Exploiting Shared Representations for Personalized Federated Learning

Liam Collins, Hamed Hassani, Aryan Mokhtari, Sanjay Shakkottai

Abstract

Deep neural networks have shown the ability to extract universal feature representations from data such as images and text that have been useful for a variety of learning tasks. However, the fruits of representation learning have yet to be fully-realized in federated settings. Although data in federated settings is often non-i.i.d. across clients, the success of centralized deep learning suggests that data often shares a global feature representation, while the statistical heterogeneity across clients or tasks is concentrated in the *labels*. Based on this intuition, we propose a novel federated learning framework and algorithm for learning a shared data representation across clients and unique local heads for each client. Our algorithm harnesses the distributed computational power across clients to perform many local-updates with respect to the low-dimensional local parameters for every update of the representation. We prove that this method obtains linear convergence to the ground-truth representation with near-optimal sample complexity in a linear setting, demonstrating that it can efficiently reduce the problem dimension for each client. This result is of interest beyond federated learning to a broad class of problems in which we aim to learn a shared low-dimensional representation among data distributions, for example in meta-learning and multi-task learning. Further, extensive experimental results show the empirical improvement of our method over alternative personalized federated learning approaches in federated environments with heterogeneous data.

^{*}Department of Electrical and Computer Engineering, The University of Texas at Austin, Austin, TX, USA. {liamc@utexas.edu, mokhtari@austin.utexas.edu, sanjay.shakkottai@utexas.edu}.

[†]Department of Electrical and Systems Engineering, University of Pennsylvania, Philadelphia, PA, USA. {hassani@seas.upenn.edu}.

1 Introduction

Many of the most heralded successes of modern machine learning have come in *centralized* settings, wherein a single model is trained on a large amount of centrally-stored data. The growing number of data-gathering devices, however, calls for a distributed architecture to train models. Federated learning aims at addressing this issue by providing a platform in which a group of clients collaborate to learn effective models for each client by leveraging the local computational power, memory, and data of all clients [McMahan et al., 2017]. The task of coordinating between the clients is fulfilled by a central server that combines the models received from the clients at each round and broadcasts the updated information to them. Importantly, the server and clients are restricted to methods that satisfy communication and privacy constraints, preventing them from directly applying centralized techniques.

However, one of the most important challenges in federated learning is the issue of data heterogeneity, where the underlying data distribution of client tasks could be substantially different from each other. In such settings, if the server and clients learn a single shared model (e.g., by minimizing average loss), the resulting model could perform poorly for many of the clients in the network (and also not generalize well across diverse data [Jiang et al., 2019]). In fact, for some clients, it might be better to simply use their own local data (even if it is small) to train a local model; see Figure 1. Finally, the (federated) trained model may not generalize well to unseen clients that have not participated in the training process. These issues raise the question:

"How can we exploit the data and computational power of all clients in data heterogeneous settings to learn a personalized model for each client?"

We address this question by taking advantage of the common representation among clients. Specifically, we view the data heterogeneous federated learning problem as n parallel learning tasks that they possibly have some common structure, and our goal is to learn and exploit this common representation to improve the quality of each client's model. This approach draws inspiration from centralized learning, where we have witnessed success in training multiple tasks or learning multiple classes simultaneously by leveraging a common (low-dimensional) representation (e.g. in image classification, next-word prediction) [Bengio et al., 2013, LeCun et al., 2015].

Main Contributions. We introduce a novel federated learning framework and an associated algorithm for data heterogeneous settings. We summarize our main contributions below.

- (i) FedRep Algorithm. Federated Representation Learning (FedRep) leverages all of the data stored across clients to learn a global low-dimensional representation using gradient-based updates. Further, it enables each client to compute a personalized, low-dimensional classifier, which we term as the client's head, that accounts for the unique labeling of each client's local data.
- (ii) Optimization for linear representation learning. We show that FedRep converges to the ground-truth representation at an exponentially fast rate in the case that each client aims to solve a linear regression problem with a two-layer linear neural network. In this special case, we reduce FedRep to alternating minimization (for the heads)-descent (for the representation). Our analysis shows that this simple algorithm requires only $\mathcal{O}((d/n + \log(n)) \log(1/\epsilon))$ samples

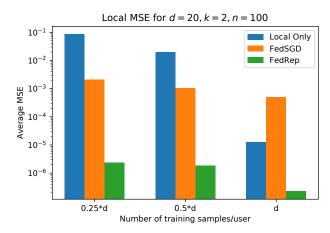


Figure 1: Local only training suffers in small-training data regimes, whereas training a single global model with FedSGD cannot overcome client heterogeneity even when the number of training samples is large. FedRep exploits a common representation of the clients to achieve small error in all cases.

per client to reach an ϵ -accurate representation, where n is the number of clients and d is the dimension of the data. This result is of interest beyond federated learning since it shows that alternating minimization-descent efficiently solves the linear multi-task representation learning problem considered in Du et al. [2020], Maurer et al. [2016], Tripuraneni et al. [2020a].

(iii) **Empirical Results.** Through a combination of synthetic and real datasets (CIFAR10, CIFAR100, FEMNIST, Sent140) we show the benefits of FedRep in: (a) leveraging many local updates, (b) robustness to different levels of heterogeneity, and (c) generalization to new clients. Our experiments indicate that FedRep outperforms several important baselines in heterogeneous settings that share a global representation.

Benefits of FedRep. FedRep has numerous advantages over standard federated learning (in which a single model is learned):

- (I) Provable gains of cooperation. From our sample complexity bounds, it follows that with FedRep, the sample complexity per client scales as $\Theta(d/n + \log(n))$. On the other hand, local learning (without any collaboration) has a sample complexity that scales as $\Theta(d)$. Thus, if $1 \ll n \ll e^{\Theta(d)}$ (see Section 4.2 for details), we expect benefits of collaboration through federation. When d is large (as is typical in practice), $e^{\Theta(d)}$ is exponentially larger, and federation helps each client. To the best of our knowledge, this is the first sample-complexity-based result for personalized federated learning that demonstrates the benefit of cooperation.
- (II) Generalization to new clients. For a new client, since a ready-made representation is available, the client only needs to learn a head with a low-dimensional representation of dimension k. Thus, its sample complexity scales only as $\Theta(k)$ instead of $\Theta(d)$ if no representation is learned.
- (III) More local updates. By reducing the problem dimension, each client can make many local updates at each communication round, which is beneficial in learning its own individual head. This is unlike standard federated learning where multiple local updates in a heterogeneous setting moves each client away from the best averaged representation, and thus hurts performance.

1.1 Related Work.

Personalized Federated Learning. A variety of recent works have studied personalization in federated learning using, for example, local fine-tuning [Wang et al., 2019, Yu et al., 2020], metalearning [Chen et al., 2018, Fallah et al., 2020, Jiang et al., 2019, Khodak et al., 2019], additive mixtures of local and global models [Deng et al., 2020, Hanzely and Richtárik, 2020, Mansour et al., 2020, and multi-task learning [Smith et al., 2017]. In all of these methods, each client's subproblem is still full-dimensional - there is no notion of learning a dimensionality-reduced set of local parameters. More recently, Liang et al. [2020] also proposed a representation learning method for federated learning, but their method attempts to learn many local representations and a single global head as opposed to a single global representation and many local heads. Earlier, Arivazhagan et al. [2019] presented an algorithm to learn local heads and a global network body, but their local procedure jointly updates the head and body (using the same number of updates), and they did not provide any theoretical justification for their proposed method. Meanwhile, another line of work has studied federated learning in heterogeneous settings [Haddadpour et al., 2020, Karimireddy et al., 2020, Mitra et al., 2021, Pathak and Wainwright, 2020, Reddi et al., 2020, Reisizadeh et al., 2020, Wang et al., 2020, and the optimization-based insights from these works may be used to supplement our formulation and algorithm.

Linear representation learning. The idea to learn a shared representation of tasks is a classical approach in multi-task learning [Ando et al., 2005, Balcan et al., 2015, Baxter, 2000, Bengio et al., 2013, Bullins et al., 2019, Denevi et al., 2018, Kong et al., 2020, LeCun et al., 2015, Pontil and Maurer, 2013, Rish et al., 2008, Tripuraneni et al., 2020b]. In particular, we aim to learn a low-dimensional subspace in which the ground-truth regressors for a collection of linear regression tasks lie. This problem is most similar to the linear representation learning problem considered by Du et al. [2020], Maurer et al. [2016], Tripuraneni et al. [2020a]. All three of these works show statistical rates of convergence of solutions to the ERM objective to the ground-truth representation, with Tripuraneni et al. [2020a] and Du et al. [2020] improving the $\mathcal{O}(d/n)$ rate from Maurer et al. [2016] (in the realizable case) to $\mathcal{O}(d/mn)$. Du et al. [2020] also provide similar complexity-based results for learning nonlinear representations with access to an ERM oracle, but their results in the linear case require $m = \Omega(d)$ samples per task, mitigating the benefit of cooperation. Tripuraneni et al. [2020a] further present and analyze a Method-of-Moments-based algorithm to solve the ERM problem, which achieves sample complexity per task with efficient dimension-dependence $(\Theta(d/n))$ but requires $m = \Omega(1/n\epsilon^2)$ samples per task to find an ϵ -accurate representation. In contrast, we show that alternating minimization-descent requires only $m = \Omega((d/n + \log(n)) \log(1/\epsilon))$ samples per client to obtain a representation with ϵ -accuracy.

2 Problem Formulation

The generic form of federated learning with n clients is

$$\min_{(q_1,\dots,q_n)\in\mathcal{Q}_n} \frac{1}{n} \sum_{i=1}^n f_i(q_i),\tag{1}$$

where f_i and q_i are the error function and learning model for the *i*-th client, respectively, and Q_n is the space of feasible sets of n models. We consider a supervised setting in which the data for the

i-th client is generated by a distribution $(\mathbf{x}_i, y_i) \sim \mathcal{D}_i$. The learning model $q_i : \mathbb{R}^d \to \mathcal{Y}$ maps inputs $\mathbf{x}_i \in \mathbb{R}^d$ to predicted labels $q_i(\mathbf{x}_i) \in \mathcal{Y}$, which we would like to resemble the true labels y_i . The error f_i is in the form of an expected risk over \mathcal{D}_i , namely $f_i(q_i) := \mathbb{E}_{(\mathbf{x}_i, y_i) \sim \mathcal{D}_i}[\ell(q_i(\mathbf{x}_i), y_i)]$, where $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ is a loss function that penalizes the distance of $q_i(\mathbf{x}_i)$ from y_i .

In order to minimize f_i , the *i*-th client accesses a dataset of M_i labeled samples $\{(\mathbf{x}_i^j, y_i^j)\}_{j=1}^{M_i}$ from \mathcal{D}_i for training. Federated learning addresses settings in which the M_i 's are typically small relative to the problem dimension while the number of clients n is large. Thus, clients may not be able to obtain solutions q_i with small expected risk by training completely locally on *only* their M_i local samples. Instead, federated learning enables the clients to cooperate, by exchanging messages with a central server, in order to learn models using the cumulative data of all the clients.

Standard approaches to federated learning aim at learning a single shared model $q = q_1 = \cdots = q_n$ that performs well on average across the clients [Li et al., 2018, McMahan et al., 2017]. In this way, the clients aim to solve a special version of Problem (1), which is to minimize $(1/n) \sum_i f_i(q)$ over the choice of the shared model q. However, this approach may yield a solution that performs poorly in heterogeneous settings where the data distributions \mathcal{D}_i vary across the clients. Indeed, in the presence of data heterogeneity, the error functions f_i will have different forms and their minimizers are not the same. Hence, learning a shared model q may not provide good solution to Problem (1) This necessities the search for more personalized solutions $\{q_i\}$ that can be learned in a federated manner using the clients' data.

Learning a Common Representation. We are motivated by insights from centralized machine learning that suggest that heterogeneous data distributed across tasks may share a common representation despite having different labels [Bengio et al., 2013, LeCun et al., 2015]; e.g., shared features across many types of images, or across word-prediction tasks. Using this common (low-dimensional) representation, the labels for each client can be simply learned using a linear classifier or a shallow neural network.

Formally, we consider a setting consisting of a global representation $\phi : \mathbb{R}^d \to \mathbb{R}^k$, which maps data points to a lower space of size k, and client-specific heads $h_i : \mathbb{R}^k \to \mathcal{Y}$. The model for the i-th client is the composition of the client's local parameters and the representation: $q_i(\mathbf{x}) = (h_i \circ \phi)(\mathbf{x})$. Critically, $k \ll d$, meaning that the number of parameters that must be learned locally by each client is small. Thus, we can assume that any client's optimal classifier for any fixed representation is easy to compute, which motivates the following re-written global objective:

$$\min_{\phi \in \Phi} \frac{1}{n} \sum_{i=1}^{n} \min_{h_i \in \mathcal{H}} f_i(h_i \circ \phi), \tag{2}$$

where Φ is the class of feasible representations and \mathcal{H} is the class of feasible heads. In our proposed scheme, clients cooperate to learn the global model using all clients' data, while they use their local information to learn their personalized head. We discuss this in detail in Section 3.

Formally, we consider a setting consisting of a global representation $q_{\phi}: \mathbb{R}^{d} \to \mathbb{R}^{k}$, which is a function parameterized by $\phi \in \Phi$ that maps data points to a lower space of dimension k, and client-specific heads $q_{h_{i}}: \mathbb{R}^{k} \to \mathcal{Y}$, which are functions parameterized by $h_{i} \in \mathcal{H}$ for $i \in [n]$ that map from the low-dimensional representation space to the label space. The model for the i-th client is the composition of the client's local parameters and the representation: $q_{i}(\mathbf{x}) = (q_{h_{i}} \circ q_{\phi})(\mathbf{x})$. Critically, $k \ll d$, meaning that the number of parameters that must be learned locally by each client may be

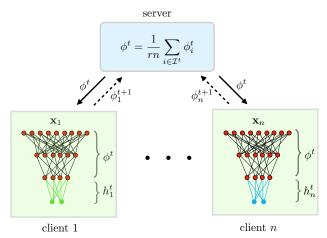


Figure 2: Federated representation learning structure where clients and the server aim at learning a global representation ϕ together, while each client i learns its unique head h_i locally.

small. Thus, we can assume that any client's optimal classifier for any fixed representation is easy to compute, which motivates the following re-written global objective:

$$\min_{\phi \in \Phi} \frac{1}{n} \sum_{i=1}^{n} \min_{h_i \in \mathcal{H}} f_i(h_i, \phi), \tag{3}$$

where we have used the shorthand $f_i(h_i, \phi) := f_i(q_{h_i} \circ q_{\phi})$ for ease of notation. In our proposed scheme, clients cooperate to learn the global model using all clients' data, while they use their local information to learn their personalized head. We discuss this in detail in Section 3.

2.1 Comparison with Standard Federated Learning

To formally demonstrate the advantage of our formulation over the standard (single-model) federated learning formulation in heterogeneous settings with a shared representation, we study a linear representation setting with quadratic loss. As we will see below, standard federated learning cannot recover the underlying representation in the face of heterogeneity, while our formulation does indeed recover it.

