w \in \mathbb{R}^{n\times 1}$ be a zero-mean random vector with covariance matrix $\Sigma_w$, statistically independent with the regressors. The average generalization error in iteration $t = 1$ of Algorithm 1, can be expressed as*
$$\mathbb{E}_{A,w}[\kappa(\hat{x}^1)] = \sum_{k=1}^K \|x_k\|^2 \alpha_k + \frac{\text{tr}(\Sigma_w)}{K^2}\gamma_k, \quad (37)$$

*where $\alpha_k, \gamma_k, k = 1, \ldots, K$, are given by*
$$\alpha_k = \frac{1}{K^2}(K^2 + (1-2K)\frac{r_{\min,k}}{p_k} + \sum_{\substack{i=1\\i\neq k}}^K \gamma_i), \quad (38)$$
$$\gamma_k = \begin{cases} \frac{r_{\min,k}}{r_{\max,k}-r_{\min,k}-1} & \text{for } p_k \notin \{n-1,n,n+1\}, \quad (39a)\\ +\infty & \text{otherwise}, \quad (39b) \end{cases}$$

*and $r_{\min,i} = \min\{p_i, n\}$ and $r_{\max,i} = \max\{p_i, n\}$.*
Proof: See Section XI-F. Similar to Theorem 1, the average error diverges for $p_i \in \{n-1, n, n+1\}$ due to pseudo-inverses, see [32] and the related discussions in the longer version [40, Section VII-C].

**Remark 2.** *Both Theorem 1 and Lemma 4 suggest that we may have a large generalization error when the number of unknowns at a node is close to the number of observations, i.e., at least one of the local system of equations is approximately square, regardless of being under- or over-parameterized.*

The following corollary illustrates how Lemma 4 can be used to provide an error expression for all iterations in the overparametrized case:

**Corollary 3.** *Let $p_k \geq n, \forall k$. The average generalization error for any iteration $t \geq 1$ is given by*
$$\mathbb{E}_{A,w}[\kappa(\hat{x}^t)] = \sum_{k=1}^K \|x_k\|^2 \alpha_k + \frac{\text{tr}(\Sigma_w)}{K^2}\gamma_k, \quad (40)$$

*where $\alpha_k, \gamma_k$ are given by Lemma 4.*
Proof: By combining Lemma 4 and Corollary 1, i.e., $\kappa(\hat{x}^t) = \kappa(\hat{x}^1)$, we obtain the desired result. □

The numerical results of [32] as well as the results in Section VIII (see Figure 2) suggest that (40) provides not only the error for the overparametrized case but also reveals the general approximate behaviour of the algorithm even if the $p_k \geq n$ condition is not satisfied for $k = 1, \ldots, K$.

*B. Comparison with the Average Generalization Error*

We now consider an example where we first study the expectation results from Lemma 4 and Corollary 3, and then compare them to the probabilistic results in Lemma 3 and Corollary 2. We here consider the setting with $w = 0$.

Let $\|x_k\|^2 = \frac{1}{K}$. By Lemma 4, we have $\mathbb{E}_A[\kappa(\hat{x}^1)] = \frac{1}{K}\sum_{k=1}^K \alpha_k$. Using $\frac{r_{\min,k}}{p_k} \leq 1$, we obtain

$$\mathbb{E}_A[\kappa(\hat{x}^1)] \leq \frac{1}{K}\left(\frac{(K-1)^2}{K} + \frac{1}{K^2}\sum_{k=1}^K\sum_{\substack{i=1\\i\neq k}}^K \gamma_i\right), \quad (41)$$

$$\leq 1 + \frac{1}{K^2} + \frac{1}{K^3}\sum_{k=1}^K\sum_{\substack{i=1\\i\neq k}}^K \gamma_i. \quad (42)$$

Under $p_k \geq n$, using Corollary 3 we observe that
$$\mathbb{E}_A[\kappa(\hat{x}^t)] = \mathbb{E}_A[\kappa(\hat{x}^1)]. \quad (43)$$

We now consider the probabilistic results in two separate cases:

i) Let $n \geq p_k, \forall k$. By Lemma 3, the following holds with probability at least $\bar{\rho}_G$ for $t = 1$:

$$\kappa(\hat{\boldsymbol{x}}^1) \leq \bar{\beta}_G^2 = \frac{(K-1)^2}{K} + \frac{1}{K^2} \sum_{k=1}^{K} \sum_{\substack{i=1 \\ i \neq k}}^{K} \bar{\gamma}_{k,i}. \quad (44)$$

Comparing (44) with (41), we observe that the expressions have a shared algebraic form where the expectation result in (41) has a scaling of $\frac{1}{K}$ compared to the probabilistic result in (44), under $\bar{\gamma}_{k,i} = \gamma_i$. Both results reveal how the dimensions of the local data matrices $\boldsymbol{A}_k$ affect the error: the expectation results in (41) through $\gamma_i$ and the probability results in (44) through $\bar{\gamma}_{k,i}$.

ii) Let $p_k \geq n, \forall k$. Using Corollary 2 and (30), we observe that the following holds with probability at least $\rho_G$

$$\kappa(\hat{\boldsymbol{x}}^t) \leq \beta_G^2 \leq 2 + \frac{2}{K} + \frac{2}{K^2} \sum_{k=1}^{K} \sum_{\substack{i=1 \\ i \neq k}}^{K} \bar{\gamma}_{k,i}, \quad (45)$$

where we used $(a+b)^2 \leq 2a^2 + 2b^2$ on (30). Using (43), we compare (45) to (42): The two bounds again quantify the dependence of the error on the partitioning through $\gamma_i$ and $\bar{\gamma}_{k,i}$ and they have the same shared form under $\gamma_i = \bar{\gamma}_{k,i}$.

This example emphasizes the common algebraic structure in the expectation and the high-probability results. Although the bounds in (44)/(45) and (41)/(42) contain different constant additive terms, they all heavily depend on the terms $\bar{\gamma}_{k,i}$ and $\gamma_i$ which are the main factors characterizing the behaviour of the generalization error with respect to the partitioning.

## VI. CORRELATED GAUSSIAN REGRESSORS

This section generalizes the results of the preceding section to correlated Gaussian regressors. The regressors $\boldsymbol{a}_i$ (i.e., the rows of $\boldsymbol{A}$) are now i.i.d. zero-mean random vectors drawn from the Gaussian distribution $\mathcal{N}(0, \boldsymbol{\Sigma})$, with $\boldsymbol{\Sigma} = \mathbb{E}_{\boldsymbol{a}_i}[\boldsymbol{a}_i \boldsymbol{a}_i^T]$, i.e., each row of $\boldsymbol{A}$ is independently drawn from $\mathcal{N}(0, \boldsymbol{\Sigma})$. Our main result in this section is given by Theorem 2:

**Theorem 2.** *Let $\boldsymbol{B}$ be defined as in* (14)*, and the rows of $\boldsymbol{A}$ be i.i.d. with $\mathcal{N}(0, \boldsymbol{\Sigma})$, $\boldsymbol{\Sigma} \succ 0$, and let $\mathcal{K} = \{k : n < p_k\}$. Let*

$$\beta_{G_c} = 1 + \frac{1}{K} \sqrt{K + \sum_{k=1}^{K} \sum_{\substack{i=1 \\ i \neq k}}^{K} \frac{n\sigma_{\max}(\boldsymbol{\Sigma}_k) + \ell_k(q_k)}{\eta_i}}, \quad (46)$$

*where $\boldsymbol{\Sigma}_k$ is the $k^{th}$ principal submatrix of $\boldsymbol{\Sigma}$ and*

$$\eta_k = \begin{cases} (n\sigma_{\min}(\boldsymbol{\Sigma}_k) - \ell_k(q_k))_+, & n \geq p_k & (47a) \\ \sigma_{\min}(\boldsymbol{\Sigma}_k)(\sqrt{p_k} - \sqrt{n} - \bar{q}_k)_+^2, & n < p_k & (47b) \end{cases}$$

*and*

$$\ell_k(q_k) = \tfrac{8}{3} n \sigma_{\max}(\boldsymbol{\Sigma}_k) C \left( \sqrt{\tfrac{p_k+q_k}{n}} + \tfrac{p_k+q_k}{n} \right), \quad (48)$$

*where $C$ is an absolute constant, and $q_k, \bar{q}_k \geq 0$, $\forall k$. Then, the following bound*

$$\|\boldsymbol{B}\| \leq \beta_{G_c}, \quad (49)$$

*holds with probability at least $\rho_{G_c} = 1 - \sum_{k=1}^{K} 2e^{-q_k} - \sum_{k \in \mathcal{K}} 2e^{-\bar{q}_k^2/2}$.*

Proof: See Section XI-G. Similar to Theorem 1, Theorem 2 illustrates the connection between the dimensions of the local matrices $\boldsymbol{A}_k$ and the norm $\|\boldsymbol{B}\|$.

We also present a result for the special case $n \geq p_k$, analogous to Lemma 3 but for the correlated Gaussian setting:

**Lemma 5.** *Consider the setting of Theorem 2, under $n \geq p_k, \forall k$. With probability at least $1 - 2\sum_{k=1}^{K} e^{-q_k}$, (35) holds with the following redefinition $\bar{\gamma}_{k,i} = (n\sigma_{\max}(\boldsymbol{\Sigma}_k) + \ell_k(q_k))/(n\sigma_{\min}(\boldsymbol{\Sigma}_k) - \ell_i(q_i))_+$.*
Proof: The proof follows the same line of argument as Lemma 3, where $\bar{\gamma}_{k,i}$ is defined using (47a). $\square$

We now combine Theorem 2 with (22) and obtain the following upper bound on the generalization error

$$\kappa(\hat{\boldsymbol{x}}^t) \leq \beta_{G_c}^{2t} \|\boldsymbol{\Sigma}\| \|\boldsymbol{x}\|^2. \quad (50)$$

**Remark 3.** *Theorem 2 is consistent with Theorem 1, also illustrating how the generalization error of the solution produced by COCOA (Algorithm 1) is affected by the partitioning scheme: if $n$ and $p_k$ are sufficiently far apart and $\boldsymbol{\Sigma}_k$'s are well-conditioned, similar to the case of $\boldsymbol{\Sigma} = \boldsymbol{I}_p$, then a low generalization error is guaranteed with high probability.*

We note that extending the results from the isotropic to the correlated Gaussian setting introduced more complexity in the expressions. In Theorem 2, the bound on $\|\boldsymbol{B}\|$ is expressed up to an absolute constant $C$, whereas in the isotropic setting of Theorem 1, a more refined bound was presented based on the stronger results for typical behaviour of such Gaussian matrices, as in Lemma 2.

Theorem 2 emphasizes the relation between the partitioning and the generalization error, and points out a nontrivial dependency on the regressors' covariance matrix through the dependence on $\sigma_{\min}(\boldsymbol{\Sigma}_k)$ and $\sigma_{\max}(\boldsymbol{\Sigma}_k)$. These results are consistent with the results in the centralized setting [13], [14], [12, Section 5.1], which also illustrate that the performance is connected to both the dimensions of the problem as well as the spectral properties of $\boldsymbol{\Sigma}$. In [13] and [12, Section 5.1], it is emphasized that the generalization performance depends not only on the dimensions of the problem, but also on the relative geometry of the regressors' covariance matrix and the unknowns. In [14], the decay of the covariance matrix's singular values is emphasized as a key indicator of whether a small generalization error can be achieved with sub-gaussian regressors, for which Gaussian regressors is a special case. An important point here is the distinction between the error in $\boldsymbol{x}$ and the generalization error. For instance, in (47b), a large discrepancy between $\sigma_{\max}(\boldsymbol{\Sigma}_k)$ and $\sigma_{\min}(\boldsymbol{\Sigma}_k)$ will lead to a large bound on $\|\boldsymbol{B}\|$ for $p_k > n$ whereas by (21) the generalization error can be potentially small, for instance, when there is only one large eigenvalue, hence typical regressor realizations are approximately the same; and hence the generalization error (error in $\boldsymbol{a}^T \boldsymbol{x}$) is small.