Consider a setting in which the functions f_i are quadratic losses, the representation ϕ is a projection onto a k-dimensional subspace of \mathbb{R}^d given by matrix $\mathbf{B} \in \mathbb{R}^{d \times k}$, and the i-th client's local head h_i is a vector $\mathbf{w}_i \in \mathbb{R}^k$. In this setting, we model the local data of clients $\{\mathcal{D}_i\}_i$ such that $y_i = \mathbf{w}_i^{*\top} \mathbf{B}^{*\top} \mathbf{x}_i$ for some ground-truth representation $\mathbf{B}^* \in \mathbb{R}^{d \times k}$ and local heads $\mathbf{w}_i^* \in \mathbb{R}^k$. This setting will be described in detail in Section 4. In particular, one can show that the expected error over the data distribution \mathcal{D}_i has the following form: $f_i(\mathbf{w}_i \circ \mathbf{B}) := \frac{1}{2} \|\mathbf{B} \mathbf{w}_i - \mathbf{B}^* \mathbf{w}_i^*\|_2^2$. Consequently, Problem (3) becomes

$$\min_{\mathbf{B} \in \mathbb{R}^{d \times k}, \mathbf{w}_{i}, \dots, \mathbf{w}_{n} \in \mathbb{R}^{k}} \frac{1}{2n} \sum_{i=1}^{n} \|\mathbf{B} \mathbf{w}_{i} - \mathbf{B}^{*} \mathbf{w}_{i}^{*}\|_{2}^{2}.$$

$$(4)$$

In contrast, standard federated learning methods, which aim to learn a shared model (\mathbf{B}, \mathbf{w}) for all

the clients, solve

$$\min_{\mathbf{B} \in \mathbb{R}^{d \times k}, \mathbf{w} \in \mathbb{R}^k} \frac{1}{2n} \sum_{i=1}^n \|\mathbf{B}\mathbf{w} - \mathbf{B}^* \mathbf{w}_i^*\|_2^2.$$
 (5)

Let $(\hat{\mathbf{B}}, \{\hat{\mathbf{w}}_i\}_i)$ denote a global minimizer of (4). We thus have $\hat{\mathbf{B}}\hat{\mathbf{w}}_i = \mathbf{B}^*\mathbf{w}_i^*$ for all $i \in [n]$. Also, it is not hard to see that $(\mathbf{B}^{\diamond}, \mathbf{w}^{\diamond})$ is a global minimizer of (5) if and only if $\mathbf{B}^{\diamond}\mathbf{w}^{\diamond} = \mathbf{B}^*(\frac{1}{n}\sum_{i=1}^n \mathbf{w}_i^*)$. Thus, our formulation finds an exact solution with zero global error, whereas standard federated learning has global error of $\frac{1}{2n}\sum_{i=1}^n \|\frac{1}{n}\mathbf{B}^*\sum_{i'=1}^n (\mathbf{w}_{i'}^* - \mathbf{w}_i^*)\|_2^2$, which grows with the heterogeneity of the \mathbf{w}_i^* . Moreover, since solving our formulation provides n matrix equations, we can fully recover the column space of \mathbf{B}^* as long as \mathbf{w}_i^* 's span \mathbb{R}^k . In contrast, solving (5) yields only one matrix equation, so there is no hope to recover the column space of \mathbf{B}^* for any k > 1.

3 FedRep Algorithm

FedRep solves Problem (3) by distributing the computation across clients. The server and clients aim to learn the parameters of the global representation together, while the *i*-th client aims to learn its unique local head locally (see Figure 2). To do so, FedRep alternates between client updates and a server update on each communication round.

Client Update. On each round, a constant fraction $r \in (0, 1]$ of the clients are selected to execute a client update. In the client update, client i makes τ_h local gradient-based updates to solve for its optimal head given the current global representation ϕ^t communicated by the server. Namely, for $s = 1, \ldots, \tau_h$, client i updates its head as follows:

$$h_i^{t,s} = \text{GRD}(f_i(h_i^{t,s-1}, \phi^t), h_i^{t,s-1}, \alpha),$$

where $GRD(f, h, \alpha)$ is generic notation for an update of the variable h using a gradient of function f with respect to h and the step size α . For example, $GRD(f_i(h_i^{t,s-1}, \phi^t), h_i^{t,s-1}, \alpha)$ can be a step of gradient descent, stochastic gradient descent (SGD), SGD with momentum, etc. Typically, we will choose τ_h to be large, since more local epochs for the head means that we come closer to solving the inner minimization in (3), which means that the updates for the representation are more accurate.

Next, the client executes τ_{ϕ} local updates for its representation, starting from the global representation ϕ^{t-1} :

$$\phi_i^{t,s} = \text{GRD}(f_i(h_i^{t,\tau_h}, \phi_i^{t,s-1}), \phi_i^{t,s-1}, \alpha),$$

for $s=1,\ldots,\tau_{\phi}$.

Server Update. Once the local updates with respect to the head and representation finish, the client participates in the server update by sending its locally-updated representation $\phi_i^{t,\tau_{\phi}}$ to the server. The server then averages the local updates to compute the next representation ϕ^t . The entire procedure is outlined in Algorithm 1.

Algorithm 1 FedRep

```
Parameters: Participation rate r, step size \alpha; number of local updates for the head \tau_h and for
the representation \tau_{\phi}; number of communication rounds T.
Initialize \phi^0, h_1^0, \dots, h_n^0
for t = 1, 2, ..., T do
    Server receives a batch of clients \mathcal{I}^t of size rn
    Server sends current representation \phi^t to these clients
    for each client i in \mathcal{I}^t do
        Client i initializes h_i^{t,0} \leftarrow h_i^{t-1,\tau_h}
         Client i makes \tau_h updates to its head:
        \begin{aligned} & \mathbf{for} \ s = 1 \ \mathbf{to} \ \tau_h \ \mathbf{do} \\ & h_i^{t,s} \leftarrow \mathtt{GRD}(f_i(h_i^{t,s-1},\phi^{t-1}),h_i^{t,s-1},\alpha) \end{aligned}
        Client i initializes \phi_i^{t,0} \leftarrow \phi^{t-1}
         Client i makes \tau_{\phi} updates to its representation:
        \begin{aligned} & \mathbf{for} \ s = 1 \ \mathbf{to} \ \tau_{\phi} \ \mathbf{do} \\ & \phi_i^{t,s} \leftarrow \mathtt{GRD}(f_i(h_i^{t,\tau_h}, \phi_i^{t,s-1}), \phi_i^{t,s-1}, \alpha) \end{aligned}
        Client i sends updated representation \phi_i^{t,\tau_{\phi}} to server
     \begin{array}{l} \textbf{for each client } i \text{ not in } \mathcal{I}^t, \, \textbf{do} \\ \text{Set } h_i^{t,\tau_h} \leftarrow h_i^{t-1,\tau_h} \end{array} 
    end for
    Server computes the new representation as \phi^t = \frac{1}{rn} \sum_{i \in \mathcal{I}^t} \phi_i^{t,\tau_\phi}
end for
```

4 Low-Dimensional Linear Representation

In this section, we analyze an instance of Problem (3) with quadratic loss functions and linear models, as discussed in Section 2.1. Here, each client's problem is to solve a linear regression with a two-layer linear neural network. In particular, each client i attempts to find a shared global projection onto a low-dimension subspace $\mathbf{B} \in \mathbb{R}^{d \times k}$ and a unique regressor $\mathbf{w}_i \in \mathbb{R}^k$ that together accurately map its samples $\mathbf{x}_i \in \mathbb{R}^d$ to labels $y_i \in \mathbb{R}$. The matrix \mathbf{B} corresponds to the representation ϕ , and \mathbf{w}_i corresponds to local head h_i for the i-th client. We thus have $(q_{h_i} \circ q_{\phi})(\mathbf{x}_i) = \mathbf{w}_i^{\top} \mathbf{B}^{\top} \mathbf{x}_i$. Hence, the loss function for client i is given by:

$$f_i(\mathbf{w}_i, \mathbf{B}) := \frac{1}{2} \mathbb{E}_{(\mathbf{x}_i, y_i) \sim \mathcal{D}_i} \left[(y_i - \mathbf{w}_i^\top \mathbf{B}^\top \mathbf{x}_i)^2 \right]$$
 (6)

meaning that the global objective is:

$$\min_{\mathbf{B} \in \mathbb{R}^{d \times k} \atop \mathbf{W} \in \mathbb{R}^{n \times k}} F(\mathbf{B}, \mathbf{W}) \coloneqq \frac{1}{2n} \sum_{i=1}^{n} \mathbb{E}_{(\mathbf{x}_{i}, y_{i})} \left[(y_{i} - \mathbf{w}_{i}^{\top} \mathbf{B}^{\top} \mathbf{x}_{i})^{2} \right], \tag{7}$$

where $\mathbf{W} = [\mathbf{w}_1^\top, \dots, \mathbf{w}_n^\top] \in \mathbb{R}^{n \times k}$ is the concatenation of client-specific heads. To evaluate the ability of FedRep to learn an accurate representation, we model the local datasets $\{\mathcal{D}_i\}_i$ such that,

for
$$i = 1 \dots, n$$

$$y_i = \mathbf{w}_i^{*\top} \mathbf{B}^{*\top} \mathbf{x}_i,$$

for some ground-truth representation $\mathbf{B}^* \in \mathbb{R}^{d \times k}$ and local heads $\mathbf{w}_i^* \in \mathbb{R}^k$ —i.e. a standard regression setting. In other words, all of the clients' optimal solutions live in the same k-dimensional subspace of \mathbb{R}^d , where k is assumed to be small. Moreover, we make the following standard assumption on the samples \mathbf{x}_i .

Assumption 1 (Sub-gaussian design). The samples $\mathbf{x}_i \in \mathbb{R}^d$ are i.i.d. with mean $\mathbf{0}$, covariance \mathbf{I}_d , and are \mathbf{I}_d -sub-gaussian, i.e. $\mathbb{E}[e^{\mathbf{v}^{\top}\mathbf{x}_i}] \leq e^{\|\mathbf{v}\|_2^2/2}$ for all $\mathbf{v} \in \mathbb{R}^d$.

4.1 FedRep

We next discuss how FedRep tries to recover the optimal representation in this setting. First, the server and clients execute the Method of Moments to learn an initial representation. Then, client and server updates are executed in an alternating fashion as follows.

Client Update. As in Algorithm 1, rn clients are selected on round t to update their current local head \mathbf{w}_i^t and the global representation \mathbf{B}^t . Each selected client i samples a fresh batch $\{\mathbf{x}_i^{t,j}, y_i^{t,j}\}_{j=1}^m$ of m samples according to its local data distribution \mathcal{D}_i to use for updating both its head and representation on each round t that it is selected. That is, within the round, client i considers the batch loss

$$\hat{f}_i^t(\mathbf{w}_i^t, \mathbf{B}^t) := \frac{1}{2m} \sum_{j=1}^m (y_i^{t,j} - \mathbf{w}_i^{t^\top} \mathbf{B}^{t^\top} \mathbf{x}_i^{t,j})^2.$$
 (8)

Since \hat{f}_i^t is strongly convex with respect to \mathbf{w}_i^t , the client can find an update for a local head that is ϵ -close to the global minimizer of (8) after at most $\log(1/\epsilon)$ local gradient updates. Alternatively, since the function is also quadratic, the client can solve for the optimal \mathbf{w} directly in only $\mathcal{O}(mk^2+k^3)$ operations. Thus, since FedRep calls for many local updates for the head, to simplify the analysis we assume each selected client obtains $\mathbf{w}_i^{t+1} = \operatorname{argmin}_{\mathbf{w}} \hat{f}_i^t(\mathbf{w}, \mathbf{B}^t)$ during each round of local updates.

Server Update. After updating its head, client i updates the global representation with one step of gradient descent using the same m samples and sends the update to the server, as outlined in Algorithm 2. Note that in practice, each client may execute multiple gradient-based updates before sending its updated representation back to the server, but here we consider the case that they make one step of gradient descent for simplicity. Once the server receives the representations, it averages them and orthogonalizes the resulting matrix to compute the new representation.

4.2 Analysis

As mentioned earlier, in FedRep, each client i perform an alternating minimization-descent method to solve its nonconvex objective in (8). This means the global loss over all clients at round t is given by

$$\frac{1}{n} \sum_{i=1}^{n} \hat{f}_{i}^{t}(\mathbf{w}_{i}^{t}, \mathbf{B}^{t}) \coloneqq \frac{1}{2mn} \sum_{i=1}^{n} \sum_{j=1}^{m} (y_{i}^{t,j} - \mathbf{w}_{i}^{t^{\top}} \mathbf{B}^{t^{\top}} \mathbf{x}_{i}^{t,j})^{2}.$$
(9)

Algorithm 2 FedRep for linear regression

end for

Input: Step size η ; number of rounds T, participation rate r. Initialization: Each client $i \in [n]$ sends $\mathbf{Z}_i := \frac{1}{m} \sum_{j=1}^m (y_i^{0,j})^2 \mathbf{x}_i^{0,j} (\mathbf{x}_i^{0,j})^{\top}$ to server, server computes

Server initializes
$$\mathbf{B}^0 \leftarrow \mathbf{U}$$

for $t = 1, 2, ..., T$ do
Server receives a subset \mathcal{I}^t of clients of size rn
Server sends current representation \mathbf{B}^t to these clients
for $i \in \mathcal{I}^t$ do
Client update:
Client i samples a fresh batch of m samples
Client i updates \mathbf{w}_i :
 $\mathbf{w}_i^{t+1} \leftarrow \operatorname{argmin}_{\mathbf{w}} \hat{f}_i^t(\mathbf{w}, \mathbf{B}^t)$
Client i updates representation:
 $\mathbf{B}_i^{t+1} \leftarrow \mathbf{B}^t - \eta \nabla_{\mathbf{B}} \hat{f}_i^t(\mathbf{w}_i^{t+1}, \mathbf{B}^t)$
Client i sends \mathbf{B}_i^{t+1} to the server
end for
Server update: $\bar{\mathbf{B}}^{t+1} \leftarrow \frac{1}{rn} \sum_{i \in \mathcal{I}^t} \mathbf{B}_i^{t+1}$; $\mathbf{B}^{t+1}, \mathbf{R}^{t+1} \leftarrow \operatorname{QR}(\bar{\mathbf{B}}^{t+1})$

This objective has many global minima, including all pairs of matrices $(\mathbf{Q}^{-1}\mathbf{W}^*, \mathbf{B}^*\mathbf{Q}^{\top})$ where $\mathbf{Q} \in \mathbb{R}^{k \times k}$ is invertible, eliminating the possibility of exactly recovering the ground-truth factors $(\mathbf{W}^*, \mathbf{B}^*)$. Instead, the ultimate goal of the server is to recover the ground-truth *representation*, i.e., the column space of \mathbf{B}^* . To evaluate how closely the column space is recovered, we define the distance between subspaces as follows.

Definition 1. The principal angle distance between the column spaces of $\mathbf{B}_1, \mathbf{B}_2 \in \mathbb{R}^{d \times k}$ is given by

$$\operatorname{dist}(\mathbf{B}_1, \mathbf{B}_2) \coloneqq \|\hat{\mathbf{B}}_{1,\perp}^{\top} \hat{\mathbf{B}}_2\|_2,\tag{10}$$

where $\hat{\mathbf{B}}_{1,\perp}$ and $\hat{\mathbf{B}}_2$ are orthonormal matrices satisfying $\operatorname{span}(\hat{\mathbf{B}}_{1,\perp}) = \operatorname{span}(\mathbf{B}_1)^{\perp}$ and $\operatorname{span}(\hat{\mathbf{B}}_2) = \operatorname{span}(\mathbf{B}_2)$.

The principal angle distance is a typical metric for measuring the distance between subspaces (e.g. Jain et al. [2013]). Next, we make two standard assumptions.

Assumption 2 (Client diversity). Let $\bar{\sigma}_{\min,*} := \min_{\mathcal{I} \in [n], |\mathcal{I}| = rn} \sigma_{\min}(\frac{1}{\sqrt{rn}} \mathbf{W}_{\mathcal{I}}^*)$, i.e. $\bar{\sigma}_{\min,*}$ is the minimum singular value of any matrix that can be obtained by taking rn rows of $\frac{1}{\sqrt{rn}} \mathbf{W}^*$. Then $\bar{\sigma}_{\min,*} > 0$.

Assumption 2 states that if we select any rn clients, their optimal heads span \mathbb{R}^k . Indeed, this assumption is weak as we expect the number of participating clients rn to be substantially larger than k. Note that if we do not have client solutions that span \mathbb{R}^k , recovering \mathbf{B}^* would be impossible because the samples (\mathbf{x}_i^j, y_i^j) may never contain any information about one or more features of \mathbf{B}^* .

Assumption 3 (Client normalization). The ground-truth client-specific parameters satisfy $\|\mathbf{w}_i^*\|_2 = \sqrt{k}$ for all $i \in [n]$, and \mathbf{B}^* has orthonormal columns.

Assumption 2 ensures that the ground-truth matrix $\mathbf{W}^*\mathbf{B}^{*\top}$ is row-wise *incoherent*, i.e. its row norms have similar magnitudes. We define this formally in Appendix B. Incoherence of the ground-truth matrices is a key property required for efficient matrix completion and other sensing problems with sparse measurements [Chi et al., 2019]. Since our measurement matrices are row-wise sparse, we require the row-wise incoherence of the ground truth. Note that Assumption 3 can be relaxed to allow $\|\mathbf{w}_i^*\|_2 \leq O(\sqrt{k})$, as the exact normalization is only for simplicity of analysis.

Our main result shows that the iterates $\{\mathbf{B}^t\}_t$ generated by FedRep in this setting linearly converge to the optimal representation \mathbf{B}^* in principal angle distance.

Theorem 1. Define $E_0 := 1 - \operatorname{dist}^2(\mathbf{B}^0, \mathbf{B}^*)$ and $\bar{\sigma}_{\max,*} := \max_{\mathcal{I} \in [n], |\mathcal{I}| = rn} \sigma_{\max}(\frac{1}{\sqrt{rn}} \mathbf{W}_{\mathcal{I}}^*)$ and $\bar{\sigma}_{\min,*} := \min_{\mathcal{I} \in [n], |\mathcal{I}| = rn} \sigma_{\min}(\frac{1}{\sqrt{rn}} \mathbf{W}_{\mathcal{I}}^*)$, i.e. the maximum and minimum singular values of any matrix that can be obtained by taking rn rows of $\frac{1}{\sqrt{rn}} \mathbf{W}^*$. Let $\kappa := \bar{\sigma}_{\max,*}/\bar{\sigma}_{\min,*}$. Suppose that $m \geq c(\kappa^4 k^2 d/(E_0^2 rn) + \kappa^4 k^3 \log(rn)/E_0^2)$ for some absolute constant c. Then for any t and any $\eta \leq 1/(4\bar{\sigma}_{\max,*}^2)$, we have

$$\operatorname{dist}(\mathbf{B}^T, \mathbf{B}^*) \le \left(1 - \eta E_0 \bar{\sigma}_{\min,*}^2 / 2\right)^{T/2} \operatorname{dist}(\mathbf{B}^0, \mathbf{B}^*), \tag{11}$$

with probability at least $1 - Te^{-100\min(k^2\log(rn),d)}$.

From Assumption 2, we have that $\bar{\sigma}_{\min,*}^2 > 0$, so the RHS of (11) strictly decreases with T for appropriate step size. Considering the complexity of m and the fact that the algorithm converges exponentially fast, the total number of samples required per client to reach an ϵ -accurate solution in principal angle distance is $\Theta(m \log(1/\epsilon))$, which is

$$\Theta\left(\left[\kappa^4 k^2 \left(\frac{d}{rn} + k \log(rn)\right)\right] \log\left(\frac{1}{\epsilon}\right)\right). \tag{12}$$

Next, a few remarks about this sample complexity follow.

When and whom does federation help? Observe that for a single client with no collaboration, the sample complexity scales as $\Theta(d/n + \log(n))$, treating k, κ and r as constants. Thus, so long as $d/n + \log(n) \ll d$, federation helps. This holds in several settings, for instance when $1 \ll n \ll e^{\Theta(d)}$. In practical scenarios, d (the data dimension) is large, and thus $e^{\Theta(d)}$ is exponentially larger; thus collaboration helps each individual client. Furthermore, new clients who enter the system later have a representation available for free, so these new clients' sample complexity is only $\Theta(k)$ because they each only need to solve a k-dimensional linear regression problem [Hsu et al., 2012]. Thus, both the overall system benefits (a representation has been learned, which is useful for the new client because it now only needs to learn a head), and each individual client that took part in the federated training also benefits.

Connection to matrix sensing. The problem in (7) is an instance of matrix sensing; see the proof in Appendix B for more details. Considering this connection, our theoretical results also contribute to the theoretical study of matrix sensing. Although matrix sensing is a well-studied problem, our setting presents two new analytical challenges: (i) due to row-wise sparsity in the measurements, the sensing operator does not satisfy the commonly-used Restricted Isometry Property (RIP) within

an efficient number of samples, i.e., it does not efficiently concentrate to an identity operation on all rank-k matrices, and (ii) FedRep executes a novel non-symmetric procedure. We further discuss these challenges in Appendix B.5. To the best of our knowledge, Theorem 1 provides the first convergence result for an alternating minimization-descent procedure to solve a matrix sensing problem. It is also the first result to show sample-efficient linear convergence of any solution to a matrix sensing with rank-one, row-wise sparse measurements. The state-of-the-art result for the closest matrix sensing setting to ours is given by Zhong et al. [2015] for rank-1, independent Gaussian measurements, which our result matches up to an $\mathcal{O}(\kappa^2)$ factor. However, our setting is more challenging as we have rank-1 and row-wise sparse measurements, and dependence on κ^4 has been previously observed in settings with sparse measurements, e.g. matrix completion [Jain et al., 2013].

Representation learning, dimensionality reduction and new users. Theorem 1 concerns a linear representation learning setting that is of interest beyond federated learning to representation learning problems more broadly, such as in meta-learning and multi-task learning. This setting has garnered significant attention recently in large part due to empirical evidence that representation learning can explain the success of meta-learning methods on few-shot learning tasks [Raghu et al., 2019]. As shown by Maurer et al. [2016], Du et al. [2020] and Tripuraneni et al. [2020a], learning an accurate k-dimensional representation during training (or meta-training) reduces the sample complexity of solving a new task from $\Theta(d)$ to $\Theta(k)$ in the linear case, enabling strong few-shot performance if k is small. Theorem 1 shows that FedRep learns an accurate k-dimensional representation during training in the linear case, so these prior results imply that FedRep also needs only $\Theta(k)$ samples to learn the model for the new client. Further, Theorem 1 shows that alternating minimization-descent (FedRep in the linear case) efficiently learns the representation compared to the methods studied in prior works (see Section 1.1 for a detailed comparison).

Remark on initialization. Theorem 1 requires that the initial principal angle distance dist($\mathbf{B}^0, \mathbf{B}^*$) is bounded away from 1 by a constant. This can be efficiently achieved by the Method of Moments without increasing the sample complexity for each client up to log factors [Tripuraneni et al., 2020a]. In turn, each user must send the server a polynomial of their data, namely $\sum_{j=1}^{m} (y_i^j)^2 \mathbf{x}_i^j (\mathbf{x}_i^j)^{\top}$ at the start of the learning procedure, which does not compromise privacy. We discuss the details of this in Appendix B.

5 Experiments

We focus on three points in our experiments: (i) the effect of many local updates for the local head in FedRep (ii) the quality of the global representation learned by FedRep and (iii) the applicability of FedRep to a wide range of datasets. Full experimental details are provided in Appendix A.

5.1 Synthetic Data

We start by experimenting with an instance of the multi-linear regression problem analyzed in Section 4. Consistent with this formulation, we generate synthetic samples $\mathbf{x}_i^j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and labels $y_i^j \sim \mathcal{N}(\mathbf{w}_i^{*^{\top}} \mathbf{B}^{*^{\top}} \mathbf{x}_i^j, 10^{-3})$ (here we include an additive Gaussian noise). The ground-truth heads $\mathbf{w}_i^* \in \mathbb{R}^k$ for clients $i \in [n]$ and the ground-truth representation $\mathbf{B}^* \in \mathbb{R}^{d \times k}$ are generated randomly

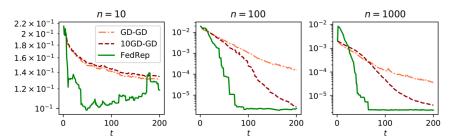


Figure 3: Comparison of (principal angle) distances between the ground-truth and estimated representations by FedRep and alternating gradient descent algorithms for different numbers of clients n. In all plots, d = 10, k = 2, m = 5, and r = 0.1.

by sampling and normalizing Gaussian matrices.

Benefit of finding the optimal head. We first demonstrate that the convergence of FedRep improves with larger number of clients n, making it highly applicable to federated settings. Further, we give evidence showing that this improvement is augmented by the minimization step in FedRep, since methods that replace the minimization step in FedRep with 1 and 10 steps of gradient descent (GD-GD and 10GD-GD, respectively) do not scale properly with n. In Figure 3, we plot convergence trajectories for FedRep, GD-GD, and 10GD-GD for four different values of n and fixed m, d, k and r. As we observe in Figure 3, by increasing the number of nodes n, clients converge to the true representation faster. Also, running more local updates for finding the local head accelerates the convergence speed of FedRep. In particular, FedRep which exactly finds the optimal local head at each round has the fastest rate compared to GD-GD and 10GD-GD that only run 1 and 10 local updates, respectively, to learn the head.

Generalization to new clients. Next, we evaluate the effectiveness of the representation learned by FedRep in reducing the sample complexity for a new client which has not participated in training. We compare against FedSGD, which executes distributed SGD to learn a single model (\mathbf{B} , \mathbf{w}). We first train FedRep and FedSGD on a fixed set of n = 100 clients as in Figure 1, where (d, k) = (20, 2). The new client has access to m_{new} labeled local samples. It will use the representation $\mathbf{B}^* \in \mathbb{R}^{d \times k}$ learned from the training clients, and learns a personalized head using this representation and its local training samples. For both FedRep and FedSGD, we solve for the optimal head given these samples and the representation learned during training. We compare the MSE of the resulting model on the new client's test data to that of a model trained by only using the m_{new} labeled samples from the new client (Local Only) in Figure 4. The large error for FedSGD demonstrates that it does not learn the ground-truth representation. Meanwhile, the representation learned by FedRep allows an accurate model to be found for the new client as long as $m_{\text{new}} \geq k$, which drastically improves over the complexity for Local Only ($m_{\text{new}} = \Omega(d)$).

5.2 Real Data

We next investigate whether these insights apply to nonlinear models and real datasets.

Datasets and Models. We use four real datasets: CIFAR10 and CIFAR100 [Krizhevsky et al., 2009], FEMNIST [Caldas et al., 2018, Cohen et al., 2017] and Sent140 [Caldas et al., 2018]. The first

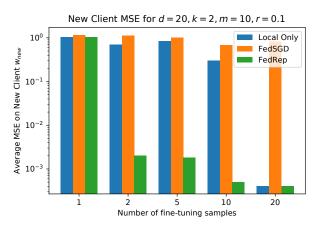


Figure 4: MSE on new clients sharing the representation after fine-tuning using various numbers of samples from the new client.

three are image datasets and the last is a text dataset for which the goal is to classify the sentiment of a tweet as positive or negative. We control the heterogeneity of CIFAR10 and CIFAR100 by assigning different numbers S of classes per client, from among 10 and 100 total classes, respectively. Each client is assigned the same number of training samples, namely 50000/n. For FEMNIST, we restrict the dataset to 10 handwritten letters and assign samples to clients according to a log-normal distribution as in Li et al. [2019]. We consider a partition of n = 150 clients with an average of 148 samples/client. For Sent140, we use the natural assignment of tweets to their author, and use n = 183 clients with an average of 72 samples per client. We use 5-layer CNNs for the CIFAR datasets, a 2-layer MLP for FEMNIST, and an RNN for Sent140 (details provided in Appendix A).

Baselines. We compare against a variety of personalized federated learning techniques as well as methods for learning a single global model and their fine-tuned analogues. Among the personalized methods, FedPer [Arivazhagan et al., 2019] is most similar to ours, as it also learns a global representation and personalized heads, but makes simultaneous local updates for both sets of parameters, therefore makes the same number of local updates for the head and the representation on each local round. Fed-MTL [Smith et al., 2017] learns local models and a regularizer to encode relationships among the clients, PerFedAvg [Fallah et al., 2020] leverages meta-learning to learn a single model that performs well after adaptation on each task, and LG-FedAvg [Liang et al., 2020] learns local representations and a global head. APFL [Deng et al., 2020] interpolates between local and global models, and L2GD [Hanzely and Richtárik, 2020] and Ditto [Li et al., 2020] learn local models that are encouraged to be close together by global regularization. For global FL methods, we consider FedAvg [McMahan et al., 2017], SCAFFOLD [Karimireddy et al., 2020], and FedProx [Li et al., 2018]. To obtain fine-tuning results, we first train the global model for the full training period, then each client then fine-tunes only the head on its local training data for 10 epochs of SGD before computing the final test accuracy.

Implementation. In each experiment we sample a ratio r = 0.1 of all the clients on every round. We initialize all models randomly and train for T = 100 communication rounds for the CIFAR datasets, T = 50 for Sent140, and T = 200 for FEMNIST. In each case, for each local updates FedRep executes ten local epochs of SGD with momentum to train the local head, followed by one epoch

Table 1: Average test accuracies on various partitions of CIFAR10, CIFAR100, Sent140 and FEMNIST with participation rate r=0.1.

		CIFAR10)	CIFA	AR100	Sent140	FEMNIST
(# clients n, # classes per client S)	(100, 2)	(100, 5)	(1000, 2)	(100, 5)	(100, 20)	(183, 2)	(150, 3)
Local Only	89.79	70.68	78.30	75.29	41.29	69.88	60.86
FedAvg [McMahan et al., 2017]	42.65	51.78	44.31	23.94	31.97	52.75	51.64
FedAvg+FT	87.65	73.68	82.04	79.34	55.44	71.92	72.41
FedProx [Li et al., 2018]	39.92	50.99	21.93	20.17	28.52	52.33	18.89
FedProx+FT	85.81	72.75	75.41	78.52	55.09	71.21	53.54
SCAFFOLD [Karimireddy et al., 2020]	37.72	47.33	33.79	20.32	22.52	51.31	17.65
SCAFFOLD+FT	86.35	68.23	78.24	78.88	44.34	71.49	52.11
Fed-MTL [Smith et al., 2017]	80.46	58.31	76.53	71.47	41.25	71.20	54.11
PerFedAvg [Fallah et al., 2020]	82.27	67.20	67.36	72.05	52.49	68.45	71.51
LG-Fed [Liang et al., 2020]	84.14	63.02	77.48	72.44	38.76	70.37	62.08
L2GD [Hanzely and Richtárik, 2020]	81.04	59.98	71.96	72.13	42.84	70.67	66.18
APFL [Deng et al., 2020]	83.77	72.29	82.39	78.20	55.44	69.87	70.74
Ditto [Li et al., 2020]	85.39	70.34	80.36	78.91	56.34	71.04	68.28
FedPer [Arivazhagan et al., 2019]	87.13	73.84	81.73	76.00	55.68	72.12	76.91
FedRep (Ours)	87.70	75.68	83.27	79.15	56.10	72.41	78.56

for the representation in the case of CIFAR10 with $n\!=\!100$ and 5 epochs in all other cases. All other methods use the same number of local epochs as FedRep does for updating the representation. Accuracies are computed by taking the average local accuracies for all users over the final 10 rounds of communication, except for the fine-tuning methods. These accuracies are computed after locally training the head of the fully-trained global model for ten epochs for each client.