## VII. SUB-GAUSSIAN REGRESSORS

In this section, we consider regressors drawn from sub-gaussian distributions. The family of sub-gaussian distributions include the Gaussian, uniform and the Bernoulli random variables as well as any other bounded random variable [41], hence it allows us to investigate a large range of data distributions.

## A. Preliminaries on Sub-gaussian Random Variables

This section provides preliminaries on sub-gaussian random variables [41].

**Definition VII.1.** (Sub-gaussian random variables) *A random variable $z \in \mathbb{R}$ is called sub-gaussian if there exists a constant $L > 0$ so that the following is satisfied*

$$\mathbb{E}\left[e^{z^2/L^2}\right] \leq 2. \qquad (51)$$

*The smallest $L$ defines the sub-gaussian norm $\|z\|_{\psi_2}$ as follows*

$$\|z\|_{\psi_2} = \inf\{L > 0 : \ \mathbb{E}\left[e^{z^2/L^2}\right] \leq 2\}. \qquad (52)$$

This definition can be extended to higher dimensions:

**Definition VII.2.** (Sub-gaussian random vectors) *A random vector $z \in \mathbb{R}^{p \times 1}$ is called sub-gaussian if for all $h \in \mathbb{R}^{p \times 1}$, $z^{\mathrm{T}}h$ is a sub-gaussian random variable.*

With a slight abuse of notation, we use $a \sim \mathcal{S}(\Sigma)$ to denote that the random vector $a$ comes from some zero-mean sub-gaussian distribution $\mathcal{S}$, and has the covariance matrix $\Sigma$. We introduce the following notation for the sub-gaussian norm

$$\psi_a(\Sigma, h) = \|h^{\mathrm{T}}a\|_{\psi_2}, \qquad (53)$$

where $h \in \mathbb{R}^{p \times 1}$.

## B. Generalization Error under Sub-gaussian Regressors

We now present our main results for the sub-gaussian case. In the following theorem, we assume that the submatrices $A_k \in \mathbb{R}^{n \times p_k}$ are generated from a matrix $Z_k \in \mathbb{R}^{n \times p_k}$ where each entry of $Z_k$ is drawn i.i.d. from $\mathcal{S}(1)$, $\forall k$. This way of generating $A_k$'s renders the matrices $A_k$ statistically independent, and the covariance matrix of the rows of $A$ block-diagonal.

**Theorem 3.** *Let the matrix $B$ be defined as in* (14), *with each $A_k$ generated as $A_k = Z_k \Lambda_k^{1/2} U_k^{\mathrm{T}} \in \mathbb{R}^{n \times p_k}$, where the entries of $Z_k \in \mathbb{R}^{n \times p_k}$ are i.i.d. with $\mathcal{S}(1)$, $\forall k$, $\Lambda_k \in \mathbb{R}^{p_k \times p_k}$ is diagonal and positive definite, and $U_k \in \mathbb{R}^{p_k \times p_k}$ is unitary. Let $\Sigma_k$ denote the associated covariance matrix for the rows of $A_k$. Let $\beta_S$ be defined as*

$$\beta_S = 1 + \frac{1}{K}\sqrt{K + \sum_{k=1}^{K}\sum_{\substack{i=1 \\ i \neq k}}^{K} \frac{n\sigma_{\max}(\Sigma_k) + \ell_k(q_k)}{\eta_i}}, \quad (54)$$

*where*

$$\eta_k = \begin{cases} (n\sigma_{\min}(\Sigma_k) - \ell_k(q_k))_+, & n \geq p_k, & (55a) \\ \sigma_{\min}(\Sigma_k)(\sqrt{p_k} - CL_k^2(\sqrt{n} + \bar{q}_k))_+^2, & n < p_k, & (55b) \end{cases}$$

$$\ell_k(q_k) = CL_k^2\left(\sqrt{\frac{p_k + q_k}{n}} + \frac{p_k + q_k}{n}\right), \qquad (56)$$

*and $C$ is an absolute constant, $q_k$, $\bar{q}_k \geq 0$, $\forall k$, and $L_k \geq 1$ are constants such that, for all $h \in \mathbb{R}^{p_k \times 1}$*

$$\psi_{a_{i,k}}(\Sigma_k, h) \leq L_k\sqrt{h^{\mathrm{T}}\Sigma_k h}, \qquad (57)$$

*where $a_{i,k} \in \mathbb{R}^{p_k \times 1}$ comes from the same distribution as the rows of $A_k$. Then, the following bound holds*

$$\|B\| \leq \beta_S, \qquad (58)$$

*with probability at least $\rho_S = 1 - \sum_{k=1}^{K} 2e^{-q_k} - \sum_{k \in \mathcal{K}} 2e^{-\bar{q}_k^2}$.*
Proof: See Section XI-J. Note that we use the subscript $i, k$ on $a_{i,k}$ in (57) to emphasize that $a_{i,k}$ is i.i.d. with the rows of $A_k$.

In our results with Gaussian regressors, we utilize the fact that for the partitions $A_k$ with i.i.d. Gaussian rows, there is always a decomposition with $Z_k$ which has entries from $\mathcal{N}(0, 1)$ (but $Z_k$'s are not necessarily i.i.d.). With sub-gaussian rows, this type of inverse relationship (i.e. from $A$ with sub-gaussian rows with a certain sub-gaussian norm to $Z_k$ with i.i.d. sub-gaussian elements with a given norm) is not straightforward. Hence, we here focus on covariance structures enabling such a relationship, constructing $A_k = Z_k \Lambda_k^{1/2} U_k^{\mathrm{T}}$ in Theorem 3, and assuming $n \geq p_k$ in the following theorem,

**Theorem 4.** *Let the matrix $B$ be defined as in* (14) *and the rows of $A$ be i.i.d. with $\mathcal{S}(\Sigma)$, with $n \geq p_k$, $\forall k$, and $\Sigma \succ 0$. Let $\Sigma_k$ denote the $k^{th}$ principal submatrix [42, Sec. 0.7.1] of $\Sigma$. Let $\beta_{S_c}$ be defined as*

$$\beta_{S_c} = 1 + \frac{1}{K}\sqrt{K + \sum_{k=1}^{K}\sum_{\substack{i=1 \\ i \neq k}}^{K} \frac{n\sigma_{\max}(\Sigma_k) + \ell_k(q_k)}{(n\sigma_{\min}(\Sigma_i) - \ell_i(q_i))_+}}, \quad (59)$$

*where $\ell_k(q_k)$ is defined in* (56)-(57). *Then, the following bound holds*

$$\|B\| \leq \beta_{S_c}, \qquad (60)$$

*with probability at least $\rho_{S_c} = 1 - \sum_{k=1}^{K} 2e^{-q_k}$.*
Proof: See Section XI-K.

Similar to the previous results with Gaussian regressors, we obtain the bounds on the generalization error as $\kappa(\hat{x}^t) \leq \beta_S^{2t}\|\Sigma\|\|x\|^2$ and $\kappa(\hat{x}^t) \leq \beta_{S_c}^{2t}\|\Sigma\|\|x\|^2$, by Theorem 3 and Theorem 4, respectively.

Theorem 3 and 4 provide analogous insights as Theorem 1 and 2 in the sense that the bounds $\beta_S$ and $\beta_{S_c}$ depend on the dimensions of the local regressor matrices $A_k$ and on the corresponding covariance matrices $\Sigma_k$.

**Remark 4.** *Theorem 3 and 4 are consistent with Theorem 1 and 2: all of these results provide bounds on the generalization error that can be guaranteed to have smaller values if $n$ and $p_k$ are further apart compared to the case when they are closer.*

Theorem 3 is consistent with the centralized setting of [14] with sub-gaussian regressors. In [14], the dimensions $p$ and $n$ as well as the spectral properties of the regressors' covariance matrix $\Sigma$ are pointed out as important factors determining the generalization error. We correspondingly highlight the local dimensions $p_k$ and $n$; and the local covariance matrices $\Sigma_k$. In [15], bounds on the generalization error with sub-gaussian regressors, focusing on the effects of training noise, are derived in the centralized setting. While our results focus on the noise-free distributed setting, the implications of the noisy interpolation results of [15] are considered as an important line of future work.

The family of sub-gaussian distributions includes a large variety of distributions, including Gaussian regressors of Section V–VI and Bernoulli regressors, popular in compressive sensing [43]. Furthermore, all bounded distributions are sub-gaussian distributions. For instance, the regressors formed by random Fourier features [44], used in various classification tasks and also studied in Section VIII-B, are sub-gaussian since the magnitude of the elements of these regressors are bounded by 1. We note that our results for the Gaussian settings are more refined than those in Theorem 3 and 4, due to
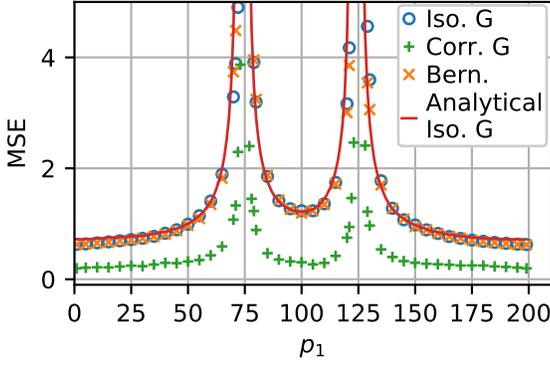
Fig. 2: The generalization error for the three synthetic datasets together with the analytical expectation from Lemma 4.



Fig. 3: The generalization error in terms of MSE and zero-one loss for the MNIST example.

the existence of more precise results for Gaussian distribution than for the broad family of sub-gaussian distributions. On the other hand, Theorem 3 and 4 cover the very general setting of sub-gaussian distributions, and could be further specialized for other special cases of sub-gaussian distributions, such as for applications with Bernoulli regressors.

## VIII. NUMERICAL RESULTS

We now illustrate the behaviour of the generalization error with data from the distributions discussed in the preceding sections, as well as image data from the MNIST dataset [37].

We first explain the experimental setup for with the synthetic datasets. We consider the following distributions for $a$: a) Isotropic Gaussian (Iso. G.) with $\mathcal{N}(0, I_p)$; b) Correlated Gaussian (Corr. G.) with $\mathcal{N}(0, \bar{\Sigma})$ with a non-diagonal $\bar{\Sigma} \in \mathbb{R}^{p \times p}$; c) Bernoulli (Bern.) distribution on $\{-1, 1\}$ with $\Sigma = I_p$, i.e., $a_{ij}$ is $-1$ or $1$ with probability $1/2$. These constitute examples for the settings of Section V, Section VI and Section VII, respectively. As our example for the sub-gaussian distributions, we consider the Bernoulli distribution which is commonly used, for example, in compressive sensing literature [43]. The covariance matrix $\bar{\Sigma} = U \Lambda U^{\mathrm{T}}$ is fixed throughout the experiments and chosen as follows: $U \in \mathbb{R}^{p \times p}$ is sampled from a Haar distribution [22] and the eigenvalues are given by $\Lambda = \mathrm{diag}(\mu_i) \in \mathbb{R}^{p \times p}$, with $\tilde{\mu}_{i+1} = 0.9631 \tilde{\mu}_i$ and $\mu_i = p\tilde{\mu}_i / \sum_{i=0}^{p-1} \tilde{\mu}_i$. The parameter vector $x$ is fixed for all experiments, randomly chosen with i.i.d. uniform elements on $[-1, 1]$ and normalized so that $\|x\| = 1$. We set $n = 75$, $p = 200$ and use a network of $K = 2$ nodes, hence $p = p_1 + p_2$. Algorithm 1 is run for $T = 1000$ iterations with $\lambda = 0$ unless otherwise stated. The generalization error is reported as the emprical mean-squared error (MSE) which is calculated as $\text{MSE} \triangleq \frac{1}{\bar{n}N} \sum_{i=1}^{N} \left\| A_{test,(i)}(x - \hat{x}_{(i)}^T)) \right\|^2$. Average simulation results for $N = 100$ realizations of the training data $A_{(i)}$, $i = 1, \ldots, N$ are reported. Here, $A_{test,(i)} \in \mathbb{R}^{\bar{n} \times p}$ denotes the test data matrix for experiment $i$, $\hat{x}_{(i)}^T = \hat{x}_{(i)}^T(A_{(i)})$ is the solution found by Algorithm 1 after its final iteration $T$ under $y_{(i)} = A_{(i)}x$, and $\bar{n} = 100n$ is the number of observations (i.e. rows) in each $A_{test,(i)}$. Unless otherwise stated, all plots provide the performance of the algorithm after convergence.