Benefit of more local updates. As mentioned in Section 1, a key advantage of our formulation is that it enables clients to run many local updates without causing divergence from the global optimal solution. We demonstrate an example of this in Figure 5. Here, there are n=100 clients where each has S=2 classes of images. For FedAvg, we observe running more local updates does not necessarily improve the performance. In contrast, FedRep's performance is monotonically non-decreasing with the number of local epochs for the heads, i.e., FedRep is never hurt by more local computation on the heads.

Robustness to varying levels of heterogeneity, number of clients and number of samples per client. We show the average local test errors for all of the algorithms for a variety of settings in Table 1. Recall that for the CIFAR datasets, the number of training samples per client is equal to 50000/n, so the columns with 100 clients have 500 training samples per client, and the column with 1000 clients has only 50 training samples per client. In all cases, FedRep is either the top-performing method or is very close to the top-performing method. Surprisingly, the fine-tuning methods perform very well, especially FedAvg+FT. The superior performance of FedAvg relative SCAFFOLD is likely because all settings involve partial client participation.

Generalization to new clients. We also evaluate the strength of the representation learned by FedRep in terms of adaptation for new users. To do so, we first train FedRep, FedAvg, PerFedAvg, LG-FedAvg, APFL, L2GD and FedProx in the usual setting on the partition of FEMNIST containing

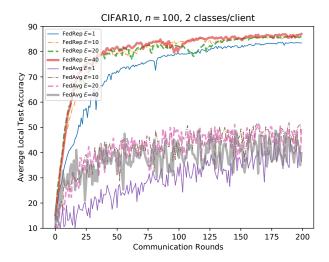


Figure 5: CIFAR10 local test errors for different numbers of local epochs E for FedAvg and for the heads in FedRep.

images of 10 handwritten letters (FEMNIST-letters). Then, we encounter clients with data from a different partition of the FEMNIST dataset, containing images of handwritten digits. We assume we have access to a dataset of 500 samples at this new client to fine tune the head. Using these, with each of the algorithms, we fine tune the head over multiple epochs while keeping the representation fixed. In Figure 6, we repeatedly sweep over the same 500 samples over multiple epochs to further refine the head, and plot the corresponding local test accuracy. As is apparent, FedRep has significantly better performance than these baselines.

6 Discussion

We introduce a novel representation learning framework and algorithm for federated learning, and we provide both theoretical and empirical justification for its utility in federated settings. In particular, our proposed framework exploits the structure of federating learning by (i) leveraging all clients' data to learn a global representation that enhances each client's model and can generalize to new users and (ii) leveraging the computational power of clients to run multiple local updates for learning their local heads. Our analysis further shows that alternating minimization-descent efficiently learns linear representations, and is therefore relevant beyond federated learning. Future work remains to analyze the representation learning capabilities of FedRep in non-linear settings.

7 Acknowledgements

The research of Liam Collins is supported through ARO Grant W911NF-11-1-0265 and NSF Grant 2019844. The research of Sanjay Shakkottai is supported by ONR Grant N00014-19-1-2566 and NSF Grant 2019844. The research of Aryan Mokhtari is supported in part by NSF Grant

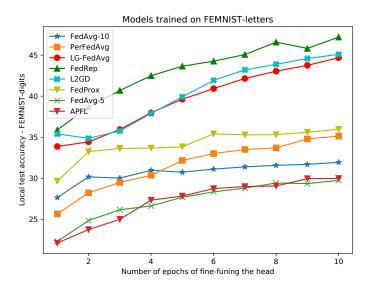


Figure 6: Test accuracy on handwritten digits from FEMNIST after fine-tuning the head of models trained on FEMNIST-letters.

 $2007668,\,\mathrm{ARO}$ Grant W911NF2110226, and the Machine Learning Laboratory at UT Austin. The research of Hamed Hassani is supported by NSF Grants 1837253, 1943064, 1934876, AFOSR Grant FA9550-20-1-0111, and DCIST-CRA.

A Additional Experimental Results

A.1 Synthetic Data: Further comparison with GD-GD

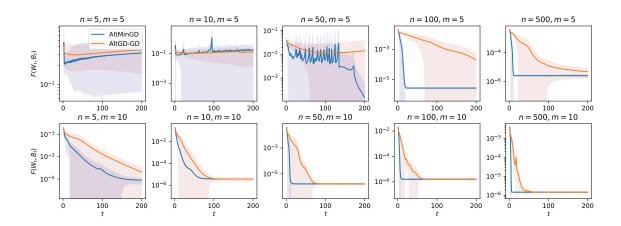


Figure 7: Function values for FedRep and GD-GD. The value of m is fixed in each row and n is fixed in each column. Here r=1 (full participation) and the average trajectories over 10 trials are plotted along with 95% confidence intervals. Principal angle distances are not plotted as the results are very similar. We see that the relative improvement of FedRep over GD-GD increases with n, highlighting the advantage of FedRep in settings with many clients.

Further experimental details. In the synthetic data experiments, the ground-truth matrices \mathbf{W}^* and \mathbf{B}^* were generated by first sampling each element as an i.i.d. standard normal variable, then taking the QR factorization of the resulting matrix, and scaling it by \sqrt{k} in the case of \mathbf{W}^* . The clients each trained on the same m samples throughout the entire training process. Test samples were generated identically as the training samples but without noise. Both the iterates of FedRep and GD-GD were initialized with the SVD of the result of 10 rounds of projected gradient descent on the unfactorized matrix sensing objective as in Algorithm 1 in Tu et al. [2016]. We would like to note that FedRep exhibited the same convergence trajectories regardless of whether its iterates were initialized with random Gaussian samples or with the projected gradient descent procedure, whereas GD-GD was highly sensitive to its initialization, often not converging when initialized randomly.

A.2 Real Data: Further experimental details

Datasets. The CIFAR10 and CIFAR100 datasets [Krizhevsky et al., 2009] were generated by randomly splitting the training data into Sn shards with 50,000/(Sn) images of a single class in each shard, as in McMahan et al. [2017]. The full Federated-EMNIST (FEMNIST) dataset contains 62 classes of handwritten letters, but in Table 1 we use a subset with only 10 classes of handwritten letters. In particular, we followed the same dataset generation procedure as in Li et al. [2019], but used 150 clients instead of 200. When testing on new clients as in Figure 6, we use samples from 10 classes of handwritten digits from FEMNIST, i.e., the MNIST dataset. In this phase there are 100

new clients, each with 500 samples from 5 different classes for fine-tuning. The fine-tuned models are then evaluated on 100 testing samples from these same 5 classes. For Sent140, we randomly sample 183 clients (Twitter users) that each have at least 50 samples (tweets). Each tweet is either positive sentiment or negative sentiment. Statistics of both the FEMNIST and Sent140 datasets we use are given in Table 2. For both FEMNIST and Sent140 we use the LEAF framework [Caldas et al., 2018].

Hyperparameters. As in Liang et al. [2020], all methods use SGD with momentum with parameter equal to 0.5. In Table 1, for CIFAR10, CIFAR100, and FEMNIST the local sample batch size is 10 and for Sent140 it is 4. The participation rate r is always 0.1, besides in the fine-tuning phases in Figure 6, in which all clients are sampled in each round. For each dataset learning rates were tuned in $\{0.001, 0.01, 0.01, 0.1\}$. We observed that the optimal learning rates for FedAvg were also typically the optimal base learning rates for the other methods, so we used the same base learning rates for all methods for each dataset, which was 0.01 in all cases, unless stated otherwise. Note that the batch size and learning rate for CIFAR10 used in Table 1 differs from the standard setting of a batch size of 50 and learning rate of 0.1 [McMahan et al., 2017], but we observed improved performance for all methods by using (10,0.01) instead. In particular, the simulation in Figure 5, the standard setting of (50,0.1) is used, but the accuracies are worse than those reported in Table 1 for both FedAvg and FedRep. Additionally, in Table 1, for CIFAR10 with (n, S) = (100, 2) and (n, S) = (100, 5), we executed 1 local epoch of SGD with momentum for the representation for FedRep and 1 local epoch for all other methods. For all other datasets we executed 5 local epochs for the representation for FedRep and for the local updates for all other methods.

Evaluation. As mentioned in the main body, in Table 1, we initialize all methods randomly and train for T=100 communication rounds for the CIFAR datasets, T=200 for FEMNIST, and T=50 for Sent140. The accuracy shown is the average local test accuracy over all users over the final ten communication rounds, besides for the fine-tuning results, in which case we report the average local test accuracies of the locally fine-tuned models over all users, after the global model has been fully trained. We repeat the entire training and evaluation process five times for each model and dataset and report the averages in Table 1.

Implementations. Our code is available at https://github.com/lgcollins/FedRep. We adapt the Pytorch codebase from Liang et al. [2020], and used the implementations of FedAvg, Fed-MTL and LG-FedAvg from this repository. For consistency we use this same codebase to implement FedRep, FedPer, SCAFFOLD, FedProx, APFL, Ditto, L2GD, and Per-FedAvg. As in the experiments in Liang et al. [2020], we used a 5-layer CNN with two convolutional layers for CIFAR10 and CIFAR100 followed by three fully-connected layers. For FEMNIST, we use an MLP with two hidden layers, and for Sent140 we use a pre-trained 300-dimensional GloVe embedding and train RNN with an LSTM module followed by two fully-connected decoding layers.

For FedRep, we treated the head as the weights and biases of the final fully-connected layer in each of the models. For LG-FedAvg, we treated the first two convolutional layers of the model for CIFAR10 and CIFAR100 as the local representation, and the fully-connected layers as the global parameters, and the input layer and hidden layers as the global parameters. For FEMNIST, we set all parameters besides those in the output layer as the local representation parameters. For Sent140, we set the RNN module to be the local representation and the decoder to be the global

¹Pennington, J., Socher, R., and Manning, C. D. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532–1543, 2014.

parameters. Unlike in the paper introducing LG-FedAvg [Liang et al., 2020], we did not initialize the models for all methods with the solution of many rounds of FedAvg (instead, we initialized randomly) and we computed the local test accuracy as the average local test accuracy over the final ten communication rounds, rather than the average of the maximum local test accuracy for each client over the entire training procedure.

For L2GD we executed multiple epochs of local SGD (discussed above) instead of one step of GD in the local update in order for reasonable comparison with the other methods. We also set p=0.9, thus the local parameters are trained on 10% of the communication rounds. We tuned α in $\{0.05,0.1,0.25,0.5,0.5,0.75\}$ and we tuned λ over $\{1,0.5\}$. We used $(\alpha,\lambda)=(0.25,1)$ in all cases besides the (n,S)=(100,5) case for CIFAR100, for which we used $\alpha=0.1$. Also, for FEMNIST we improved performance by using a learning rate of 0.001 instead of 0.01. For APFL, we used a fixed α that we tuned in $\{0.1,0.25,0.5,0.75\}$, and chose $\alpha=0.25$ for all cases besides the most heterogeneous CIFAR versions, namely (n,S)=(100,2) for CIFAR10 and (n,S)=(100,25) for CIFAR100. For Ditto we tuned λ among $\{0.25,0.5,0.75,1\}$, and used $\lambda=0.75$ for all cases besides CIFAR100, for which we used $\lambda=1$. For PerFedAvg, we used an inner learning rate of 10^{-4} and 8 samples as the support set and 2 samples as the target set in each local meta-gradient update. We used the Hessian-free version. For FedProx we tuned μ among $\{0.05,0.1,0.25,0.5\}$, and used $\mu=0.1$ for CIFAR and $\mu=0.25$ for FEMNIST and Sent140. For SCAFFOLD we used a global learning rate of 1 in all cases besides FEMNIST, for which 0.5 was superior.

Table 2: Dataset statistics.

DATASET	Number of users (n)	AVG SAMPLES/USER	Min samples/user
FEMNIST	150	148	50
SENT140	183	72	50

B Proof of Main Result

B.1 Preliminaries.

We start by defining some notions used throughout the proof.

Definition 2. For a random vector $\mathbf{x} \in \mathbb{R}^d$ and a fixed matrix $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$, the vector $\mathbf{A}^\top \mathbf{x}$ is called $\|\mathbf{A}\|_2$ -sub-gaussian if $\mathbf{y}^\top \mathbf{A}^\top \mathbf{x}$ is sub-gaussian with sub-gaussian norm $\mathcal{O}(\|\mathbf{A}\|_2 \|\mathbf{y}\|_2)$ for all $\mathbf{y} \in \mathbb{R}^{d_2}$, i.e. $\mathbb{E}[\exp(\mathbf{y}^\top \mathbf{A}^\top \mathbf{x})] \leq \exp(\|\mathbf{y}\|_2^2 \|\mathbf{A}\|_2^2/2)$.

Definition 3. A rank-k matrix $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$ is μ -row-wise incoherent if $\max_{i \in [d_1]} \|\mathbf{m}_i\|_2 \leq (\mu \sqrt{d_2}/\sqrt{d_1}) \|\mathbf{M}\|_F$, where $\mathbf{m}_i \in \mathbb{R}^{d_2}$ is the *i*-th row of \mathbf{M} .

Note that Assumption 3 implies that \mathbf{W}^* is row-wise incoherent with parameter 1.

We use hats to denote orthonormal matrices (a matrix is called orthonormal if its set of columns is an orthonormal set). By Assumption 3, the ground truth representation \mathbf{B}^* is orthonormal, so from now on we will write it as $\hat{\mathbf{B}}^*$. Likewise, we will denote the iterates \mathbf{B}^t as $\hat{\mathbf{B}}^t$.

For a matrix $\mathbf{W} \in \mathbb{R}^{n \times k}$ and a random set of indices $\mathcal{I} \in [n]$ of cardinality rn, define $\mathbf{W}_{\mathcal{I}} \in \mathbb{R}^{rn \times k}$ as the matrix formed by taking the rows of \mathbf{W} indexed by \mathcal{I} . Define $\bar{\sigma}_{\max,*} \coloneqq \max_{\mathcal{I} \in [n], |\mathcal{I}| = rn} \sigma_{\max}(\frac{1}{\sqrt{rn}}\mathbf{W}_{\mathcal{I}}^*)$ and $\bar{\sigma}_{\min,*} \coloneqq \min_{\mathcal{I} \in [n], |\mathcal{I}| = rn} \sigma_{\min}(\frac{1}{\sqrt{rn}}\mathbf{W}_{\mathcal{I}}^*)$, i.e. the maximum and minimum singular values of any matrix that can be obtained by taking rn rows of $\frac{1}{\sqrt{rn}}\mathbf{W}^*$. Note that by Assumption 3, each row of \mathbf{W}^* has norm \sqrt{k} , so $\frac{1}{\sqrt{rn}}$ acts as a normalizing factor such that $\|\frac{1}{\sqrt{rn}}\mathbf{W}_{\mathcal{I}}^*\|_F = \sqrt{k}$. In addition, define $\kappa = \bar{\sigma}_{\max,*}/\bar{\sigma}_{\min,*}$.

Let i now be an index over [rn], and let i' be an index over [n]. For random batches of samples $\{\{(\mathbf{x}_i^j, y_i^j)\}_{j=1}^m\}_{i=1}^{rn}$, define the random linear operator $\mathcal{A}: \mathbb{R}^{rn \times d} \to \mathbb{R}^{rnm}$ as $\mathcal{A}(\mathbf{M}) = [\langle \mathbf{A}_{i,j}, \mathbf{M} \rangle]_{1 \leq i \leq rn, 1 \leq j \leq m} \in \mathbb{R}^{rnm}$. Here, $\mathbf{A}_{i,j} := \mathbf{e}_i(\mathbf{x}_i^j)^{\top}$, where \mathbf{e}_i is the i-th standard vector in \mathbb{R}^{rn} , and $\mathbf{M} \in \mathbb{R}^{rn \times d}$. Then, the loss function in (7) is equivalent to

$$\min_{\mathbf{B} \in \mathbb{R}^{d \times k}, \mathbf{W} \in \mathbb{R}^{n \times k}} \{ F(\mathbf{B}, \mathbf{W}) := \frac{1}{2rnm} \mathbb{E}_{\mathcal{A}, \mathcal{I}} \left[\| \mathbf{Y} - \mathcal{A}(\mathbf{W}_{\mathcal{I}} \mathbf{B}^{\top}) \|_{2}^{2} \right] \},$$
(13)

where $\mathbf{Y} = \mathcal{A}(\mathbf{W}_{\mathcal{I}}^* \hat{\mathbf{B}}^{*^{\top}}) \in \mathbb{R}^{rnm}$ is a concatenated vector of labels. It is now easily seen that the problem of recovering $\mathbf{W}^* \hat{\mathbf{B}}^{*^{\top}}$ from finitely-many measurements $\mathcal{A}(\mathbf{W}_{\mathcal{I}}^* \hat{\mathbf{B}}^{*^{\top}})$ is an instance of matrix sensing. Moreover, the updates of FedRep satisfy the following recursion:

$$\mathbf{W}_{\mathcal{I}^{t}}^{t+1} = \underset{\mathbf{W}_{\mathcal{I}^{t}} \in \mathbb{R}^{rn \times k}}{\operatorname{argmin}} \frac{1}{2rnm} \| \mathcal{A}^{t} (\mathbf{W}_{\mathcal{I}^{t}}^{*} \hat{\mathbf{B}}^{*^{\top}} - \mathbf{W}_{\mathcal{I}^{t}} \hat{\mathbf{B}}^{t^{\top}}) \|_{2}^{2}$$

$$(14)$$

$$\bar{\mathbf{B}}^{t+1} = \hat{\mathbf{B}}^t - \frac{\eta}{rnm} \left((\mathcal{A}^t)^{\dagger} \mathcal{A}^t (\mathbf{W}_{\mathcal{I}^t}^{t+1} \hat{\mathbf{B}}^{t^{\top}} - \mathbf{W}_{\mathcal{I}^t}^* \hat{\mathbf{B}}^{*^{\top}}) \right)^{\top} \mathbf{W}_{\mathcal{I}^t}^{t+1}$$
(15)

$$\hat{\mathbf{B}}^{t+1}, \mathbf{R}^{t+1} = \mathrm{QR}(\bar{\mathbf{B}}^t) \tag{16}$$

where \mathcal{A}^t is an instance of \mathcal{A} , $(\mathcal{A}^t)^{\dagger}$ is the adjoint operator of \mathcal{A}^t , i.e. $(\mathcal{A}^t)^{\dagger}(\mathbf{M}) = \sum_{i=1}^{rn} \sum_{j=1}^{m} (\langle \mathbf{A}_i^{t,j}, \mathbf{M} \rangle) \mathbf{A}_i^{t,j}$, and $\mathrm{QR}(\cdot)$ is the QR factorization. Note that for the purposes of analysis, it does not matter how $\mathbf{w}_{i'}^{t+1}$ is computed for all $i' \notin \mathcal{I}^t$, as these vectors do not affect the computation of \mathbf{B}^{t+1} . Moreover, our analysis does not rely on any particular properties of the batches $\mathcal{I}^1, \ldots, \mathcal{I}^T$ other than the fact that they have cardinality rn, so without loss of generality we assume $\mathcal{I}^t = [rn]$ for all $t = 1, \ldots T$ and drop the subscripts \mathcal{I}^t on \mathbf{W}^t . Further, since our analysis focuses on a particular iteration t, we will drop the superscript t on \mathcal{A}^t and each $\mathbf{A}_i^{t,j}$ and $(\mathbf{x}_i^{t,j}, y_i^{t,j})$ for ease of notation (while noting that each iteration requires a new batch of i.i.d. data).

B.2 Auxilliary Lemmas

We start by computing the update for \mathbf{W} .

Lemma 1. In the linear version of FedRep, update for **W** is:

$$\mathbf{W}^{t+1} = \mathbf{W}^* \hat{\mathbf{B}}^{*^{\top}} \hat{\mathbf{B}}^t - \mathbf{F}$$
 (17)

where \mathbf{F} is defined in equation (22) below.

Proof. We adapt the argument from Lemma 4.5 in [Jain et al., 2013] to compute the update for \mathbf{W}^{t+1} , and borrow heavily from their notation.

Let \mathbf{w}_p^{t+1} (respectively $\hat{\mathbf{b}}_p^{t+1}$) be the p-th column of \mathbf{W}^t (respectively $\hat{\mathbf{B}}^t$). Since \mathbf{W}^{t+1} minimizes $\tilde{F}(\mathbf{W}, \hat{\mathbf{B}}^t) \coloneqq \frac{1}{2rnm} \|\mathcal{A}(\mathbf{W}^*(\hat{\mathbf{B}}^*)^\top - \mathbf{W}(\mathbf{B}^t)^\top)\|_2^2$ with respect to \mathbf{W} , we have $\nabla_{\mathbf{w}_p} \tilde{F}(\mathbf{W}^{t+1}, \hat{\mathbf{B}}^t) = \mathbf{0}$ for all $p \in [k]$. Thus, for any $p \in [k]$, we have

$$\begin{aligned} \mathbf{0} &= \nabla_{\mathbf{w}_p} \tilde{F}(\mathbf{W}^{t+1}, \hat{\mathbf{B}}^t) \\ &= \frac{1}{rnm} \sum_{i=1}^{rn} \sum_{j=1}^{m} \left(\langle \mathbf{A}_{i,j}, \mathbf{W}^{t+1} (\hat{\mathbf{B}}^t)^\top - \mathbf{W}^* (\hat{\mathbf{B}}^*)^\top \rangle \right) \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t \\ &= \frac{1}{rnm} \sum_{i=1}^{rn} \sum_{j=1}^{m} \left(\sum_{q=1}^{k} (\hat{\mathbf{b}}_q^t)^\top \mathbf{A}_{i,j}^\top \mathbf{w}_q^{t+1} - \sum_{q=1}^{k} (\hat{\mathbf{b}}_q^*)^\top \mathbf{A}_{i,j}^\top \mathbf{w}_q^* \right) \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t \end{aligned}$$

This implies

$$\frac{1}{m} \sum_{q=1}^{k} \left(\sum_{i=1}^{rn} \sum_{j=1}^{m} \mathbf{A}_{i,j} \hat{\mathbf{b}}_{p}^{t} (\hat{\mathbf{b}}_{q}^{t})^{\top} \mathbf{A}_{i,j}^{\top} \right) \mathbf{w}_{q}^{t+1} = \frac{1}{m} \sum_{q=1}^{k} \left(\sum_{i=1}^{rn} \sum_{j=1}^{m} \mathbf{A}_{i,j} \hat{\mathbf{b}}_{p}^{t} (\hat{\mathbf{b}}_{q}^{*})^{\top} \mathbf{A}_{i,j}^{\top} \right) \mathbf{w}_{q}^{*}$$
(18)

To solve for \mathbf{w}^{t+1} , we define \mathbf{G} , \mathbf{C} , and \mathbf{D} as rnk-by-rnk block matrices, as follows:

$$\mathbf{G} \coloneqq \begin{bmatrix} \mathbf{G}_{11} & \cdots & \mathbf{G}_{1k} \\ \vdots & \ddots & \vdots \\ \mathbf{G}_{k1} & \cdots & \mathbf{G}_{kk} \end{bmatrix}, \mathbf{C} \coloneqq \begin{bmatrix} \mathbf{C}_{11} & \cdots & \mathbf{C}_{1k} \\ \vdots & \ddots & \vdots \\ \mathbf{C}_{k1} & \cdots & \mathbf{C}_{kk} \end{bmatrix}, \mathbf{D} \coloneqq \begin{bmatrix} \mathbf{D}_{11} & \cdots & \mathbf{D}_{1k} \\ \vdots & \ddots & \vdots \\ \mathbf{D}_{k1} & \cdots & \mathbf{D}_{kk} \end{bmatrix}$$
(19)

where, for $p, q \in [k]$: $\mathbf{G}_{pq} \coloneqq \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^{m} \mathbf{A}_{i,j} \hat{\mathbf{b}}_{p}^{t} \hat{\mathbf{b}}_{q}^{\mathsf{T}} \mathbf{A}_{i,j}^{\mathsf{T}} \in \mathbb{R}^{rn \times rn}$, $\mathbf{C}_{pq} \coloneqq \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^{m} \mathbf{A}_{i,j} \hat{\mathbf{b}}_{p}^{t} (\hat{\mathbf{b}}_{q}^{*})^{\mathsf{T}} \mathbf{A}_{i,j}^{\mathsf{T}} \in \mathbb{R}^{rn \times rn}$, and, $\mathbf{D}_{pq} \coloneqq \langle \hat{\mathbf{b}}_{p}^{t}, \hat{\mathbf{b}}_{q}^{*} \rangle \mathbf{I}_{rn} \in \mathbb{R}^{rn \times rn}$. Recall that $\hat{\mathbf{b}}_{p}^{t}$ is the p-th column of $\hat{\mathbf{B}}^{t}$ and $\hat{\mathbf{b}}_{q}^{*}$ is the q-th column of $\hat{\mathbf{B}}^{t}$. Further, define

$$\widetilde{\mathbf{w}}^{t+1} = \begin{bmatrix} \mathbf{w}_1^{t+1} \\ \vdots \\ \mathbf{w}_k^{t+1} \end{bmatrix} \in \mathbb{R}^{rnk}, \quad \widetilde{\mathbf{w}}^* = \begin{bmatrix} \mathbf{w}_1^* \\ \vdots \\ \mathbf{w}_k^* \end{bmatrix} \in \mathbb{R}^{rnk}.$$

Then, by (18), we have

$$\begin{split} \widetilde{\mathbf{w}}^{t+1} &= \mathbf{G}^{-1} \mathbf{C} \widetilde{\mathbf{w}}^* \\ &= \mathbf{D} \widetilde{\mathbf{w}}^* - \mathbf{G}^{-1} \left(\mathbf{G} \mathbf{D} - \mathbf{C} \right) \widetilde{\mathbf{w}}^* \end{split}$$

where we can invert **G** conditioned on the event that its minimum singular value is strictly positive, which Lemma 2 shows holds with high probability. Now consider the p-th block of $\widetilde{\mathbf{w}}^{t+1}$, and let $((\mathbf{G}\mathbf{D} - \mathbf{C})\mathbf{w}^*)_p$ denote the p-th block of $(\mathbf{G}\mathbf{D} - \mathbf{C})\mathbf{w}^*$. We have

$$\widetilde{\mathbf{w}}_{p}^{t+1} = \sum_{q=1}^{k} \langle \hat{\mathbf{b}}_{p}^{t}, \hat{\mathbf{b}}_{q}^{*} \rangle \mathbf{w}_{q}^{*} - (\mathbf{G}^{-1} (\mathbf{G}\mathbf{D} - \mathbf{C}) \mathbf{w}^{*})_{p}$$

$$= \left(\sum_{q=1}^{k} \mathbf{w}_{q}^{*} (\hat{\mathbf{b}}_{p}^{*})^{\top} \right) \hat{\mathbf{b}}_{q}^{t} - (\mathbf{G}^{-1} (\mathbf{G}\mathbf{D} - \mathbf{C}) \mathbf{w}^{*})_{p}$$

$$= \left(\mathbf{W}^{*} (\hat{\mathbf{B}}^{*})^{\top} \right) \hat{\mathbf{b}}_{q}^{t} - (\mathbf{G}^{-1} (\mathbf{G}\mathbf{D} - \mathbf{C}) \mathbf{w}^{*})_{p}$$

$$(20)$$

By constructing \mathbf{W}^{t+1} such that the p-th column of \mathbf{W}^{t+1} is \mathbf{w}_p^{t+1} for all $p \in [k]$, we obtain

$$\mathbf{W}^{t+1} = \mathbf{W}^* \hat{\mathbf{B}}^* (\hat{\mathbf{B}}^t)^\top - \mathbf{F}$$
 (21)

where

$$\mathbf{F} = [(\mathbf{G}^{-1}(\mathbf{G}\mathbf{D} - \mathbf{C})\widetilde{\mathbf{w}}^*)_1, \dots, (\mathbf{G}^{-1}(\mathbf{G}\mathbf{D} - \mathbf{C})\widetilde{\mathbf{w}}^*)_k]$$
(22)

and $(\mathbf{G}^{-1}(\mathbf{G}\mathbf{D} - \mathbf{C})\widetilde{\mathbf{w}}^*)_p$ is the *p*-th *n*-dimensional block of the rnk-dimensional vector $\mathbf{G}^{-1}(\mathbf{G}\mathbf{D} - \mathbf{C})\widetilde{\mathbf{w}}^*$.

Next we bound the Frobenius norm of the matrix \mathbf{F} , which requires multiple steps. First, we establish some helpful notations. We drop superscripts indicating the iteration number t for simplicity.

Again let \mathbf{w}^* be the rnk-dimensional vector formed by stacking the columns of \mathbf{W}^* , and let $\hat{\mathbf{b}}_p$ (respectively $\hat{\mathbf{b}}_q^*$) be the p-th column of $\hat{\mathbf{B}}$ (respectively the q-th column of $\hat{\mathbf{B}}_*$). Recall that \mathbf{F} can be obtained by stacking $\mathbf{G}^{-1}(\mathbf{G}\mathbf{D} - \mathbf{C})\mathbf{w}^*$ into k columns of length n, i.e. $\text{vec}(\mathbf{F}) = \mathbf{G}^{-1}(\mathbf{G}\mathbf{D} - \mathbf{C})\mathbf{w}^*$. Further, $\mathbf{G} \in \mathbb{R}^{rnk \times rnk}$ is a block matrix whose blocks $\mathbf{G}_{pq} \in \mathbf{R}^{rn \times rn}$ for $p, q \in [k]$ are given by:

$$\mathbf{G}_{pq} = \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^{m} \mathbf{A}_{i,j} \hat{\mathbf{b}}_{p} \hat{\mathbf{b}}_{q}^{\top} \mathbf{A}_{i,j}^{\top}$$

$$= \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^{m} \mathbf{e}_{i} (\mathbf{x}_{i}^{j})^{\top} \hat{\mathbf{b}}_{p} \hat{\mathbf{b}}_{q}^{\top} \mathbf{x}_{i}^{j} \mathbf{e}_{i}^{\top}$$
(23)

So, each \mathbf{G}_{pq} is diagonal with diagonal entries

$$(\mathbf{G}_{pq})_{ii} = \frac{1}{m} \sum_{j=1}^{m} (\mathbf{x}_i^j)^\top \hat{\mathbf{b}}_p \hat{\mathbf{b}}_q^\top \mathbf{x}_i^j = \hat{\mathbf{b}}_p^\top \left(\frac{1}{m} \sum_{j=1}^{m} \mathbf{x}_i^j (\mathbf{x}_i^j)^\top \right) \hat{\mathbf{b}}_q$$
(24)

Define $\Pi^i := \frac{1}{m} \sum_{j=1}^m \mathbf{x}_i^j (\mathbf{x}_i^j)^{\top}$ for all $i \in [rn]$. Similarly as above, each block \mathbf{C}_{pq} of \mathbf{C} is diagonal with entries

$$(\mathbf{C}_{pq})_{ii} = \hat{\mathbf{b}}_p^{\top} \mathbf{\Pi}^i \hat{\mathbf{b}}_{*,q} \tag{25}$$

Analogously to the matrix completion analysis in [Jain et al., 2013], we define the following matrices, for all $i \in [rn]$:

$$\mathbf{G}^{i} := \left[\hat{\mathbf{b}}_{p}^{\top} \mathbf{\Pi}^{i} \hat{\mathbf{b}}_{q}\right]_{1 \leq p, q \leq k} = \hat{\mathbf{B}}^{\top} \mathbf{\Pi}^{i} \hat{\mathbf{B}}, \quad \mathbf{C}^{i} := \left[\hat{\mathbf{b}}_{p}^{\top} \mathbf{\Pi}^{i} \hat{\mathbf{b}}_{*, q}\right]_{1 \leq p, q \leq k} = \hat{\mathbf{B}}^{\top} \mathbf{\Pi}^{i} \hat{\mathbf{B}}_{*}$$
(26)

In words, \mathbf{G}^i is the $k \times k$ matrix formed by taking the *i*-th diagonal entry of each block \mathbf{G}_{pq} , and likewise for \mathbf{C}^i . Recall that \mathbf{D} also has diagonal blocks, in particular $\mathbf{D}_{pq} = \langle \hat{\mathbf{B}}_p, \hat{\mathbf{B}}_q^* \rangle \mathbf{I}_d$, thus we also define $\mathbf{D}^i := [\langle \hat{\mathbf{B}}_p, \hat{\mathbf{B}}_q^* \rangle]_{1 \leq p,q \leq k} = \hat{\mathbf{B}}^\top \hat{\mathbf{B}}_*$.