In addition to the above, we also consider the digit classification problem from the MNIST dataset [37], [45] in order
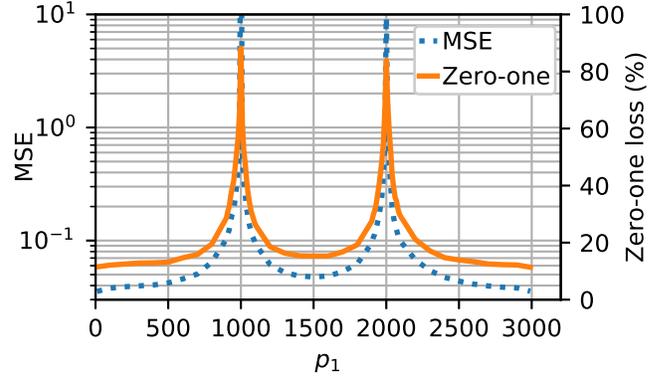
to further illustrate the practical implications of our results. This dataset poses a classification problem consisting of ten classes, i.e., digits. We convert each 28-by-28 image to a 784-by-1 vector $z_i$ and transform the data using the following random features [44]: $a_i = [\cos(z_i^{\mathrm{T}} \omega_1); \ldots; \cos(z_i^{\mathrm{T}} \omega_p)] \in \mathbb{R}^{p \times 1}$, $p = 3 \times 10^3$ where $\omega_i \sim \mathcal{N}(0, \zeta^2 I_{784})$ with $\zeta = 0.2$. The matrix of regressors $A \in \mathbb{R}^{n \times p}$ is obtained by using $a_i^{\mathrm{T}}$ as its rows. We train one classifier for each class and apply a one-v.s.-rest classification strategy [46]. We subsample the training dataset with a factor of 60, resulting in $n = 10^3$ samples. For the test, we use the full test dataset with $\bar{n} = 10^4$. We report both the MSE and the classification error on the test data.

### A. Generalization Error and the Partitioning of the Model

In Figure 2, we present the empirical generalization error associated with the solution of CoCoA (Algorithm 1) as $p_1$ is varied from 1 to $p - 1$. The results for the three synthetic datasets together with the theoretical expected generalization error from Lemma 4 for the isotropic Gaussian case (Analytical Iso. G) are provided. These plots illustrate that the generalization error depends significantly on the partitioning. For all datasets, the average generalization error blows up as either $p_1$ or $p_2$ approaches $n$, and it is relatively low when $p_1$ and $p_2$ are both far from $n$. In particular, the peak MSEs are given by $2.2 \times 10^9$, $8.1 \times 10^2$ and $1.2 \times 10^6$ for the cases a) – c), respectively. Note that these values are comparably large and far outside the range of the plot, hence they are truncated, in Figure 2.

These observations are consistent with Theorems 1 – 4, demonstrating that the generalization error is small with high probability when $p_1$ and $p_2$ are far from $n$, while small values cannot be guaranteed when $p_1$ or $p_2$ are close to $n$. We now report the performance of the centralized solution in (7). For all the cases, i.e., Iso. G, Corr. G and Bern., the training error is below $10^{-21}$ for the centralized solution as well as for the distributed solution for all values of $p_1$ (values are not included in the plots). The generalization error for the centralized solution is $0.63$, $0.20$ and $0.62$ for the cases Iso.G, Corr.G, Bern., respectively.

These results illustrate that the partitioning can greatly affect the generalization error, making it significantly larger than what the centralized solution achieves, while the training performance is on the same level as the centralized solution.
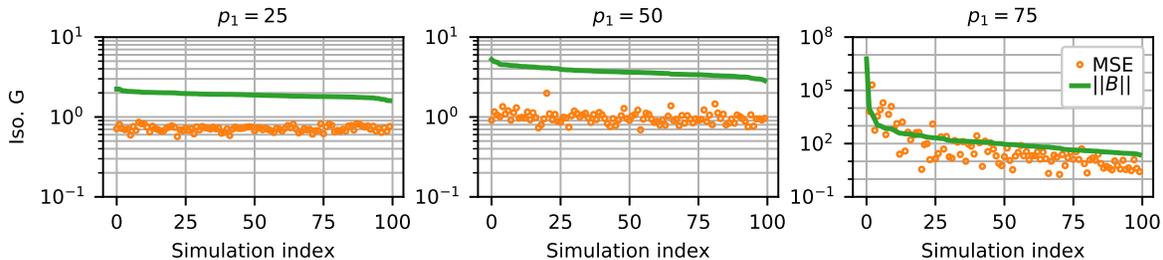
Fig. 4: Generalization error (MSE) versus $\|\boldsymbol{B}\|$ for each realization, i.e., simulation, of the training data. The simulation indices are ordered so that $\|\boldsymbol{B}\|$ decreases monotonically with increasing simulation indices.

The plots for the Iso. G. data in Figure 2 illustrates a close match between the empirical average generalization error and the expectation results in Lemma 4. This observation emphasizes that the result in Corollary 3 can be relevant even if $p_k \geq n$ is not fulfilled for all nodes.

### B. Generalization Error on MNIST data with CoCoA

In Figure 3, we plot the MSE and the zero-one loss, i.e., the percentage of incorrect classifications, for the MNIST test data. Similary as in Figure 2, the generalization error significantly depends on the partitioning at the nodes, both in terms of MSE and zero-one loss. In particular, we see an extremely large error if any $p_k$ is close to $n$ compared to the case where $p_k$ and $n$ are significantly different. In particular, the generalization error (in terms of MSE) with $p_k = n$ is given by $2.5 \times 10^3$ whereas the training error is below $10^{-27}$ for all choices of $p_1$. The centralized solution in (7) also achieves a training MSE below $10^{-27}$ while the corresponding generalization error is $3.5 \times 10^{-2}$.

These results highlight practical consequences of design choices in distributed learning. In particular, it is not only for data coming exactly from certain probability distributions, but also for practical real-world datasets that the generalization error significantly depends on the partitioning over the nodes. Furthermore, the results here together with the results for the synthetic data in Section VIII-A suggest that in order to have a low generalization error one should avoid a partitioning where $p_k$ is close to $n$ for any node.

### C. Generalization Error and Spectral Norm of $\boldsymbol{B}$

Theorems 1-4 highlight the dependence of the generalization error on the spectral norm of $\boldsymbol{B}$. We now further investigate this relationship. For ease of disposition, we consider only the isotropic Gaussian data. In Figure 4, we plot $\|\boldsymbol{B}\|$ and the generalization error (the MSE) for each of the 100 different realizations of the training dataset, i.e., $\boldsymbol{A}$, that we have averaged over in Section VIII-A. Each simulation index corresponds to one realization of the training dataset, i.e., one realization of $\boldsymbol{A}$. For each simulation index, the corresponding spectral norm $\|\boldsymbol{B}\|$ and the MSE is provided. The simulation indices are arranged so that $\|\boldsymbol{B}\|$ is monotonically decreasing from left to right.

Comparing the plots for $p_1 \in \{25, 50, 75\}$, we observe that the MSE level depends on $\|\boldsymbol{B}\|$ in a consistent manner. The partitioning $p_1 = 25$ gives the lowest values of $\|\boldsymbol{B}\|$, as well
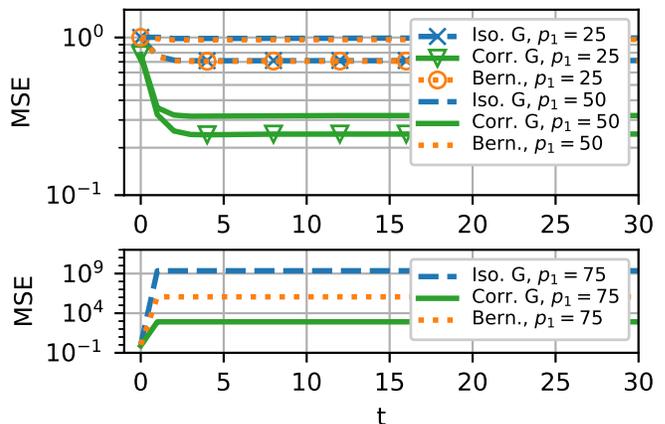


Fig. 5: The average generalization error of the unregularized CoCoA evaluated at each iteration.

as the lowest values of the MSE. When $p_1$ is increased to $p_1 = 50$, both $\|\boldsymbol{B}\|$ and the MSE increases slightly. Consistent with Remark 1, with $p_1 = 75$ (hence $p_1 = n = 75$, where $n$ is the number of observations) both $\|\boldsymbol{B}\|$ and the MSE start to take extremely large values, such as up to $10^5$ for the MSE.

For $p_1 = 25$ and $p_1 = 50$, both $\|\boldsymbol{B}\|$ and the MSE are concentrated around their mean over different simulation indices. This illustrates that with this type of partitioning, it is possible to obtain reliable performance over different training datasets. On the other hand, we observe an extremely large spread in the MSE (from 1 to over $10^5$) over different simulation indices when $p_1 = 75$, illustrating how under this partitioning, the generalization performance can vary substantially over different realizations of the training dataset.

### D. Behaviour of the Generalization Error over Iterations

We now investigate convergence of the generalization error over iterations of CoCoA. Furthermore, we verify the analytical result from Lemma 1.

In Figure 5, the average generalization error associated with the solution produced by each iteration of CoCoA is plotted for the first 30 iterations. There is no visual change in the error values on the plot in the later iterations, hence this range is chosen to be able to better illustrate the transient behaviour.

We observe that the algorithm on average converges quickly, within the first few iterations for all cases. The curves for $p_1 = 75$ are consistent with the result of Lemma 1: here all nodes have $p_k \geq n$, hence the algorithm converges in one iteration. (Although the plots shows only the average, this is
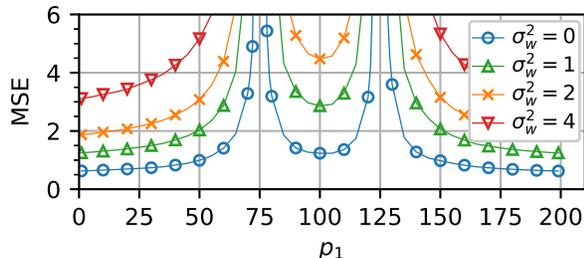
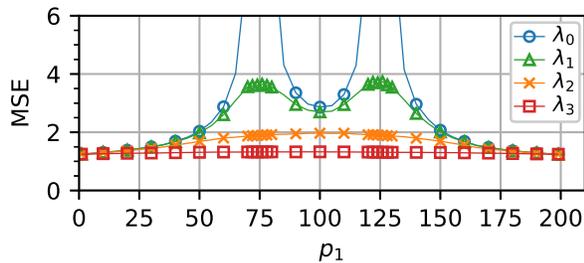Fig. 6: The average generalization error for the Iso. Gaussian case for $\lambda = 0$ and varying noise levels $\sigma_w^2$.



Fig. 7: The average generalization error for the Iso. Gaussian case evaluated at $t = 1000$ for varying values of $\lambda$ Here, $(\lambda_0, \lambda_1, \lambda_2, \lambda_3) = (0, 10^{-3}, 10^{-2}, 10^{-1})$, and $\sigma_w^2 = 1$.

also true for the individual runs.) For the cases of $p_1 = 25$ and $p_1 = 50$, although the results of Lemma 1 does not directly apply, the quick convergence suggests that $B^t$ becomes an approximate projection matrix in the first iterations.

### E. Generalization Error and Noisy Training Data

We now focus on the effect of noise. In particular, we consider the setting with $\lambda = 0$ with additive Gaussian noise in the training data, i.e., with $w_i \sim \mathcal{N}(0, \sigma_w^2)$ in the model (1). Note that the test data is noise free, and that the generalization error is defined as in (5). We consider the case with isotropic Gaussian regressors.

In Figure 6, we present the generalization error for four different noise levels, $\sigma_w^2 = \{0, 1, 2, 4\}$. We observe that the overall level of the generalization error increases with the noise level. There are again large peaks in the error for $p_k$ values close to $n$, suggesting that noisy data does not dampen the peaks, further supporting the insights gained from our analytical results in Theorem $1 - 4$. We emphasize that Figure 6 shows the generalization error for $t = 1000$, which is well after convergence. Our simulations show that the generalization error in the unregularized, noisy setting converges in the same quick fashion as the unregularized noise free setting as illustrated in Figure 5. Due to space limitations, these convergence plots are not included here.

### F. Generalization Error and Regularization

In the preceding sections, we have considered the unregularized scenario. In this section, we provide results that illustrate that even with regularization the partitioning scheme can have a large impact on the generalization error. In particular, the generalization error heavily depends on the partitioning before convergence.
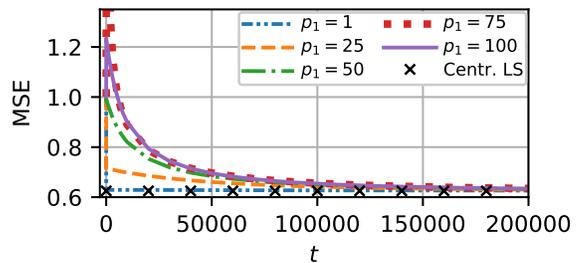


Fig. 8: The average generalization error evaluated in each iteration for CoCoA with regularization $\lambda = 10^{-3}$ and $\sigma_w^2 = 0$.

We consider the scenario with $w_i \sim \mathcal{N}(0, 1)$ in (1) and with varying values of the regularization parameter $\lambda$ in (2). In Figure 7, we compare the generalization error associated with the solutions produced by CoCoA at iteration $t = 1000$, with Iso. G. regressors. We consider the following choices of $\lambda$: $\lambda_0 = 0$, $\lambda_1 = 10^{-3}$, $\lambda_2 = 10^{-2}$ and $\lambda_3 = 10^{-1}$. The figure illustrates that the generalization error peaks are dampened for all three choices of $\lambda > 0$, compared to $\lambda = 0$: For $\lambda_0 = 0$, the peak generalization error is $\approx 10^9$, whereas for $\lambda_1 - \lambda_3$ the peak error values are between $\approx 1$ and $\approx 4$. The minimum generalization error over all partitioning choices is around 1. Hence, with larger values of $\lambda$, dependence on the partitioning can become weaker or completely vanish. We note that, for $\lambda_1$, there is still a strong dependence on the partitioning, and the error have significant peaks, although bounded, around $p_k \approx n$. As we will discuss more in detail in the remainder of this section, convergence rate of CoCoA heavily depends on $\lambda$ and the algorithm has not yet converged for all values of $\lambda$ in Figure 7.

We now focus on the transient behaviour of the regularized CoCoA algorithm. In Figure 8, we plot the average generalization error with $\lambda = 10^{-3}$ and $\sigma_w^2 = 0$ over the iterations $t$ for five partitioning choices $p_1 = \{1, 25, 50, 75, 100\}$ as well as that of the centralized regularized least-squares solution in (7). We note two main effects: Firstly, for all choices of $p_1$, the generalization error converges to that of the centralized LS solution, i.e., $\approx 0.63$. This is consistent with the fact that the regularized CoCoA converges to the centralized LS solution with the same regularization, see the discussions in Section IV. Secondly, the partitioning can greatly affect the convergence rate of the regularized algorithm. For values of $p_1$ closer to $n = 75$, the convergence rate is much slower than for values smaller than $p_1 = 25$. For instance, for $p_1 = 25$ the MSE reaches $\approx 0.66$ at $t = 5 \cdot 10^4$, while it takes almost twice as many iterations to reach the same MSE for $p_1 = 75$.

We conclude this discussion by highlighting the manner in which our results presented in Theorem 1–4 are relevant for the regularized setting, although they are derived for the unregularized version of CoCoA. A key observation in Figure 8 is that while the partitioning does not affect the performance after convergence, it has a considerable effect on the conditioning of the problem; which in turn affects the practical performance of the algorithm significantly. If one would set $T$ in the order of $10^4$, which is a very large number of iterations in terms of convergence in relation to the unregularized setting, then the choice of partitioning can severely affect the generalization performance of the final
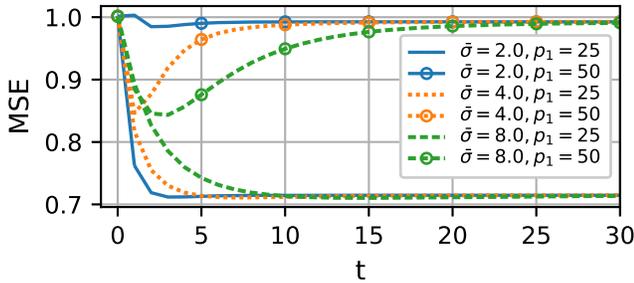
Fig. 9: Convergence of COCOA for three choices of $\bar{\sigma}$.

solution $\hat{x}^T$, as seen in Figure 7. In other words, if we are not willing or able to run the algorithm for a substantially larger number of iterations compared to the unregularized case, then we should also avoid having $p_k$ close to $n$ for the regularized setting, in order to avoid a large generalization error, similarly as we should avoid $p_k$ close to $n$ in the unregularized setting.

### G. Hyperparameter Study for COCOA parameters

We now conduct a hyperparameter study for the aggregation parameter $\bar{\varphi}$ and the subproblem parameter $\bar{\sigma}$ of COCOA. Choosing these parameters as $\bar{\varphi} \in (0,1]$ and $\bar{\sigma} \geq \bar{\varphi}K$ facilitates formal convergence guarantees [25, Sec. 3.1]. As discussed in Section IV, with $\lambda = 0$, it is only the ratio $\bar{\varphi}/\bar{\sigma}$ that affects our expressions, hence our previous results cover all admissible values of $\bar{\varphi}$ and $\bar{\sigma}$ as long as $\bar{\varphi}/\bar{\sigma} = K$. To study scenarios that have not been covered with the previous plots but still with formal convergence guarantees, we use $\bar{\varphi} = 1$ and $\bar{\sigma} \in \{4, 8\}$ together with the earlier case of $\bar{\sigma} = K = 2$ for comparison. In Figure 9, we present the convergence under isotropic Gaussian setting, with $\lambda = 0$, $\boldsymbol{w} = 0$ and $p_1 \in \{25, 50\}$. While the error converges to the same value for all choices of $\bar{\sigma}$, the rate of convergence is slower for $\bar{\sigma}$ with $\bar{\sigma} > K\bar{\varphi}$, as expected. Hence, these results support the expectation that the parameters chosen for the previous studies in this section provide relatively fast convergence.

## IX. DISCUSSIONS

We now discuss the practical guidelines our results provide. Our results emphasize the relation between the partitioning of the model unknowns over the network and the generalization error. Although the training error is low and at the same level as that of the centralized solution, a low generalization error is not guaranteed when the number of unknowns $p_k$ at any node is close to the number of observations $n$ in the training data. Furthermore, if any $p_k$ is close to $n$, then the generalization performance can vary significantly over different realizations of the training data. Hence, the partitioning should not be chosen so that $p_k$ is close to $n$, for any node.

Explicit regularization, i.e., $\lambda > 0$ in (2), can improve the generalization error of the distributed scheme. Nevertheless, the choice of the regularization parameter $\lambda$ is not straightforward. The algorithm can need significantly different numbers of iterations (such as $10^4$ times more) for different values of $\lambda$. If $\lambda$ is chosen too small, then the generalization error can still be relatively large if any $p_k$ is close to $n$, compared to other possible data partitionings. Hence, for a fixed number

of iterations, one should choose a large enough $\lambda$ in order to guarantee that the generalization error does not depend on the data partitioning over nodes.

While our results are restricted to the convex setting, extensions into the non-convex formulations are considered an important line of future research. A key step for this challenging set-up could involve relaxation of the definition of convergence, as done in the centralized case [47], [48].

## X. CONCLUSIONS

We have focused on the generalization error associated with solutions produced by the distributed learning algorithm COCOA for the linear regression problem. We have presented upper bounds on the generalization error that hold with high probability for isotropic Gaussian, correlated Gaussian and sub-gaussian data. We have compared our probabilistic bounds with the results on the expected generalization error. With our numerical results, we have illustrated the generalization performance of the algorithm with both synthetic and real data.

In existing works, there is a lack of efforts for determining how the generalization error in the distributed setting can be affected by the algorithm design. Here, we have addressed this gap by providing bounds that characterize how the partitioning of the model's unknowns over the nodes in the network affects the generalization error. Our results provide guidelines on how to partition the model over the network in order to avoid potential pitfalls. Our results show that if the number of unknowns $p_k$ in any node is close to the total number of observations $n$, then the generalization error can be very large, even though the training error is small. Hence, in order to obtain a good generalization performance, the number of unknowns in any node should be chosen to be sufficiently larger or smaller than the number of observations if possible. If one has to operate with a node with $p_k \approx n$, regularization can be used to significantly dampen the generalization error. On the other hand, choosing the regularization parameter is not straightforward. If the regularization parameter is too small, COCOA needs a relatively large amount of iterations in order to mitigate the effect that the partitioning has on the generalization error.

Extensions of our results to the fully decentralized scenarios as well model misspecification are considered as important directions for future work.

## XI. APPENDIX

### A. Preliminaries

This section provides a collection of properties that are used frequently in different proofs:

(a) Given $c_k$, $k = 1, \ldots, K$, with $\mathbb{P}(c_k) \geq 1 - \rho_k$, the probability of intersection can be bounded as

$$\mathbb{P}\left(\bigcap_{k=1}^{K} c_k\right) \geq 1 - \sum_{k=1}^{K} \rho_k. \tag{61}$$

(b) Partition a symmetric matrix $\boldsymbol{M} \succeq 0$ as $\boldsymbol{M} = \begin{bmatrix} \boldsymbol{M}_{11} & \boldsymbol{M}_{12} \\ \boldsymbol{M}_{12}^{\mathrm{T}} & \boldsymbol{M}_{22} \end{bmatrix}$. Then, $\|\boldsymbol{M}\| \leq \|\boldsymbol{M}_{11}\| + \|\boldsymbol{M}_{22}\|$.

(c) Let $\boldsymbol{A}$, $\boldsymbol{B}$ be two real square matrices. If $\|\boldsymbol{A} - \boldsymbol{B}\| \leq q$, then $\sigma_{\min}(\boldsymbol{A}) \geq \sigma_{\min}(\boldsymbol{B}) - q$.

(d) Let $\mathbf{\Lambda} \in \mathbb{R}^{p \times p}$ be diagonal with the entries $\mu_p \geq \cdots \geq \mu_1 \geq 0$, and let $\mathbf{U} \in \mathbb{R}^{p \times p}$ be unitary. Let $\mathbf{A} \in \mathbb{R}^{n \times p}$ have the following decomposition in terms of another matrix $\mathbf{Z} \in \mathbb{R}^{n \times p}$ as $\mathbf{A} = \mathbf{Z}\mathbf{\Lambda}^{1/2}\mathbf{U}^T$. Then $\sigma_{\min}^2(\mathbf{A}) \geq \mu_1 \lambda_{\min}(\mathbf{Z}\mathbf{Z}^{\mathrm{T}})$.

(e) Let $\mathbf{M} \in \mathbb{R}^{n \times p}$ be generated as $\mathbf{M} = \mathbf{Z}\mathbf{\Lambda}^{1/2}\mathbf{U}^{\mathrm{T}}$, where $\mathbf{Z} \in \mathbb{R}^{n \times p}$ has i.i.d. entries with $\mathcal{N}(0,1)$, $\mathbf{\Lambda} \in \mathbb{R}^{p \times p}$ is diagonal and positive definite, and $\mathbf{U} \in \mathbb{R}^{p \times p}$ is unitary. Then, $\mathbf{M}^+\mathbf{M} = \mathbf{I}_p$ and $\mathbf{M}\mathbf{M}^+ = \mathbf{I}_n$ with probability (w.p.) one, if $n \geq p$ or $p \geq n$, respectively.