Using this notation we can decouple $\mathbf{G}^{-1}(\mathbf{G}\mathbf{D} - \mathbf{C})\mathbf{w}^*$ into i subvectors. Namely, let $\mathbf{w}_i^* \in \mathbb{R}^k$ be the vector formed by taking the ((p-1)rn+i)-th elements of \mathbf{w}^* for p=0,...,k-1, and similarly, let \mathbf{f}_i be the vector formed by taking the ((p-1)rn+i)-th elements of $\mathbf{G}^{-1}(\mathbf{G}\mathbf{D} - \mathbf{C})\mathbf{w}^*$ for p=0,...,k-1. Then

$$\mathbf{f}_i = (\mathbf{G}^i)^{-1} (\mathbf{G}^i \mathbf{D}^i - \mathbf{C}^i) \mathbf{w}_i^* \tag{27}$$

is the *i*-th row of **F**. Now we control $\|\mathbf{F}\|_F$.

Lemma 2. Let $\delta_k = c \frac{k^{3/2} \sqrt{\log(rn)}}{\sqrt{m}}$ for some absolute constant c, then

$$\|\mathbf{G}^{-1}\|_2 \le \frac{1}{1 - \delta_k}$$

with probability at least $1 - e^{-111k^3 \log(rn)}$.

Proof. We must lower bound $\sigma_{\min}(\mathbf{G})$. For some vector $\mathbf{z} \in \mathbb{R}^{rnk}$, let $\mathbf{z}^i \in \mathbb{R}^k$ denote the vector formed by taking the ((p-1)rn+i)-th elements of \mathbf{z} for p=0,...,k-1. Since \mathbf{G} is symmetric, we have

$$\begin{split} \sigma_{\min}(\mathbf{G}) &= \min_{\mathbf{z}: \|\mathbf{z}\|_2 = 1} \mathbf{z}^{\top} \mathbf{G} \mathbf{z} \\ &= \min_{\mathbf{z}: \|\mathbf{z}\|_2 = 1} \sum_{i=1}^{rn} (\mathbf{z}^i)^{\top} \mathbf{G}^i \mathbf{z}^i \\ &= \min_{\mathbf{z}: \|\mathbf{z}\|_2 = 1} \sum_{i=1}^{rn} (\mathbf{z}^i)^{\top} \hat{\mathbf{B}}^{\top} \mathbf{\Pi}^i \hat{\mathbf{B}} \mathbf{z}^i \\ &\geq \min_{i \in [rn]} \sigma_{\min}(\hat{\mathbf{B}}^{\top} \mathbf{\Pi}^i \hat{\mathbf{B}}) \end{split}$$

Note that the matrix $\hat{\mathbf{B}}^{\top} \mathbf{\Pi}^i \hat{\mathbf{B}}$ can be written as follows:

$$\hat{\mathbf{B}}^{\top} \mathbf{\Pi}^{i} \hat{\mathbf{B}} = \sum_{j=1}^{m} \frac{1}{\sqrt{m}} \hat{\mathbf{B}}^{\top} \mathbf{x}_{i}^{j} \left(\frac{1}{\sqrt{m}} \hat{\mathbf{B}}^{\top} \mathbf{x}_{i}^{j} \right)^{\top}$$
(28)

Let $\mathbf{v}_i^j \coloneqq \frac{1}{\sqrt{m}} \hat{\mathbf{B}}^\top \mathbf{x}_i^j$ for all $i \in [rn]$ and $j \in [m]$, and note that each \mathbf{v}_i^j is i.i.d. $\frac{1}{\sqrt{m}} \hat{\mathbf{B}}$ -sub-gaussian. Thus using the one-sided version of equation (4.22) (Theorem 4.6.1) in [Vershynin, 2018], we have

$$\sigma_{min}(\hat{\mathbf{B}}^{\top} \mathbf{\Pi}^{i} \hat{\mathbf{B}}) \ge 1 - C \left(\sqrt{\frac{k}{m}} + \frac{z}{\sqrt{m}} \right)$$
 (29)

with probability at least $1 - e^{-z^2}$ for $m \ge k$, $z \ge 0$ and some absolute constant C. Now let $\delta_k = C\left(\sqrt{\frac{k}{m}} + \frac{z}{\sqrt{m}}\right)$ to obtain

$$\sigma_{min}(\hat{\mathbf{B}}^{\top} \mathbf{\Pi}^i \hat{\mathbf{B}}) \ge 1 - \delta_k \tag{30}$$

with probability at least $1 - e^{-(\delta_k \sqrt{m}/C - \sqrt{k})^2}$ for m > k. Now, choose z such that $\delta_k = \frac{12Ck^{3/2}\sqrt{\log(rn)}}{\sqrt{m}}$, we have that (30) holds with probability at least

$$1 - \exp\left(-\left(12k^{3/2}\sqrt{\log(rn)} - \sqrt{k}\right)^2\right) \ge 1 - \exp\left(-k(12\sqrt{k}\sqrt{\log(rn)} - 1)^2\right)$$
$$\ge 1 - \exp\left(121k^3\log(rn)\right) \tag{31}$$

Finally, taking a union bound over $i \in [n]$ yields $\sigma_{\min}(\mathbf{G}) \ge 1 - \delta_k$ with probability at least

$$1 - rn \exp\left(-121k^3 \log(rn)\right) \ge 1 - e^{-110k^3 \log(rn)},\tag{32}$$

completing the proof. \Box

Lemma 3. Let $\delta_k = c \frac{k^{3/2} \sqrt{\log(rn)}}{\sqrt{m}}$ for some absolute constant c, then

$$\|(\mathbf{G}\mathbf{D} - \mathbf{C})\mathbf{w}^*\|_2 \le \delta_k \|\mathbf{W}^*\|_2 \operatorname{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*)$$

with probability at least $1 - e^{-111k^2 \log(rn)}$.

Proof. For ease of notation we drop superscripts t. We define $\mathbf{H} = \mathbf{GD} - \mathbf{C}$ and

$$\mathbf{H}^{i} := \mathbf{G}^{i} \mathbf{D}^{i} - \mathbf{C}^{i} = \hat{\mathbf{B}}^{\top} \mathbf{\Pi}^{i} \hat{\mathbf{B}} \hat{\mathbf{B}}^{\top} \hat{\mathbf{B}}^{*} - \hat{\mathbf{B}}^{\top} \mathbf{\Pi}^{i} \hat{\mathbf{B}}^{*} = \hat{\mathbf{B}}^{\top} \left(\frac{1}{m} \mathbf{X}_{i}^{\top} \mathbf{X}_{i} \right) (\hat{\mathbf{B}} \hat{\mathbf{B}}^{\top} - \mathbf{I}_{d}) \hat{\mathbf{B}}^{*}, \tag{33}$$

for all $i \in [rn]$. Then we have

$$\|(\mathbf{GD} - \mathbf{C})\mathbf{w}_{*}\|_{2}^{2} = \sum_{i=1}^{rn} \|\mathbf{H}^{i}\mathbf{w}_{*}^{i}\|_{2}^{2}$$

$$\leq \sum_{i=1}^{rn} \|\mathbf{H}^{i}\|_{2}^{2} \|\mathbf{w}_{i}^{*}\|_{2}^{2}$$

$$\leq \frac{k}{rn} \|\mathbf{W}^{*}\|_{2}^{2} \sum_{i=1}^{rn} \|\mathbf{H}^{i}\|_{2}^{2}$$
(34)

where the last inequality follows almost surely from Assumption 3 (the 1-row-wise incoherence of \mathbf{W}^*), the fact that $krn = \|\mathbf{W}^*\|_F^2 \le k\|\mathbf{W}^*\|_2^2$ by Assumption 3, and the fact that \mathbf{W}^* has rank k. It remains to bound $\frac{1}{rn} \sum_{i=1}^{rn} \|\mathbf{H}^i\|_2^2$. Although $\|\mathbf{H}^i\|_2$ is sub-exponential (as we will show), $\|\mathbf{H}^i\|_2^2$ is not sub-exponential, so we cannot directly apply standard concentration results. Instead, we compute a tail bound for each $\|\mathbf{H}^i\|_2^2$ individually, then then union bound over $i \in [rn]$. Let $\mathbf{U} := \frac{1}{\sqrt{m}} \mathbf{X}_i(\hat{\mathbf{B}}\hat{\mathbf{B}}^\top - \mathbf{I}_d)\hat{\mathbf{B}}^*$, then the j-th row of \mathbf{U} is given by

$$\mathbf{u}_j = \frac{1}{\sqrt{m}} \mathbf{\hat{B}}^{*^\top} (\mathbf{\hat{B}} \mathbf{\hat{B}}^\top - \mathbf{I}_d) \mathbf{x}_i^j,$$

and is $\frac{1}{\sqrt{m}}\hat{\mathbf{B}}^{*\top}(\hat{\mathbf{B}}\hat{\mathbf{B}}^{\top}-\mathbf{I}_d)$ -sub-gaussian. Likewise, define $\mathbf{V}\coloneqq\frac{1}{\sqrt{m}}\mathbf{X}_i\hat{\mathbf{B}}$, then the j-th row of \mathbf{V} is

$$\mathbf{v}_j = \frac{1}{\sqrt{m}} \hat{\mathbf{B}}^\top \mathbf{x}_i^j,$$

therefore is $\frac{1}{\sqrt{m}}\hat{\mathbf{B}}$ -sub-gaussian. We leverage the sub-gaussianity of the rows of \mathbf{U} and \mathbf{V} to make a similar concentration argument as in Proposition 4.4.5 in Vershynin [2018]. First, let \mathcal{S}^{k-1} denote the unit sphere in k dimensions, and let \mathcal{N}_k be a $\frac{1}{4}$ -th net of cardinality $|\mathcal{N}_k| \leq 9^k$, which exists by Corollary 4.2.13 in Vershynin [2018]. Next, using equation 4.13 in Vershynin [2018], we obtain

$$\begin{split} \|(\hat{\mathbf{B}}^*)^\top (\hat{\mathbf{B}}\hat{\mathbf{B}}^\top - \mathbf{I}_d) \mathbf{X}_i^\top \mathbf{X}_i \mathbf{B}\|_2 &= \left\| \mathbf{U}^\top \mathbf{V} \right\|_2 \le 2 \max_{\mathbf{z}, \mathbf{y} \in \mathcal{N}_k} \mathbf{z}^\top \left(\mathbf{U}^\top \mathbf{V} \right) \mathbf{y} \\ &= 2 \max_{\mathbf{z}, \mathbf{y} \in \mathcal{N}_k} \mathbf{z}^\top \left(\sum_{j=1}^m \mathbf{u}_j \mathbf{v}_j^\top \right) \mathbf{y} \\ &= 2 \max_{\mathbf{z}, \mathbf{y} \in \mathcal{N}_k} \sum_{j=1}^m \langle \mathbf{z}, \mathbf{u}_j \rangle \langle \mathbf{v}_j, \mathbf{y} \rangle \end{split}$$

By definition of sub-gaussianity, $\langle \mathbf{z}, \mathbf{u}_j \rangle$ and $\langle \mathbf{v}_j, \mathbf{y} \rangle$ are sub-gaussian with norms $\frac{1}{\sqrt{m}} \| \hat{\mathbf{B}}^{*\top} (\hat{\mathbf{B}} \hat{\mathbf{B}}^{\top} - \mathbf{I}_d) \|_2 = \frac{1}{\sqrt{m}} \mathrm{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)$ and $\frac{1}{\sqrt{m}} \| \hat{\mathbf{B}} \|_2 = \frac{1}{\sqrt{m}}$, respectively. Thus for all $j \in [m]$, $\langle \mathbf{z}, \mathbf{u}_j \rangle \langle \mathbf{v}_j, \mathbf{z} \rangle$ is sub-exponential with norm $\frac{c}{m} \mathrm{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)$ for some absolute constant c. Note that for any $j \in [m]$ and any \mathbf{z} , $\mathbb{E}[\langle \mathbf{z}, \mathbf{u}_j \rangle \langle \mathbf{v}_j, \mathbf{y} \rangle] = \mathbf{z}^{\top}((\hat{\mathbf{B}}^*)^{\top}(\hat{\mathbf{B}}\hat{\mathbf{B}}^{\top} - \mathbf{I}_d)\mathbf{B})\mathbf{y} = 0$. Thus we have a sum of m mean-zero, independent sub-exponential random variables. We can now use Bernstein's inequality to obtain, for any fixed $\mathbf{z}, \mathbf{y} \in \mathcal{N}_k$,

$$\mathbb{P}\left(\sum_{j=1}^{m} \langle \mathbf{z}, \mathbf{u}_{j} \rangle \langle \mathbf{v}_{j}, \mathbf{y} \rangle \ge s\right) \le \exp\left(-c' m \min\left(\frac{s^{2}}{\operatorname{dist}^{2}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^{*})}, \frac{s}{\operatorname{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^{*})}\right)\right)$$
(35)

Now union bound over all $\mathbf{z}, \mathbf{y} \in \mathcal{N}_k$ to obtain

$$\mathbb{P}\left(\frac{1}{m}\|(\hat{\mathbf{B}}^*)^{\top}(\hat{\mathbf{B}}\hat{\mathbf{B}}^{\top} - \mathbf{I}_d)\mathbf{X}_i^{\top}\mathbf{X}_i\hat{\mathbf{B}}\|_2 \ge 2s\right) \le 9^{2k} \exp\left(-c'm\min(s^2/\mathrm{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*), s/\mathrm{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*))\right)$$
(36)

Let $\frac{s}{\operatorname{dist}(\hat{\mathbf{B}},\hat{\mathbf{B}}^*)} = \max(\varepsilon, \varepsilon^2)$ for some $\epsilon > 0$, then it follows that $\min(s^2/\operatorname{dist}^2(\hat{\mathbf{B}},\hat{\mathbf{B}}^*), s/\operatorname{dist}(\hat{\mathbf{B}},\hat{\mathbf{B}}^*)) = \varepsilon^2$. So we have

$$\mathbb{P}\left(\frac{1}{m}\|(\hat{\mathbf{B}}^*)^{\top}(\hat{\mathbf{B}}\hat{\mathbf{B}}^{\top} - \mathbf{I}_d)\mathbf{X}_i^{\top}\mathbf{X}_i\hat{\mathbf{B}}\|_2 \ge 2\mathrm{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)\max(\varepsilon, \varepsilon^2)\right) \le 9^{2k}e^{-c'm\varepsilon^2}$$
(37)