We now present the proofs for these properties.

*1) Proof of Property (a):* We observe that

$$\mathbb{P}\left(\bigcap_{k=1}^K c_k\right) = 1 - \mathbb{P}\left(\left(\bigcap_{k=1}^K c_k\right)^c\right) = 1 - \mathbb{P}\left(\bigcup_{k=1}^K c_k^c\right) \quad (62)$$

Using the union bound we have that $\mathbb{P}\left(\bigcup_{k=1}^K c_k^c\right) \leq \sum_{k=1}^K \mathbb{P}(c_k^c)$. By definition, $\mathbb{P}(c_k) = 1 - \mathbb{P}(c_k^c) \geq 1 - \rho_k$, or equivalently, $\mathbb{P}(c_k^c) \leq \rho_k$. Hence, $\mathbb{P}\left(\bigcup_{k=1}^K c_k^c\right) \leq \sum_{k=1}^K \rho_k$. Using this inequality together with (62) yields to the desired inequality after re-arranging the terms.

*2) Proof of Property (b):* See [49, Proposition 8.3].

*3) Proof of Property (c):* Let $\mathbf{u}$ be any vector such that $\|\mathbf{u}\| = 1$. Combining $\|\mathbf{A} - \mathbf{B}\| = \|\mathbf{A} - \mathbf{B}\|\|\mathbf{u}\| \geq \|\mathbf{A}\mathbf{u} - \mathbf{B}\mathbf{u}\|$ and $\|\mathbf{A} - \mathbf{B}\| \leq q$, we obtain $q \geq \|\mathbf{A}\mathbf{u} - \mathbf{B}\mathbf{u}\|$. Applying reverse triangle inequality, we have $q \geq |\|\mathbf{A}\mathbf{u}\| - \|\mathbf{B}\mathbf{u}\||$ which yields to $q \geq \|\mathbf{A}\mathbf{u}\| - \|\mathbf{B}\mathbf{u}\| \geq -q$. Rearranging the right-hand side inequality, we obtain $\|\mathbf{A}\mathbf{u}\| \geq \|\mathbf{B}\mathbf{u}\| - q$. Let $(\lambda_{\min}(\mathbf{A}^{\mathrm{T}}\mathbf{A}), \bar{\mathbf{u}})$ with $\|\bar{\mathbf{u}}\| = 1$ be the eigenpair corresponding to the smallest eigenvalue of $\mathbf{A}^{\mathrm{T}}\mathbf{A}$. Then,

$$\|\mathbf{A}\bar{\mathbf{u}}\| = \sqrt{\bar{\mathbf{u}}^{\mathrm{T}}\mathbf{A}^{\mathrm{T}}\mathbf{A}\bar{\mathbf{u}}} = \sqrt{\lambda_{\min}(\mathbf{A}^{\mathrm{T}}\mathbf{A})} \quad (63)$$
$$\geq \sqrt{\bar{\mathbf{u}}^{\mathrm{T}}\mathbf{B}^{\mathrm{T}}\mathbf{B}\bar{\mathbf{u}}} - q \geq \sqrt{\lambda_{\min}(\mathbf{B}^{\mathrm{T}}\mathbf{B})} - q. \quad (64)$$

Note that $\sqrt{\lambda_{\min}(\mathbf{M}^{\mathrm{T}}\mathbf{M})} = \sigma_{\min}(\mathbf{M})$ for any $\mathbf{M} \in \mathbb{R}^{n \times n}$, so $\sigma_{\min}(\mathbf{A}) \geq \sigma_{\min}(\mathbf{B}) - q$.

*4) Proof of Property (d):* Writing $\mathbf{A}\mathbf{A}^{\mathrm{T}}$ in terms of $\mathbf{Z}$

$$\mathbf{A}\mathbf{A}^{\mathrm{T}} = \mathbf{Z}\mathbf{\Lambda}^{1/2}\mathbf{U}^{\mathrm{T}}\mathbf{U}\mathbf{\Lambda}^{1/2}\mathbf{Z}^{\mathrm{T}} = \mathbf{Z}\mathbf{\Lambda}\mathbf{Z}^{\mathrm{T}}, \quad (65)$$

together with adding and subtracting $\mu_1 \mathbf{Z}\mathbf{Z}^{\mathrm{T}}$, we obtain

$$\mathbf{A}\mathbf{A}^{\mathrm{T}} = \mu_1 \mathbf{Z}\mathbf{Z}^{\mathrm{T}} + \mathbf{Z}(\mathbf{\Lambda} - \mu_1 \mathbf{I}_p)\mathbf{Z}^{\mathrm{T}}. \quad (66)$$

Let $(\lambda_{\min}(\mathbf{A}\mathbf{A}^{\mathrm{T}}), \mathbf{v})$ be the eigenpair corresponding to the smallest eigenvalue of $\mathbf{A}\mathbf{A}^{\mathrm{T}}$. We now evaluate $\mathbf{v}^{\mathrm{T}}\mathbf{A}\mathbf{A}^{\mathrm{T}}\mathbf{v}$ using (66) to obtain

$$\lambda_{\min}(\mathbf{A}\mathbf{A}^{\mathrm{T}}) = \mu_1 \mathbf{v}^{\mathrm{T}}\mathbf{Z}\mathbf{Z}^{\mathrm{T}}\mathbf{v} + \mathbf{v}^{\mathrm{T}}\mathbf{Z}(\mathbf{\Lambda} - \mu_1 \mathbf{I}_p)\mathbf{Z}^{\mathrm{T}}\mathbf{v}, \quad (67)$$
$$\geq \mu_1 \lambda_{\min}(\mathbf{Z}\mathbf{Z}^{\mathrm{T}}), \quad (68)$$

where we have used that $\mathbf{v}^{\mathrm{T}}\mathbf{Z}(\mathbf{\Lambda} - \mu_1 \mathbf{I}_p)\mathbf{Z}^{\mathrm{T}}\mathbf{v} \geq 0$. Note that $\sigma_{\min}^2(\mathbf{A}) \geq \sigma_{\min}(\mathbf{A}\mathbf{A}^{\mathrm{T}}) = \lambda_{\min}(\mathbf{A}\mathbf{A}^{\mathrm{T}})$, and we have $\sigma_{\min}^2(\mathbf{A}) \geq \mu_1 \lambda_{\min}(\mathbf{Z}\mathbf{Z}^{\mathrm{T}})$.

*5) Proof of Property (e):* By [39, eqn (3.2)], $\mathbf{Z}$ is full rank w.p. 1. Hence, $\mathbf{M}$ is full rank since it is the product of full rank matrices. Thus, with $n \geq p$, $\mathbf{M}^{\mathrm{T}}\mathbf{M} \in \mathbb{R}^{p \times p}$ is full rank, i.e. invertible. Hence $\mathbf{M}^+\mathbf{M} = (\mathbf{M}^{\mathrm{T}}\mathbf{M})^+\mathbf{M}^{\mathrm{T}}\mathbf{M} = \mathbf{I}_p$. A similar line of argument holds for $p \geq n$ with $\mathbf{M}\mathbf{M}^+ = \mathbf{I}_n$.

*B. Proof of Lemma 1*

Expanding $\mathbf{B}^2$, we obtain the following

$$\mathbf{B}^2 = \left(\mathbf{I}_p - \frac{1}{K}\bar{\mathbf{A}}\mathbf{A}\right)^2 = \mathbf{I}_p + \frac{1}{K^2}(\bar{\mathbf{A}}\mathbf{A})^2 - \frac{2}{K}\bar{\mathbf{A}}\mathbf{A}, \quad (69)$$

Note that $(\bar{\mathbf{A}}\mathbf{A})^2 = \bar{\mathbf{A}}\mathbf{A}\bar{\mathbf{A}}\mathbf{A}$, where $\mathbf{A}\bar{\mathbf{A}} = \sum_{k=1}^K \mathbf{A}_k\mathbf{A}_k^+$. The rows of $\mathbf{A}_k$ have positive definite covariance matrices.

Hence, by Property (e) under $n \leq p_k$, we have $\mathbf{A}_k\mathbf{A}_k^+ = \mathbf{I}_n$, and $\mathbf{A}\bar{\mathbf{A}} = K\mathbf{I}_n$. Hence, $(\bar{\mathbf{A}}\mathbf{A})^2 = K\bar{\mathbf{A}}\mathbf{A}$. Then, we have

$$\mathbf{B}^2 = \mathbf{I}_p + \frac{1}{K^2}K\bar{\mathbf{A}}\mathbf{A} - \frac{2}{K}\bar{\mathbf{A}}\mathbf{A} = \mathbf{I}_p - \frac{1}{K}\bar{\mathbf{A}}\mathbf{A} = \mathbf{B}. \quad (70)$$

*C. Proof of Theorem 1*

Using the triangle inequality and (14), we obtain

$$\|\mathbf{B}\| \leq 1 + \frac{1}{K}\|\bar{\mathbf{A}}\mathbf{A}\|. \quad (71)$$

We now present the following algebraic property of $\bar{\mathbf{A}}\mathbf{A}$ which holds regardless of the distribution of $\mathbf{A}$:

**Lemma 6.** *For any matrix $\mathbf{A} = [\mathbf{A}_1, \cdots, \mathbf{A}_K] \in \mathbb{R}^{n \times p}$ and $\bar{\mathbf{A}} = [\mathbf{A}_1^+; \cdots; \mathbf{A}_K^+] \in \mathbb{R}^{p \times n}$, the following bound holds:*

$$\|\bar{\mathbf{A}}\mathbf{A}\|^2 \leq K + \sum_{k=1}^K \sum_{\substack{i=1 \\ i \neq k}}^K \frac{\sigma_{\max}^2(\mathbf{A}_k)}{\sigma_{\min+}^2(\mathbf{A}_i)}, \quad (72)$$

*where $\sigma_{\min+}(\mathbf{A}_i)$ denotes the smallest non-zero singular value of $\mathbf{A}_i$.*

Proof: See Section XI-D. Note that we in the subsequent sections use that $\frac{1}{\sigma_{min+}} \leq \frac{1}{\sigma_{\min}}$, but when $\sigma_{\min} \to 0$, this upper bound is uninformative. The aim of our results is to find bounds on the minimum singular value which is away from zero.

The result of Theorem 1 is obtained by combining (71), (72) and the bounds on the extreme singular values of $\mathbf{A}_k$. In particular, denote the event that the singular value inequalities given in Lemma 2 holds for the partition $\mathbf{A}_k$ as $c_k$, i.e. $c_k = c_k^L \cap c_k^U$, where $c_k^L = \{r_k^L(q_k) \leq \sigma_{\min}(\mathbf{A}_k)\}$ and $c_k^U = \{\sigma_{\max}(\mathbf{A}_k) \leq r_k^U(q_k)\}$, with $r_k^L(q_k) = \sqrt{r_{\max,k}} - \sqrt{r_{\min,k}} - q_k$ and $r_k^U(q_k) = \sqrt{r_{\max,k}} + \sqrt{r_{\min,k}} + q_k$. Note that $c_k$ can be rearranged as follows

$$c_k = \left\{\frac{1}{\sigma_{\min}(\mathbf{A}_k)} \leq \frac{1}{r_k^L(q_k)} \cap \sigma_{\max}(\mathbf{A}_k) \leq r_k^U(q_k)\right\}. \quad (73)$$

For any $k \neq i$, $c_k$ and $c_i$ are statistically independent since the entries of $\mathbf{A}$ are statistically independent. Hence,

$$\mathbb{P}\left(\bigcap_{k=1}^K c_k\right) \geq \prod_{k=1}^K (1 - 2e^{-q_k^2/2})_+ \quad (74)$$

Therefore, (72), (73) and (74) yields to

$$\|\bar{\mathbf{A}}\mathbf{A}\|^2 \leq K + \sum_{k=1}^K \sum_{\substack{i=1 \\ i \neq k}}^K \frac{(\sqrt{r_{\max,k}} + \sqrt{r_{\min,k}} + q_k)^2}{(\sqrt{r_{\max,i}} - \sqrt{r_{\min,i}} - q_i)^2}, \quad (75)$$

with probability at least $\prod_{k=1}^K (1 - 2e^{-q_k^2/2})_+^K$. The desired result in Theorem 1 is obtained by combining (75) and (71). We bound $q_i < \sqrt{r_{\max,i}} - \sqrt{r_{\min,i}}$, so that the bound on $\sigma_{\min}^2(\mathbf{A}_i)$ is informative, i.e., strictly greater than zero.