Moreover, letting $\varepsilon^2 = \frac{ck^2 \log(rn)}{4m}$ for some constant c, and $m \ge ck^2 \log(rn)$, we have

$$\mathbb{P}\left(\frac{1}{m}\|(\hat{\mathbf{B}}^*)^{\top}(\hat{\mathbf{B}}\hat{\mathbf{B}}^{\top} - \mathbf{I}_d)\mathbf{X}_i^{\top}\mathbf{X}_i\hat{\mathbf{B}}\|_2 \ge \operatorname{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)\sqrt{\frac{ck^2\log(rn)}{m}}\right) \le 9^{2k}e^{-c_1k^2\log(rn)} \\
\le e^{-111k^2\log(rn)} \tag{38}$$

for large enough constant c_1 . Thus, noting that $\|\mathbf{H}^i\|_2^2 = \|\frac{1}{m}(\hat{\mathbf{B}}^*)^\top (\hat{\mathbf{B}}\hat{\mathbf{B}}^\top - \mathbf{I}_d)\mathbf{X}_i^\top \mathbf{X}_i \hat{\mathbf{B}}\|_2^2$, we obtain

$$\mathbb{P}\left(\|\mathbf{H}^i\|_2^2 \ge c \operatorname{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*) \frac{k^2 \log(rn)}{m}\right) \le e^{-111k^2 \log(rn)}$$
(39)

Thus, using (34), we have

$$\mathbb{P}\left(\|(\mathbf{G}\mathbf{D} - \mathbf{C})\mathbf{w}_*\|_2^2 \ge c\|\mathbf{W}^*\|_2^2 \operatorname{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*) \frac{k^3 \log(rn)}{m}\right) \\
\le \mathbb{P}\left(\frac{k}{rn}\|\mathbf{W}^*\|_2^2 \sum_{i=1}^{rn} \|\mathbf{H}^i\|_2^2 \ge c\|\mathbf{W}^*\|_2^2 \operatorname{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*) \frac{k^3 \log(rn)}{m}\right) \\
= \mathbb{P}\left(\frac{1}{rn} \sum_{i=1}^{rn} \|\mathbf{H}^i\|_2^2 \ge c \operatorname{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*) \frac{k^2 \log(rn)}{m}\right) \\
\le rn\mathbb{P}\left(\|\mathbf{H}^1\|_2^2 \ge c \operatorname{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*) \frac{k^2 \log(rn)}{m}\right) \\
\le e^{-110k^2 \log(rn)}$$

completing the proof.

Lemma 4. Let $\delta_k = \frac{ck^{3/2}\sqrt{\log(rn)}}{\sqrt{m}}$, then

$$\|\mathbf{F}\|_{F} \le \frac{\delta_{k}}{1 - \delta_{k}} \|\mathbf{W}^{*}\|_{2} \operatorname{dist}(\hat{\mathbf{B}}_{t}, \hat{\mathbf{B}}_{*})$$
(40)

with probability at least $1 - e^{-110k^2 \log(n)}$.

Proof. By the definition of \mathbf{F} and the Cauchy-Schwarz inequality, we have $\|\mathbf{F}\|_F = \|\mathbf{G}^{-1}(\mathbf{G}\mathbf{D} - \mathbf{C})\widetilde{\mathbf{w}}^*\|_2 \le \|\mathbf{G}^{-1}\|_2 \|(\mathbf{G}\mathbf{D} - \mathbf{C})\widetilde{\mathbf{w}}^*\|_2$. Combining the bound on $\|\mathbf{G}^{-1}\|_2$ from Lemma 2 and the bound on $\|(\mathbf{G}\mathbf{D} - \mathbf{C})\widetilde{\mathbf{w}}^*\|_2$ from Lemma 3 via a union bound yields the result.

We next focus on showing concentration of the operator $\frac{1}{m}\mathcal{A}^{\dagger}\mathcal{A}$ to the identity operator.

Lemma 5. Let $\delta'_k = ck \frac{\sqrt{d}}{\sqrt{rnm}}$ for some absolute constant c. Then for any t, if $\delta'_k \leq k$,

$$\frac{1}{rn} \left\| \left(\frac{1}{m} \mathcal{A}^* \mathcal{A}(\mathbf{Q}^t) - \mathbf{Q}^t \right)^\top \mathbf{W}^{t+1} \right\|_2 \le \delta_k' \operatorname{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*)$$
(41)

with probability at least $1 - e^{-110d} - e^{-110k^2 \log(rn)}$.

Proof. We drop superscripts t for simplicity. We first bound the norms of the rows of \mathbf{Q} and \mathbf{W} . Let $\mathbf{q}_i \in \mathbb{R}^d$ be the i-th row of \mathbf{Q} and let $\mathbf{w}_i \in \mathbb{R}^k$ be the i-th row of \mathbf{W} . Recall the computation of \mathbf{W} from Lemma 1:

$$\mathbf{W} = \mathbf{W}_* \hat{\mathbf{B}}_*^\top \hat{\mathbf{B}} - \mathbf{F} \implies \mathbf{w}_i^\top = (\hat{\mathbf{w}}_i^*)^\top \hat{\mathbf{B}}_*^\top \hat{\mathbf{B}} - \mathbf{f}_i^\top$$

Thus

$$\|\mathbf{q}_{i}\|_{2}^{2} = \|\hat{\mathbf{B}}\hat{\mathbf{B}}^{\top}\hat{\mathbf{B}}^{*}\hat{\mathbf{w}}_{i}^{*} - \hat{\mathbf{B}}\mathbf{f}_{i} - \hat{\mathbf{B}}^{*}\hat{\mathbf{w}}_{i}^{*}\|_{2}^{2}$$

$$= \|(\hat{\mathbf{B}}\hat{\mathbf{B}}^{\top} - \mathbf{I}_{d})\hat{\mathbf{B}}^{*}\hat{\mathbf{w}}_{i}^{*} - \hat{\mathbf{B}}\mathbf{f}_{i}\|_{2}^{2}$$

$$\leq 2\|(\hat{\mathbf{B}}\hat{\mathbf{B}}^{\top} - \mathbf{I}_{d})\hat{\mathbf{B}}^{*}\hat{\mathbf{w}}_{i}^{*}\|_{2}^{2} + 2\|\hat{\mathbf{B}}\mathbf{f}_{i}\|_{2}^{2}$$

$$\leq 2\|(\hat{\mathbf{B}}\hat{\mathbf{B}}^{\top} - \mathbf{I}_{d})\hat{\mathbf{B}}^{*}\|_{2}^{2}\|\hat{\mathbf{w}}_{i}^{*}\|_{2}^{2} + 2\|\mathbf{f}_{i}\|_{2}^{2}$$

$$= 2k\mathrm{dist}^{2}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^{*}) + 2\|\mathbf{f}_{i}\|_{2}^{2}$$

$$(42)$$

Also recall that $\text{vec}(\mathbf{F}) = \mathbf{G}^{-1}(\mathbf{GD} - \mathbf{C})\hat{\mathbf{w}}_*$ from Lemma 1. From equation (27), the *i*-th row of \mathbf{F} is given by:

$$\mathbf{f}_i = (\mathbf{G}^i)^{-1} (\mathbf{G}^i \mathbf{D}^i - \mathbf{C}^i) \mathbf{w}_i^*$$

Thus, using the Cauchy-Schwarz inequality and our previous bounds,

$$\|\mathbf{f}_{i}\|_{2}^{2} \leq \|(\mathbf{G}^{i})^{-1}\|_{2}^{2} \|\mathbf{G}^{i}\mathbf{D}^{i} - \mathbf{C}^{i}\|_{2}^{2} \|\mathbf{w}_{i}^{*}\|_{2}^{2}$$

$$\leq \|(\mathbf{G}^{i})^{-1}\|_{2}^{2} \|\mathbf{G}^{i}\mathbf{D}^{i} - \mathbf{C}^{i}\|_{2}^{2} k$$
(43)

where (43) follows by Assumption 3, i.e. the row-wise incoherence of \mathbf{W}^* . From (39), we have that

$$\mathbb{P}\left(\|\mathbf{G}^{i}\mathbf{D}^{i} - \mathbf{C}^{i}\|_{2}^{2} \ge \delta_{k}^{2} \operatorname{dist}^{2}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^{*})\right) \le e^{-112k^{2}\log(rn)}$$

where δ_k is defined in Lemma 2. Similarly, from equations (30) and (31), we have that

$$\mathbb{P}\left(\|(\mathbf{G}^i)^{-1}\|_2^2 \ge \frac{1}{(1-\delta_k)^2}\right) \le e^{-121k^3\log(rn)} \tag{44}$$

Now plugging this back into (43) and assuming $\delta_k \leq \frac{1}{2}$, we obtain

$$\|\mathbf{q}_i\|_2^2 \le 2k \operatorname{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*) \left(1 + \frac{\delta_k^2}{(1 - \delta_k)^2}\right) \le 4k \operatorname{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)$$
(45)

with probability at least $1 - e^{-111k^2 \log(rn)}$. Likewise, to upper bound $\|\mathbf{w}_i\|_2$ we have

$$\|\mathbf{w}_{i}\|_{2}^{2} \leq 2\|\hat{\mathbf{B}}^{\top}\hat{\mathbf{B}}^{*}\mathbf{w}_{i}^{*}\|_{2}^{2} + 2\|\mathbf{f}_{i}\|_{2}^{2}$$

$$\leq 2\|\hat{\mathbf{B}}^{\top}\hat{\mathbf{B}}^{*}\|_{2}^{2}\|\mathbf{w}_{i}^{*}\|_{2}^{2} + 2\|\mathbf{f}_{i}\|_{2}^{2}$$

$$\leq 2k + 2\frac{\delta_{k}^{2}}{(1 - \delta_{k})^{2}}\operatorname{dist}^{2}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^{*})k$$

$$\leq 4k$$
(46)

where (46) holds with probability at least $1 - e^{-111k^2 \log(rn)}$ conditioning on the same event as in (45), and (47) holds almost surely as long as $\delta_k \leq 1/2$. For the rest of the proof we condition on the event $\mathcal{E} := \bigcap_{i=1}^{rn} \left\{ \|\mathbf{q}_i\|_2^2 \leq 4k \operatorname{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*) \cap \|\mathbf{w}_i\|_2^2 \leq 4k \right\}$, which holds with probability at least $1 - e^{-110k^2 \log(rn)}$ by a union bound over $i \in [rn]$. Observe that the matrix $\frac{1}{m} \mathcal{A}^* \mathcal{A}(\mathbf{Q}) - \mathbf{Q}$ can be re-written as

$$\frac{1}{m} \mathcal{A}^* \mathcal{A}(\mathbf{Q}) - \mathbf{Q} = \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^{m} \left(\langle \mathbf{e}_i(\mathbf{x}_i^j)^\top, \mathbf{Q} \rangle \mathbf{e}_i(\mathbf{x}_i^j)^\top - \mathbf{Q} \right)$$

$$= \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^{m} \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \mathbf{e}_i(\mathbf{x}_i^j)^\top - \mathbf{Q} \tag{48}$$

Multiplying the transpose by $\frac{1}{rn}\mathbf{W}$ yields

$$\frac{1}{rn} \left(\frac{1}{m} \mathcal{A}^* \mathcal{A}(\mathbf{Q}) - \mathbf{Q} \right)^\top \mathbf{W} = \frac{1}{rnm} \sum_{i=1}^n \sum_{j=1}^m \left(\langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \ \mathbf{x}_i^j(\mathbf{w}_i)^\top - \mathbf{q}_i(\mathbf{w}_i)^\top \right)$$
(49)

where we have used the fact that $(\mathbf{Q})^{\top}\mathbf{W} = \sum_{i=1}^{n} \mathbf{q}_{i}(\mathbf{w}_{i})^{\top}$. We will argue similarly as in Proposition 4.4.5 in Vershynin [2018] to bound the spectral norm of the *d*-by-*k* matrix in the RHS of (49).

First, let \mathcal{S}^{d-1} and \mathcal{S}^{k-1} denote the unit spheres in d and k dimensions, respectively. Construct $\frac{1}{4}$ -nets \mathcal{N}_d and \mathcal{N}_k over \mathcal{S}^{d-1} and \mathcal{S}^{k-1} , respectively, such that $|\mathcal{N}_d| \leq 9^d$ and $|\mathcal{N}_k| \leq 9^k$ (which is possible by Corollary 4.2.13 in Vershynin [2018]). Then, using equation 4.13 in Vershynin [2018], we

have

$$\left\| \frac{1}{rnm} \sum_{i=1}^{rn} \sum_{j=1}^{m} \left(\langle \mathbf{x}_{i}^{j}, \mathbf{q}_{i} \rangle \mathbf{x}_{i}^{j} (\mathbf{w}_{i})^{\top} - \mathbf{q}_{i} (\mathbf{w}_{i})^{\top} \right) \right\|_{2}^{2}$$

$$\leq 2 \max_{\mathbf{u} \in \mathcal{N}_{d}, \mathbf{v} \in \mathcal{N}_{k}} \mathbf{u}^{\top} \left(\sum_{i=1}^{rn} \sum_{j=1}^{m} \left(\frac{1}{rnm} \langle \mathbf{x}_{i}^{j}, \mathbf{q}_{i} \rangle \mathbf{x}_{i}^{j} (\mathbf{w}_{i})^{\top} - \frac{1}{rnm} \mathbf{q}_{i} (\mathbf{w}_{i})^{\top} \right) \right) \mathbf{v}$$

$$= 2 \max_{\mathbf{u} \in \mathcal{N}_{d}, \mathbf{v} \in \mathcal{N}_{k}} \sum_{i=1}^{rn} \sum_{j=1}^{m} \left(\frac{1}{rnm} \langle \mathbf{x}_{i}^{j}, \mathbf{q}_{i} \rangle \langle \mathbf{u}, \mathbf{x}_{i}^{j} \rangle \langle \mathbf{w}_{i}, \mathbf{v} \rangle - \frac{1}{rnm} \langle \mathbf{u}, \mathbf{q}_{i} \rangle \langle \mathbf{w}_{i}, \mathbf{v} \rangle \right)$$

$$(50)$$

By the \mathbf{I}_d -sub-gaussianity of \mathbf{x}_i^j , the inner product $\langle \mathbf{u}, \mathbf{x}_i^j \rangle$ is sub-gaussian with norm at most $c \|\mathbf{u}\|_2 = c$ for some absolute constant c for any fixed $\mathbf{u} \in \mathcal{N}_d$. Similarly, $\langle \mathbf{x}_i^j, \mathbf{q}_i \rangle$ is sub-gaussian with norm at most $\|\mathbf{q}_i\|_2 \leq 2c\sqrt{k}$ dist $(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)$ using (45). Further, since the sub-exponential norm of the product of two sub-gaussian random variables is at most the product of the sub-gaussian norms of the two random variables (Lemma 2.7.7 in Vershynin [2018]), we have that $\langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \langle \mathbf{u}, \mathbf{x}_i^j \rangle$ is sub-exponential with norm at most $2c^2\sqrt{k}$ dist $(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)$. Further, $\frac{1}{rnm}\langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \langle \mathbf{u}, \mathbf{x}_i^j \rangle \langle \mathbf{w}_i, \mathbf{v} \rangle$ is sub-exponential with norm at most

$$\frac{2c^2\sqrt{k}}{rnm}\operatorname{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)\langle \mathbf{w}_i, \mathbf{v}\rangle \leq \frac{2c^2\sqrt{k}}{rnm}\operatorname{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)\|\mathbf{w}_i\|_2 \leq \frac{c_1k}{rnm}\operatorname{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*).$$