*D. Proof of Lemma 6*

The matrices $\bar{\mathbf{A}}\mathbf{A}$ and $\mathbf{A}^{\mathrm{T}}\bar{\mathbf{A}}^{\mathrm{T}}\bar{\mathbf{A}}\mathbf{A}$ can be seen as matrices consisting of $K \times K$ blocks. The $(k,j)^{\mathrm{th}}$ block of $\bar{\mathbf{A}}\mathbf{A}$ is of size $p_k \times p_j$ and given by

$$[\bar{\mathbf{A}}\mathbf{A}]_{k,j} = \mathbf{A}_k^+ \mathbf{A}_j. \quad (76)$$

The $(k,j)^{\mathrm{th}}$ block of $\mathbf{A}^{\mathrm{T}}\bar{\mathbf{A}}^{\mathrm{T}}\bar{\mathbf{A}}\mathbf{A}$ is of size $p_k \times p_j$ and given by

$$[\mathbf{A}^{\mathrm{T}}\bar{\mathbf{A}}^{\mathrm{T}}\bar{\mathbf{A}}\mathbf{A}]_{k,j} = \mathbf{A}_k^{\mathrm{T}}\left(\sum_{i=1}^K \mathbf{A}_i^{+\mathrm{T}}\mathbf{A}_i^+\right)\mathbf{A}_j. \quad (77)$$

Hence, we have

$$\|\bar{\boldsymbol{A}}\boldsymbol{A}\|^2 = \|\boldsymbol{A}^{\mathrm{T}}\bar{\boldsymbol{A}}^{\mathrm{T}}\bar{\boldsymbol{A}}\boldsymbol{A}\| \leq \sum_{k=1}^{K}\|\boldsymbol{A}_k^{\mathrm{T}}(\sum_{i=1}^{K}\boldsymbol{A}_i^{+\mathrm{T}}\boldsymbol{A}_i^{+})\boldsymbol{A}_k\| \quad (78)$$

$$\leq \sum_{k=1}^{K}\sum_{i=1}^{K}\|\boldsymbol{A}_k^{\mathrm{T}}\boldsymbol{A}_i^{+\mathrm{T}}\boldsymbol{A}_i^{+}\boldsymbol{A}_k\|, \quad (79)$$

where we obtained (78) using property (b) of Section XI-A repeatedly on the blocks of $\boldsymbol{A}^{\mathrm{T}}\bar{\boldsymbol{A}}^{\mathrm{T}}\bar{\boldsymbol{A}}\boldsymbol{A}$ on the diagonal (i.e. $k = j$). In (79), we used the triangle inequality.

We now consider the individual terms in the double summation of (79). For the terms with $k = i$, consider the s.v.d. $\boldsymbol{A}_k = \boldsymbol{U}_k\boldsymbol{\Lambda}_k\boldsymbol{V}_k^{\mathrm{T}}$, where $\boldsymbol{\Lambda}_k \in \mathbb{R}^{n \times p_k}$ is the (possibly rectangular) diagonal matrix of singular values, and $\boldsymbol{U}_k \in \mathbb{R}^{n \times n}$ and $\boldsymbol{V}_k \in \mathbb{R}^{p_k \times p_k}$ are unitary. Hence, we have $\boldsymbol{A}_k^{\mathrm{T}}\boldsymbol{A}_k^{+\mathrm{T}}\boldsymbol{A}_k^{+}\boldsymbol{A}_k = \boldsymbol{V}_k\boldsymbol{\Lambda}_k^{\mathrm{T}}\boldsymbol{\Lambda}_k^{+\mathrm{T}}\boldsymbol{\Lambda}_k^{+}\boldsymbol{\Lambda}_k\boldsymbol{V}_k^{\mathrm{T}} = \boldsymbol{V}_k\boldsymbol{R}_k\boldsymbol{V}_k^{\mathrm{T}}$, where $\boldsymbol{R}_k \in \mathbb{R}^{p_k \times p_k}$ is a diagonal matrix with ones and zeroes on its diagonal. Hence,

$$\|\boldsymbol{A}_k^{\mathrm{T}}\boldsymbol{A}_k^{+\mathrm{T}}\boldsymbol{A}_k^{+}\boldsymbol{A}_k\| = 1. \quad (80)$$

For the terms with $i \neq k$, we use that the spectral norm is submultiplicative and self-adjoint [42, Sec. 5.6] to obtain

$$\|\boldsymbol{A}_k^{\mathrm{T}}\boldsymbol{A}_i^{+\mathrm{T}}\boldsymbol{A}_i^{+}\boldsymbol{A}_k\| \leq \|\boldsymbol{A}_k\|^2\|\boldsymbol{A}_i^{+}\|^2 = \frac{\sigma_{\max}^2(\boldsymbol{A}_k)}{\sigma_{\min+}^2(\boldsymbol{A}_i)} \quad (81)$$

where we have used $\|\boldsymbol{A}_k\|^2 = \sigma_{\max}^2(\boldsymbol{A}_k)$ and the property $\|\boldsymbol{A}_i^{+}\| = \frac{1}{\sigma_{\min+}(\boldsymbol{A}_i)}$ where $\sigma_{\min+}(\boldsymbol{A}_i)$ is the smallest non-zero singular value of $\boldsymbol{A}_i$. Note that $\boldsymbol{A}_i^{+}$ can be written as $\boldsymbol{A}_i^{+} = \boldsymbol{V}_i\boldsymbol{\Lambda}_i^{+}\boldsymbol{U}_i^{\mathrm{T}}$. Hence, $\|\boldsymbol{A}_i^{+}\| = \sigma_{\min+}^{-1}(\boldsymbol{A}_i)$. Now combining (80), (81) and (79), we obtain (72).

### E. Proof of Lemma 3

Let $\boldsymbol{M} = \boldsymbol{B}^{\mathrm{T}}\boldsymbol{B} \in \mathbb{R}^{p \times p}$, and denote its $K$ blocks on the diagonal as $\boldsymbol{M}_{kk} \in \mathbb{R}^{p_k \times p_k}$, $\forall k$. Using Property (b) of Section XI-A,

$$\|\boldsymbol{M}\| \leq \sum_{k=1}^{K}\|\boldsymbol{M}_{kk}\|. \quad (82)$$

With $\boldsymbol{B}$ from (14), we have $\boldsymbol{M} = \boldsymbol{B}^{\mathrm{T}}\boldsymbol{B} = \boldsymbol{I}_p + \frac{1}{K^2}\boldsymbol{A}^{\mathrm{T}}\bar{\boldsymbol{A}}^{\mathrm{T}}\bar{\boldsymbol{A}}\boldsymbol{A} - \frac{1}{K}(\bar{\boldsymbol{A}}\boldsymbol{A})^{\mathrm{T}} - \frac{1}{K}\bar{\boldsymbol{A}}\boldsymbol{A}$. Using (76) and (77), we find that the blocks on the diagonal of $\boldsymbol{M}$, $\boldsymbol{M}_{kk}$ can be decomposed such that $\boldsymbol{M}_{kk} = \boldsymbol{C}_{kk} + \boldsymbol{D}_{kk}$, where $\boldsymbol{C}_{kk} = (\boldsymbol{I}_{p_k} - \frac{1}{K}\boldsymbol{A}_k^{+}\boldsymbol{A}_k)^2$ and $\boldsymbol{D}_{kk} = \frac{1}{K^2}\boldsymbol{A}_k^{\mathrm{T}}(\sum_{\substack{i=1 \\ i\neq k}}^{K}\boldsymbol{A}_i^{+\mathrm{T}}\boldsymbol{A}_i^{+})\boldsymbol{A}_k$.

By Property (e), under $n \geq p_k$, we have $\boldsymbol{A}_k^{+}\boldsymbol{A}_k = \boldsymbol{I}_{p_k}$, hence $\boldsymbol{C}_{kk} = \boldsymbol{I}_{p_k}(\frac{K-1}{K})^2$ and $\|\boldsymbol{C}_{kk}\| = (\frac{K-1}{K})^2$. From the proof of Theorem 1, we have with probability at least $\bar{\rho}_G = \prod_{k=1}^{K}(1 - 2e^{-q_k^2/2})_+$:

$$\|\boldsymbol{D}_{kk}\| = \frac{1}{K^2}\|\boldsymbol{A}_k^{\mathrm{T}}(\sum_{\substack{i=1 \\ i\neq k}}^{K}\boldsymbol{A}_i^{+\mathrm{T}}\boldsymbol{A}_i^{+})\boldsymbol{A}_k\| \leq \frac{1}{K^2}\sum_{\substack{i=1 \\ i\neq k}}^{K}\bar{\gamma}_{k,i}, \quad (83)$$

for $k = 1, \ldots, K$, with $\bar{\gamma}_{k,i} = \frac{(\sqrt{r_{\max,k}} + \sqrt{r_{\min,i}} + q_k)^2}{(\sqrt{r_{\max,i}} - \sqrt{r_{\min,i}} - q_i)^2}$. Hence, using (82) and $\|\boldsymbol{M}_{kk}\| = \|\boldsymbol{C}_{kk} + \boldsymbol{D}_{kk}\| \leq \|\boldsymbol{C}_{kk}\| + \|\boldsymbol{D}_{kk}\|$, we have

$$\|\boldsymbol{B}^{\mathrm{T}}\boldsymbol{B}\| \leq \frac{(K-1)^2}{K} + \frac{1}{K^2}\sum_{k=1}^{K}\sum_{\substack{i=1 \\ i\neq k}}^{K}\bar{\gamma}_{k,i}. \quad (84)$$

To conclude the proof, use the fact that $\|\boldsymbol{B}\| = \sqrt{\|\boldsymbol{B}^{\mathrm{T}}\boldsymbol{B}\|}$.

### F. Proof of Lemma 4

By (19) and $\boldsymbol{\Sigma} = \boldsymbol{I}_p$, we have

$$\kappa(\hat{\boldsymbol{x}}^1) = \|\boldsymbol{B}\boldsymbol{x} - \frac{1}{K}\bar{\boldsymbol{A}}\boldsymbol{w}\|^2 \quad (85)$$

$$= \|\boldsymbol{B}\boldsymbol{x}\|^2 + \|\frac{1}{K}\bar{\boldsymbol{A}}\boldsymbol{w}\|^2 - \frac{2}{K}\boldsymbol{x}^{\mathrm{T}}\boldsymbol{B}^{\mathrm{T}}\bar{\boldsymbol{A}}\boldsymbol{w}. \quad (86)$$

Since $\boldsymbol{w}$ is zero-mean and statistically independent with $\boldsymbol{A}$, we have

$$\mathbb{E}_{\boldsymbol{A},\boldsymbol{w}}[\kappa(\hat{\boldsymbol{x}}^1)] = \mathbb{E}_{\boldsymbol{A}}[\|\boldsymbol{B}\boldsymbol{x}\|^2] + \frac{1}{K^2}\mathbb{E}_{\boldsymbol{A},\boldsymbol{w}}[\|\bar{\boldsymbol{A}}\boldsymbol{w}\|^2]. \quad (87)$$

The first term is evaluated in [32, Thm. 1]. We now focus on the second term. Using the fact that $\boldsymbol{A}$ and $\boldsymbol{w}$ are statistically independent, we find that

$$\mathbb{E}_{\boldsymbol{A},\boldsymbol{w}}[\|\bar{\boldsymbol{A}}\boldsymbol{w}\|^2] = \mathbb{E}_{\boldsymbol{w}}[\boldsymbol{w}^{\mathrm{T}}\mathbb{E}_{\boldsymbol{A}}[\bar{\boldsymbol{A}}^{\mathrm{T}}\bar{\boldsymbol{A}}]\boldsymbol{w}], \quad (88)$$

where $\bar{\boldsymbol{A}}^{\mathrm{T}}\bar{\boldsymbol{A}} = \sum_{k=1}^{K}(\boldsymbol{A}_k\boldsymbol{A}_k^{\mathrm{T}})^{+}$. From [50] we have that $\mathbb{E}_{\boldsymbol{A}}[\bar{\boldsymbol{A}}^{\mathrm{T}}\bar{\boldsymbol{A}}] = \frac{1}{n}\sum_{k=1}^{K}\gamma_k\boldsymbol{I}_n$, with $\gamma_k$ as in (39). Hence,

$$\mathbb{E}_{\boldsymbol{A},\boldsymbol{w}}[\|\bar{\boldsymbol{A}}\boldsymbol{w}\|^2] = \frac{1}{n}\sum_{k=1}^{K}\gamma_k\mathbb{E}_{\boldsymbol{w}}[\boldsymbol{w}^{\mathrm{T}}\boldsymbol{w}] = \frac{1}{n}\sum_{k=1}^{K}\gamma_k n\sigma_w^2. \quad (89)$$

Combining (89), [32, Thm. 1] and (87) concludes the proof.

### G. Proof of Theorem 2

We first consider an intermediate general result on sub-gaussian variables. Background information on sub-gaussian variables and the notation can be found in Section VII-A.

**Lemma 7.** *Let $\boldsymbol{M} \in \mathbb{R}^{n \times p}$ have rows i.i.d. with $\mathcal{S}(\boldsymbol{\Sigma})$, $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$. The following bounds hold with probability at least $1 - 2e^{-q}$, $q \geq 0$:*

$$n\sigma_{\min}(\boldsymbol{\Sigma}) - \ell(q) \leq \sigma_{\min}^2(\boldsymbol{M})$$
$$\leq \sigma_{\max}^2(\boldsymbol{M}) \leq n\sigma_{\max}(\boldsymbol{\Sigma}) + \ell(q), \quad (90)$$

*where*

$$\ell(q) = CL^2\left(\sqrt{\frac{p+q}{n}} + \frac{p+q}{n}\right)n\sigma_{\max}(\boldsymbol{\Sigma}), \quad (91)$$

*where $C$ is an absolute constant, and $L \geq 1$ is constant such that $\psi_{\boldsymbol{m}}(\boldsymbol{\Sigma}, \boldsymbol{h}) \leq L\sqrt{\boldsymbol{h}^{\mathrm{T}}\boldsymbol{\Sigma}\boldsymbol{h}}$, $\forall \boldsymbol{h} \in \mathbb{R}^{p \times 1}$ where $\boldsymbol{m}$ comes from the same distribution as the rows of $\boldsymbol{M}$.*

Proof: See Section XI-H.

In particular, we have the following for the Gaussian case:

**Lemma 8.** *Let $\boldsymbol{M} \in \mathbb{R}^{n \times p}$ have rows i.i.d. with $\mathcal{N}(0, \boldsymbol{\Sigma})$, $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$. Then (90) holds with $L = \sqrt{8/3}$ in (91), with probability at least $1 - 2e^{-q}$, $q \geq 0$. Additionally, the following holds with probability at least $1 - 2e^{-\bar{q}^2/2}$, $\bar{q} \geq 0$:*

$$\sigma_{\min}(\boldsymbol{M}) \geq \sqrt{\sigma_{\min}(\boldsymbol{\Sigma})}(\sqrt{p} - \sqrt{n} - \bar{q}). \quad (92)$$

Proof: See Section XI-I.

To prove Theorem 2, we apply Lemma 6 and Lemma 8 to (71). Note that the rows of $\boldsymbol{A}_k$ are i.i.d. with $\mathcal{N}(0, \boldsymbol{\Sigma}_k)$, where $\boldsymbol{\Sigma}_k$ is the $k^{\text{th}}$ principal submatrix of $\boldsymbol{\Sigma}$. Applying (90) for $\boldsymbol{A}_k$, we have

$$n\sigma_{\min}(\boldsymbol{\Sigma}_k) - \ell_k(q_k) \leq \sigma_{\min}^2(\boldsymbol{A}_k)$$
$$\leq \sigma_{\max}^2(\boldsymbol{A}_k) \leq n\sigma_{\max}(\boldsymbol{\Sigma}_k) + \ell_k(q_k), \quad (93)$$

with probability at least $1 - 2e^{-q_k}$, $q_k \geq 0$, where $\ell_k(q_k) = \frac{8}{3}n\sigma_{\max}(\boldsymbol{\Sigma}_k)C(\sqrt{\frac{p_k+q_k}{n}} + \frac{p_k+q_k}{n})$. Additionally, by (92) the following holds with probability at least $1 - 2e^{-\bar{q}_k^2/2}$, $\bar{q}_k \geq 0$

$$\sigma_{\min}(\boldsymbol{A}_k) \geq \sqrt{\sigma_{\min}(\boldsymbol{\Sigma}_k)}(\sqrt{p_k} - \sqrt{n} - \bar{q}_k). \quad (94)$$

For the lower bounds in Theorem 2, we use (94) for all broad matrices $\boldsymbol{A}_k$. For tall matrices $\boldsymbol{A}_k$, (94) is uninformative, then we use (93). For (46), we use (93) for all matrices $\boldsymbol{A}_k$, $\forall k$, because we use the upper bound on the largest singular value for all $k$. Consider $\mathcal{K} = \{k : n < p_k\}$ which denotes the set of partition indices $k$ for which the matrices $\boldsymbol{A}_k$ are broad. Then, by property (a) in Section XI-A, the bounds in (93) for all $\boldsymbol{A}_k$ and the bound in (92) for $k \in \mathcal{K}$ simultaneously hold with probability at least $1 - \sum_{k=1}^{K} 2e^{-q_k} - \sum_{k \in \mathcal{K}} 2e^{-\bar{q}_k^2/2}$, which concludes the proof.

### H. Proof of Lemma 7

We use the following result from [41]:

**Lemma 9.** [41, Sec. 4.7] *Consider the setting of Lemma 7. Then, the following bound holds with probability at least $1 - 2e^{-q}$, $q \geq 0$:*

$$\|\boldsymbol{M}^{\mathrm{T}}\boldsymbol{M} - n\boldsymbol{\Sigma}\| \leq \ell(q), \qquad (95)$$

*where $\ell(q)$ is defined as in Lemma 7.*

Using Property (c) from Section XI-A, we observe that if $\|\boldsymbol{M}^{\mathrm{T}}\boldsymbol{M} - n\boldsymbol{\Sigma}\| \leq \ell(q)$, then $\sigma_{\min}^2(\boldsymbol{M}) \geq \sigma_{\min}(\boldsymbol{M}^{\mathrm{T}}\boldsymbol{M}) \geq n\sigma_{\min}(\boldsymbol{\Sigma}) - \ell(q)$, which constitutes the lower bound of Lemma 7. To find the upper bound in (90), we apply the reverse triangle inequality to (95) to obtain $\ell(q) \geq \|\boldsymbol{M}^{\mathrm{T}}\boldsymbol{M}\| - n\sigma_{\max}(\boldsymbol{\Sigma})$, and use that $\sigma_{\max}^2(\boldsymbol{M}) = \|\boldsymbol{M}^{\mathrm{T}}\boldsymbol{M}\|$. The upper and lower bounds hold with the same probability as in (95).

### I. Proof of Lemma 8

Using (52), (53), the properties of Gaussian integral and the fact that the rows $\boldsymbol{M}$ are i.i.d. with $\mathcal{N}(0, \boldsymbol{\Sigma})$, we arrive at $\psi_{\boldsymbol{m}}(\boldsymbol{\Sigma}, \boldsymbol{h}) = \sqrt{\frac{8}{3}}\sqrt{\boldsymbol{h}^{\mathrm{T}}\boldsymbol{\Sigma}\boldsymbol{h}}$, hence $L$ can be chosen as $\sqrt{\frac{8}{3}}$. (Details are omitted due to space constraints.) We now derive (92). Denote the s.v.d. of $\boldsymbol{\Sigma}$ as $\boldsymbol{\Sigma} = \boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{U}^T$, where $\boldsymbol{\Lambda} \in \mathbb{R}^{p \times p}$ is the diagonal matrix of singular values and $\boldsymbol{U} \in \mathbb{R}^{p \times p}$ is unitary. Since $\boldsymbol{M}$ has i.i.d. rows with $\mathcal{N}(0, \boldsymbol{\Sigma})$, it can be decomposed as $\boldsymbol{M} = \boldsymbol{Z}\boldsymbol{\Lambda}^{1/2}\boldsymbol{U}^T$, where entries of $\boldsymbol{Z} \in \mathbb{R}^{n \times p}$ are i.i.d. Gaussian with $\mathcal{N}(0, 1)$. Using Property (d) from Section XI-A, we obtain

$$\sigma_{\min}^2(\boldsymbol{M}) \geq \sigma_{\min}(\boldsymbol{\Sigma})\lambda_{\min}(\boldsymbol{Z}\boldsymbol{Z}^{\mathrm{T}}). \qquad (96)$$

Note that if $\boldsymbol{Z}$ is broad, i.e., $n < p$, then $\lambda_{\min}(\boldsymbol{Z}\boldsymbol{Z}^{\mathrm{T}}) = \sigma_{\min}^2(\boldsymbol{Z}) = \sigma_{\min}^2(\boldsymbol{Z}^{\mathrm{T}})$. Applying Lemma 2 to $\boldsymbol{Z}$, we obtain

$$\lambda_{\min}(\boldsymbol{Z}\boldsymbol{Z}^{\mathrm{T}}) \geq (\sqrt{p} - \sqrt{n} - \bar{q})_+^2, \qquad (97)$$

with probability at least $1 - 2e^{-\bar{q}^2/2}$, $\bar{q} \geq 0$. Note that the bound holds even if $n \geq p$, but it is then informative. Using (96) and (97), we obtain the desired result in (92).

### J. Proof of Theorem 3

To find the bound on $\|\boldsymbol{B}\|$ in Theorem 3, we combine (71), Lemma 6, Lemma 7 and the following result:

**Lemma 10.** *Let $\boldsymbol{M} \in \mathbb{R}^{n \times p}$ be generated as $\boldsymbol{M} = \boldsymbol{Z}\boldsymbol{\Lambda}^{1/2}\boldsymbol{U}^{\mathrm{T}}$, where $\boldsymbol{Z} \in \mathbb{R}^{n \times p}$ has entries i.i.d. with $\mathcal{S}(1)$, $\boldsymbol{\Lambda} \in \mathbb{R}^{p \times p}$ is diagonal with nonnegative entries and $\boldsymbol{U} \in \mathbb{R}^{p \times p}$ is unitary. Let $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$ denote the corresponding covariance matrix of the rows in $\boldsymbol{M}$. Under $n < p$, the following bound holds with probability at least $1 - 2e^{-\bar{q}^2}$, $\bar{q} \geq 0$:*

$$\sigma_{\min}^2(\boldsymbol{M}) \geq \sigma_{\min}(\boldsymbol{\Sigma})(\sqrt{p} - CL^2(\sqrt{n} + \bar{q}))_+^2, \qquad (98)$$

*with $C$ and $L$ as in Lemma 7.*

Proof: By Property (d) of Section XI-A, we have

$$\sigma_{\min}^2(\boldsymbol{M}) \geq \sigma_{\min}(\boldsymbol{\Sigma})\lambda_{\min}(\boldsymbol{Z}\boldsymbol{Z}^{\mathrm{T}}). \qquad (99)$$

Note that $\sigma_{\min}^2(\boldsymbol{Z}^{\mathrm{T}}) = \lambda_{\min}(\boldsymbol{Z}\boldsymbol{Z}^{\mathrm{T}})$ for $n < p$. Using (99) with [41, Thm. 4.6.1], we obtain the desired result in (98). $\square$

We now continue with the proof of Theorem 3. Since $\boldsymbol{A}_k = \boldsymbol{Z}_k\boldsymbol{\Lambda}_k^{1/2}\boldsymbol{U}_k$, the rows of $\boldsymbol{A}_k$ are zero-mean subgaussian random vectors with the covariance matrix $\boldsymbol{\Sigma}_k = \boldsymbol{U}_k\boldsymbol{\Lambda}_k\boldsymbol{U}_k^{\mathrm{T}}$ (see Definition VII.1 and VII.2). Hence, (90) holds with $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_k$, $\boldsymbol{M} = \boldsymbol{A}_k$, $\ell(q) = \ell_k(q_k) = CL_k^2\left(\sqrt{\frac{p_k+q_k}{n}} + \frac{p_k+q_k}{n}\right)n\sigma_{\max}(\boldsymbol{\Sigma}_k)$, with probability at least $1 - 2e^{-q_k}$ and $L_k \geq 1$ such that $\psi_{\boldsymbol{a}_{i,k}}(\boldsymbol{\Sigma}_k, \boldsymbol{h}) \leq L_k\sqrt{\boldsymbol{h}^{\mathrm{T}}\boldsymbol{\Sigma}_k\boldsymbol{h}}$, any $\boldsymbol{h} \in \mathbb{R}^{p_k \times 1}$, where $\boldsymbol{a}_{i,k} \sim \mathcal{S}(\boldsymbol{\Sigma}_k)$. By (98), we also have $\sigma_{\min}^2(\boldsymbol{A}_k) \geq \sigma_{\min}(\boldsymbol{\Sigma}_k)(\sqrt{p_k} - CL_k^2(\sqrt{n} + \bar{q}_k))_+^2$, with probability at least $1 - 2e^{-\bar{q}_k^2}$, if $n < p_k$. To find the bound on $\|\boldsymbol{B}\|$ in (58), we plug in the upper bound on $\sigma_{\max}^2(\boldsymbol{A}_k)$ for each $\boldsymbol{A}_k$, and the respective lower bound for $\sigma_{\min}^2(\boldsymbol{A}_k)$, depending on whether $\boldsymbol{A}_k$ is broad or tall. Using $\mathcal{K} = \{k : n < p_k\}$, and Property (a) of Section XI-A, we find the desired probability bound $1 - \sum_{k=1}^{K} 2e^{-q_k} - \sum_{k \in \mathcal{K}} 2e^{-\bar{q}_k^2}$.

### K. Proof of Theorem 4

The proof follows a similar line of argument with the proof of Theorem 3. In particular, we use (90) with $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_k$, $\boldsymbol{M} = \boldsymbol{A}_k$. The probability expression is found using the probability bound for (90) for each $k$, i.e., $1 - 2e^{-q_k}$, and Property (a) of Section XI-A. We omit the details due to space constraints.

### REFERENCES

[1] S. Niknam, H. S. Dhillon, and J. H. Reed, "Federated learning for wireless communications: Motivation, opportunities, and challenges," *IEEE Commun. Mag.*, vol. 58, no. 6, pp. 46–51, 2020.

[2] X. Wang, H. Ishii, L. Du, P. Cheng *et al.*, "Privacy-preserving distributed machine learning via local randomization and ADMM perturbation," *IEEE Trans. Signal Process.*, vol. 68, pp. 4226–4241, 2020.

[3] H. Chen, Y. Ye, M. Xiao, M. Skoglund *et al.*, "Coded Stochastic ADMM for Decentralized Consensus Optimization with Edge Computing," *arXiv:2010.00914*, Oct. 2020.

[4] S. Wang, T. Tuor, T. Salonidis, K. K. Leung *et al.*, "Adaptive federated learning in resource constrained edge computing systems," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 6, pp. 1205–1221, 2019.

[5] L. Bottou, F. E. Curtis, and J. Nocedal, "Optimization methods for large-scale machine learning," *SIAM Review*, vol. 60, no. 2, pp. 223–311, 2018.

[6] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba *et al.*, "Towards federated learning at scale: System design," *Proc. of the 2nd SysML Conf.*, 2019.

[7] J. Dean, G. Corrado, R. Monga, K. Chen *et al.*, "Large scale distributed deep networks," *Adv. Neural Inf. Process. Syst.*, pp. 1223–1231, 2012.

[8] S. M. Kay, *Fundamentals of Stat. Signal Process.* Prentice Hall, 1993.

[9] M. Belkin, D. Hsu, S. Ma, and S. Mandal, "Reconciling modern machine-learning practice and the classical bias–variance trade-off," *Proc. of the Nat. Acad. of Sciences*, vol. 116, no. 32, 2019.

[10] L. Breiman and D. Freedman, "How many variables should be entered in a regression equation?" *J. Amer. Stat. Assoc.*, vol. 78, no. 381, pp. 131–136, 1983.

[11] M. Belkin, D. Hsu, and J. Xu, "Two models of double descent for weak features," *SIAM Journal on Mathematics of Data Science*, vol. 2, no. 4, pp. 1167–1180, 2020.

[12] P. Nakkiran, P. Venkat, S. Kakade, and T. Ma, "Optimal regularization can mitigate double descent," *arXiv:2003.01897*, 2020.

[13] T. Hastie, A. Montanari, S. Rosset, and R. J. Tibshirani, "Surprises in high-dimensional ridgeless least squares interpolation," *arXiv:1903.08560*, 2020.

[14] P. L. Bartlett, P. M. Long, G. Lugosi, and A. Tsigler, "Benign overfitting in linear regression," *arxiv:1906.11300*, 2020.

[15] V. Muthukumar, K. Vodrahalli, V. Subramanian, and A. Sahai, "Harmless interpolation of noisy data in regression," *IEEE J. on Sel. Areas in Inf. Theory*, vol. 1, no. 1, pp. 67–83, 2020.

[16] A. Koppel, F. Y. Jakubiec, and A. Ribeiro, "A saddle point algorithm for networked online convex optimization," *IEEE Trans. Signal Process.*, vol. 63, no. 19, pp. 5149–5164, 2015.

[17] A. S. Bedi, A. Koppel, and K. Rajawat, "Asynchronous online learning in multi-agent systems with proximity constraints," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 5, no. 3, pp. 479–494, 2019.

[18] S. Samar, S. Boyd, and D. Gorinevsky, "Distributed estimation via dual decomposition," *European Control Conference*, 2007.

[19] H. Terelius, U. Topcu, and R. M. Murray, "Decentralized multi-agent optimization via dual decomposition," *18th IFAC World Congress*, vol. 44, no. 1, pp. 11 245–11 251, 2011.

[20] S. Boyd, N. Parikh, E. Chu, B. Peleato *et al.*, *Distributed optimization and statistical learning via the alternating direction method of multipliers*. Now Publishers Inc, 2011.

[21] S. Paternain, S. Lee, M. M. Zavlanos, and A. Ribeiro, "Distributed constrained online learning," *IEEE Trans. Signal Process.*, vol. 68, pp. 3486–3499, 2020.

[22] X. Zhao and A. H. Sayed, "Distributed clustering and learning over networks," *IEEE Trans. Signal Process.*, vol. 63, no. 13, pp. 3285–3300, 2015.

[23] D. Alistarh, D. Grubic, J. Li, R. Tomioka *et al.*, "QSGD: Communication-efficient SGD via gradient quantization and encoding," *Adv. in Neural Inf. Process. Syst.*, pp. 1709–1720, 2017.

[24] S. Magnússon, C. Enyioha, N. Li, C. Fischione *et al.*, "Communication complexity of dual decomposition methods for distributed resource allocation optimization," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 4, pp. 717–732, Aug. 2018.

[25] V. Smith, S. Forte, C. Ma, M. Takáč *et al.*, "COCOA: A general framework for communication-efficient distributed optimization," *J. Mach. Learn. Res.*, vol. 18, no. 1, pp. 8590–8638, 2017.

[26] H. Zhang, J. M. F. Moura, and B. Krogh, "Dynamic field estimation using wireless sensor networks: Tradeoffs between estimation error and communication cost," *IEEE Trans. Signal Process.*, vol. 57, no. 6, pp. 2383–2395, 2009.

[27] U. A. Khan and J. M. F. Moura, "Distributing the Kalman filter for large-scale systems," *IEEE Trans. Signal Process.*, vol. 56, no. 10, pp. 4919–4935, 2008.

[28] C. G. Lopes and A. H. Sayed, "Diffusion least-mean squares over adaptive networks: Formulation and performance analysis," *IEEE Trans. Signal Process.*, vol. 56, no. 7, pp. 3122–3136, 2008.

[29] A. H. Sayed, "Adaptation, learning, and optimization over networks," *Foundations and Trends in Machine Learning*, vol. 7, pp. 311–801, 2014.

[30] L. Li and J. A. Chambers, "Distributed adaptive estimation based on the APA algorithm over diffusion networks with changing topology," *IEEE Workshop on Stat. Signal Process.*, pp. 757–760, 2009.

[31] K. R. Varshney, "Generalization error of linear discriminant analysis in spatially-correlated sensor networks," *IEEE Trans. Signal Process.*, vol. 60, no. 6, pp. 3295–3301, 2012.

[32] M. Hellkvist, A. Özçelikkale, and A. Ahlén, "Generalization error for linear regression under distributed learning," *IEEE Int. Workshop on Signal Process. Advances in Wireless Commun.*, May 2020.

[33] J. Tsitsiklis, D. Bertsekas, and M. Athans, "Distributed asynchronous deterministic and stochastic gradient optimization algorithms," *IEEE Trans. Automat. Control*, vol. 31, no. 9, pp. 803–812, 1986.

[34] T. Ben-Nun and T. Hoefler, "Demystifying parallel and distributed deep learning: An in-depth concurrency analysis," *ACM Computing Surveys*, vol. 52, no. 4, pp. 1–43, 2019.

[35] M. Jaggi, V. Smith, M. Takác, J. Terhorst *et al.*, "Communication-efficient distributed dual coordinate ascent," *Adv. Neural Inf. Process. Systems*, pp. 3068–3076, 2014.

[36] C. Ma, J. Konečnỳ, M. Jaggi, V. Smith *et al.*, "Distributed optimization with arbitrary local solvers," *Optimization Methods and Softw.*, vol. 32, no. 4, pp. 813–848, 2017.

[37] Y. LeCun, C. Cortes, and C. J. C. Burges, "MNIST handwritten digit database," 2010. [Online]. Available: http://yann.lecun.com/exdb/mnist/

[38] C. Zhang, S. Bengio, M. Hardt, B. Recht *et al.*, "Understanding deep learning requires rethinking generalization," *arXiv:1611.03530*, Nov. 2016.

[39] M. Rudelson and R. Vershynin, "Non-asymptotic theory of random matrices: extreme singular values," *Proc. of the Int. Congress of Mathematicians*, pp. 1576–1602, 2010.

[40] M. Hellkvist, A. Özçelikkale, and A. Ahlén, "Generalization error for linear regression under distributed learning," *arXiv:2004.14637*, Apr. 2020.

[41] R. Vershynin, *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018.

[42] R. Horn and C. Johnson, *Matrix Analysis*. Cambridge University Press, 1990.

[43] S. Foucart and H. Rauhut, *A Mathematical Introduction to Compressive Sensing*. Springer, 2013.

[44] A. Rahimi and B. Recht, "Random features for large-scale kernel machines," *Adv. Neural Inf. Process. Syst.*, pp. 1177–1184, 2008.

[45] M. Garcia, "MNIST data," Oct. 2018. [Online]. Available: https://github.com/datapythonista/mnist .

[46] C. Bishop, *Pattern Recognition and Machine Learning*, ser. Information Science and Statistics. Springer, 2006.

[47] M. Raginsky, A. Rakhlin, and M. Telgarsky, "Non-convex learning via stochastic gradient langevin dynamics: a nonasymptotic analysis," *Conference on Learning Theory, Proceedings of Machine Learning Research*, vol. 65, pp. 1674–1703, Jul 2017.

[48] S. Gelfand and S. Mitter, "Recursive stochastic algorithms for global optimization in $R^d$," *SIAM Journal on Control and Optimization*, vol. 29, pp. 999–1018, 1991.

[49] N. Halko, P.-G. Martinsson, and J. A. Tropp, "Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions," *SIAM review*, vol. 53, no. 2, pp. 217–288, 2011.

[50] R. D. Cook and L. Forzani, "On the mean and variance of the generalized inverse of a singular Wishart matrix," *Electron. J. Statist.*, vol. 5, pp. 146–158, 2011.