Finally, note that $\mathbb{E}[\frac{1}{rnm}\langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \langle \mathbf{u}, \mathbf{x}_i^j \rangle \langle \mathbf{w}_i, \mathbf{v} \rangle - \frac{1}{rnm}\langle \mathbf{u}, \mathbf{q}_i \rangle \langle \mathbf{w}_i, \mathbf{v} \rangle] = 0$. Thus, we have a sum of rnm independent, mean zero sub-exponential random variables, so we apply Bernstein's inequality.

$$\mathbb{P}\left(\sum_{i=1}^{rn}\sum_{j=1}^{m}\left(\frac{1}{rnm}\langle\mathbf{x}_{i}^{j},\mathbf{q}_{i}\rangle\langle\mathbf{u},\mathbf{x}_{i}^{j}\rangle\langle\mathbf{w}_{i},\mathbf{v}\rangle - \frac{1}{rnm}\langle\mathbf{u},\mathbf{q}_{i}\rangle\langle\mathbf{w}_{i},\mathbf{v}\rangle\right) \geq s\right)$$

$$\leq \exp\left(-c_{1}rnm\min\left(\frac{s^{2}}{k^{2}\operatorname{dist}^{2}(\hat{\mathbf{B}},\hat{\mathbf{B}}^{*})}, \frac{s}{k\operatorname{dist}(\hat{\mathbf{B}},\hat{\mathbf{B}}^{*})}\right)\right)$$

Union bounding over all $\mathbf{u} \in \mathcal{N}_d$ and $\mathbf{v} \in \mathcal{N}_k$, we obtain

$$\mathbb{P}\left(\left\|\frac{1}{rn}\left(\frac{1}{m}\mathcal{A}^*\mathcal{A}(\mathbf{Q}) - \mathbf{Q}\right)^{\top}\mathbf{W}\right\|_{2} \ge 2s \mid \mathcal{E}\right) \le 9^{d+k} \exp\left(-c_1 rnm \min\left(\frac{s^2}{k^2 \operatorname{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)}, \frac{s}{k \operatorname{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)}\right)\right)$$

Let $\frac{s}{k \operatorname{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)} = \max(\epsilon, \epsilon^2)$ for some $\epsilon > 0$, then $\epsilon^2 = \min\left(\frac{s^2}{k^2 \operatorname{dist}^2(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)}, \frac{s}{k \operatorname{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)}\right)$. Further, let $\epsilon^2 = \frac{112(d+k)}{c_1 r n m}$, then as long as $\epsilon^2 \leq 1$, we have

$$\mathbb{P}\left(\left\|\frac{1}{rn}\left(\frac{1}{m}\mathcal{A}^*\mathcal{A}(\mathbf{Q}) - \mathbf{Q}\right)^{\top}\mathbf{W}\right\|_{2} \ge c_2k \operatorname{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*)\sqrt{d/(rnm)} \mid \mathcal{E}^c\right) \le e^{-110(d+k)} \le e^{-110d}.$$

Finally, we use
$$\mathbb{P}(A \mid \mathcal{E}^c) \leq \mathbb{P}(A \mid \mathcal{E}^c) + \mathbb{P}(\mathcal{E}^c)$$
, where $A := \left\{ \left\| \frac{1}{rn} \left(\frac{1}{m} \mathcal{A}^* \mathcal{A}(\mathbf{Q}) - \mathbf{Q} \right)^\top \mathbf{W} \right\|_2 \geq c_2 k \operatorname{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*) \sqrt{d/(rnm)} \right\}$, to complete the proof.

B.3 Main Result

Now we are ready to show Theorem 1, which follows immediately from the following descent lemma.

Lemma 6. Define $E_0 := 1 - \text{dist}^2(\hat{\mathbf{B}}^0, \hat{\mathbf{B}}^*)$ and $\bar{\sigma}_{\max,*} := \max_{\mathcal{I} \in [n], |\mathcal{I}| = rn} \sigma_{\max}(\frac{1}{\sqrt{rn}} \mathbf{W}_{\mathcal{I}}^*)$ and $\bar{\sigma}_{\min,*} := \min_{\mathcal{I} \in [n], |\mathcal{I}| = rn} \sigma_{\min}(\frac{1}{\sqrt{rn}} \mathbf{W}_{\mathcal{I}}^*)$, i.e. the maximum and minimum singular values of any matrix that can be obtained by taking rn rows of $\frac{1}{\sqrt{rn}} \mathbf{W}^*$.

Suppose that $m \ge c(\kappa^4 k^3 \log(rn)/E_0^2 + \kappa^4 k^2 d/(E_0^2 rn))$ for some absolute constant c. Then for any t and any $\eta \le 1/(4\bar{\sigma}_{\max,*}^2)$, we have

$$\operatorname{dist}(\hat{\mathbf{B}}^{t+1}, \hat{\mathbf{B}}^*) \le (1 - \eta E_0 \bar{\sigma}_{\min,*}^2 / 2)^{1/2} \operatorname{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*),$$

with probability at least $1 - e^{-100 \min(k^2 \log(rn), d)}$.

Proof. Recall that $\mathbf{W}^{t+1} \in \mathbb{R}^{rn \times k}$ and $\mathbf{\bar{B}}^{t+1} \in \mathbb{R}^{d \times k}$ are computed as follows:

$$\mathbf{W}^{t+1} = \underset{\mathbf{W} \in \mathbb{R}^{rn \times k}}{\operatorname{argmin}} \frac{1}{2rnm} \| \mathcal{A}(\mathbf{W}^* \hat{\mathbf{B}}^{*^{\top}} - \mathbf{W} \hat{\mathbf{B}}^{t^{\top}}) \|_2^2$$
 (51)

$$\bar{\mathbf{B}}^{t+1} = \hat{\mathbf{B}}^t - \frac{\eta}{rnm} \left(\mathcal{A}^{\dagger} \mathcal{A} (\mathbf{W}^{t+1} \hat{\mathbf{B}}^{t^{\top}} - \mathbf{W}^* \hat{\mathbf{B}}^{*^{\top}}) \right)^{\top} \mathbf{W}^{t+1}$$
(52)

Let $\mathbf{Q}^t = \mathbf{W}^{t+1} \hat{\mathbf{B}}^{t^{\top}} - \mathbf{W}^* \hat{\mathbf{B}}^{*^{\top}}$. We have

$$\bar{\mathbf{B}}^{t+1} = \hat{\mathbf{B}}^t - \frac{\eta}{rnm} \left(\mathcal{A}^{\dagger} \mathcal{A}(\mathbf{Q}^t) \right)^{\top} \mathbf{W}^{t+1}
= \hat{\mathbf{B}}^t - \frac{\eta}{rn} \mathbf{Q}^{t^{\top}} \mathbf{W}^{t+1} - \frac{\eta}{rn} \left(\frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A}(\mathbf{Q}^t) - \mathbf{Q}^t \right)^{\top} \mathbf{W}^{t+1}$$
(53)

Now, multiply both sides by $\mathbf{\hat{B}}_{\perp}^{*^{\top}}$ to obtain

$$\hat{\mathbf{B}}_{\perp}^{*\top} \bar{\mathbf{B}}^{t+1} = \hat{\mathbf{B}}_{\perp}^{*\top} \hat{\mathbf{B}}^{t} - \frac{\eta}{rn} \hat{\mathbf{B}}_{\perp}^{*\top} \mathbf{Q}^{t^{\top}} \mathbf{W}^{t+1} - \frac{\eta}{rn} \hat{\mathbf{B}}_{\perp}^{*\top} \left(\frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A}(\mathbf{Q}^{t}) - \mathbf{Q}^{t} \right)^{\top} \mathbf{W}^{t+1}$$

$$= \hat{\mathbf{B}}_{\perp}^{*\top} \hat{\mathbf{B}}^{t} (\mathbf{I}_{k} - \frac{\eta}{rn} \mathbf{W}^{t+1^{\top}} \mathbf{W}^{t+1}) - \frac{\eta}{rn} \hat{\mathbf{B}}_{\perp}^{*\top} \left(\frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A}(\mathbf{Q}^{t}) - \mathbf{Q}^{t} \right)^{\top} \mathbf{W}^{t+1}$$
(54)

where the second equality follows because $\hat{\mathbf{B}}_{\perp}^{*^{\mathsf{T}}} \mathbf{Q}^{t^{\mathsf{T}}} = \hat{\mathbf{B}}_{\perp}^{*^{\mathsf{T}}} \hat{\mathbf{B}}^{t} \mathbf{W}^{t+1^{\mathsf{T}}} - \hat{\mathbf{B}}_{\perp}^{*^{\mathsf{T}}} \hat{\mathbf{B}}^{*} \mathbf{W}^{*^{\mathsf{T}}} = \hat{\mathbf{B}}_{\perp}^{*^{\mathsf{T}}} \hat{\mathbf{B}}^{t} \mathbf{W}^{t+1^{\mathsf{T}}}$. Then, writing the QR decomposition of $\bar{\mathbf{B}}^{t+1}$ as $\mathbf{B}^{t+1} = \hat{\mathbf{B}}^{t+1} \mathbf{R}^{t+1}$ and multiplying both sides of (54) from the right by $(\mathbf{R}^{t+1})^{-1}$ yields

$$\hat{\mathbf{B}}_{\perp}^{*^{\top}}\hat{\mathbf{B}}^{t+1} = \left(\hat{\mathbf{B}}_{\perp}^{*^{\top}}\hat{\mathbf{B}}^{t}(\mathbf{I}_{k} - \frac{\eta}{rn}(\mathbf{W}^{t+1})^{\top}\mathbf{W}^{t+1}) - \frac{\eta}{rn}\hat{\mathbf{B}}_{\perp}^{*^{\top}}\left(\frac{1}{m}\mathcal{A}^{\dagger}\mathcal{A}(\mathbf{Q}^{t}) - \mathbf{Q}^{t}\right)^{\top}\mathbf{W}^{t+1}\right)(\mathbf{R}^{t+1})^{-1}$$
(55)

Hence,

$$\operatorname{dist}(\hat{\mathbf{B}}^{t+1}, \hat{\mathbf{B}}^{*}) = \left\| \left(\hat{\mathbf{B}}_{\perp}^{*\top} \hat{\mathbf{B}}^{t} (\mathbf{I}_{k} - \frac{\eta}{rn} (\mathbf{W}^{t+1})^{\top} \mathbf{W}^{t+1}) - \frac{\eta}{rn} \hat{\mathbf{B}}_{\perp}^{*\top} \left(\frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A} (\mathbf{Q}^{t}) - \mathbf{Q}^{t} \right)^{\top} \mathbf{W}^{t+1} \right) (\mathbf{R}^{t+1})^{-1} \right\|_{2}$$

$$\leq \left\| \hat{\mathbf{B}}_{\perp}^{*\top} \hat{\mathbf{B}}^{t} (\mathbf{I}_{k} - \frac{\eta}{rn} (\mathbf{W}^{t+1})^{\top} \mathbf{W}^{t+1}) \right\|_{2} \left\| (\mathbf{R}^{t+1})^{-1} \right\|_{2}$$

$$+ \frac{\eta}{rn} \left\| \hat{\mathbf{B}}_{\perp}^{*\top} \left(\frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A} (\mathbf{Q}^{t}) - \mathbf{Q}^{t} \right)^{\top} \mathbf{W}^{t+1} \right\|_{2} \left\| (\mathbf{R}^{t+1})^{-1} \right\|_{2}$$

$$=: A_{1} + A_{2}.$$

$$(56)$$

where (56) follows by applying the triangle and Cauchy-Schwarz inequalities. We have thus split the upper bound on $\operatorname{dist}(\mathbf{B}^{t+1}, \hat{\mathbf{B}}^*)$ into two terms, A_1 and A_2 . The second term, A_2 , is small due to the concentration of $\frac{1}{m}\mathcal{A}^{\dagger}\mathcal{A}$ to the identity operator, and the first term is strictly smaller than $\operatorname{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*)$. We start by controlling A_2 :

$$A_{2} = \frac{\eta}{rn} \left\| \hat{\mathbf{A}}_{\perp}^{*\top} \left(\frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A}(\mathbf{Q}^{t}) - \mathbf{Q}^{t} \right)^{\top} \mathbf{W}^{t+1} \right\|_{2} \left\| (\mathbf{R}^{t+1})^{-1} \right\|_{2}$$

$$\leq \frac{\eta}{rn} \left\| \left(\frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A}(\mathbf{Q}^{t}) - \mathbf{Q}^{t} \right)^{\top} \mathbf{W}^{t+1} \right\|_{2} \left\| (\mathbf{R}^{t+1})^{-1} \right\|_{2}$$

$$\leq \eta \delta_{k}' \operatorname{dist}(\hat{\mathbf{B}}^{t}, \hat{\mathbf{B}}^{*}) \left\| (\mathbf{R}^{t+1})^{-1} \right\|_{2}$$

$$(58)$$

where (58) follows almost surely by Cauchy-Schwarz and the fact that $\hat{\mathbf{B}}_{\perp}^*$ is normalized, and (59) follows with probability at least $1 - e^{-110d}$ by Lemma 5. Next we control A_1 :

$$A_{1} = \left\| \hat{\mathbf{B}}_{\perp}^{*\top} \hat{\mathbf{B}}^{t} (\mathbf{I}_{k} - \frac{\eta}{rn} (\mathbf{W}^{t+1})^{\top} \mathbf{W}^{t+1}) \right\|_{2} \| (\mathbf{R}^{t+1})^{-1} \|_{2}$$

$$\leq \left\| \hat{\mathbf{B}}_{\perp}^{*\top} \hat{\mathbf{B}}^{t} \right\|_{2} \left\| \mathbf{I} - \frac{\eta}{rn} (\mathbf{W}^{t+1})^{\top} \mathbf{W}^{t+1} \right\|_{2} \| (\mathbf{R}^{t+1})^{-1} \|_{2}$$

$$= \operatorname{dist}(\hat{\mathbf{B}}^{t}, \hat{\mathbf{B}}^{*}) \left\| \mathbf{I}_{k} - \frac{\eta}{rn} (\mathbf{W}^{t+1})^{\top} \mathbf{W}^{t+1} \right\|_{2} \| (\mathbf{R}^{t+1})^{-1} \|_{2}$$

$$(60)$$

The middle factor gives us contraction. To see this, recall that $\mathbf{W}^{t+1} = \mathbf{W}^* \hat{\mathbf{B}}^{*^{\top}} \hat{\mathbf{B}}^t - \mathbf{F}$ where \mathbf{F} is defined in Lemma 1. By Lemma 4, we have that

$$\|\mathbf{F}\|_{2} \leq \frac{\delta_{k}}{1 - \delta_{k}} \|\mathbf{W}^{*}\|_{2} \operatorname{dist}(\hat{\mathbf{B}}^{t}, \hat{\mathbf{B}}^{*})$$
(61)

with probability at least $1 - e^{-110k^2 \log(rn)}$, which we will use throughout the proof. Conditioning on this event, we have

$$\lambda_{\max} \left((\mathbf{W}^{t+1})^{\top} \mathbf{W}^{t+1} \right) = \|\mathbf{W}^* \hat{\mathbf{B}}^{*^{\top}} \hat{\mathbf{B}}^t - \mathbf{F}\|_2^2$$

$$\leq 2 \|\mathbf{W}^* \hat{\mathbf{B}}^{*^{\top}} \hat{\mathbf{B}}^t \|_2^2 + 2 \|\mathbf{F}\|_2^2$$

$$\leq 2 \|\mathbf{W}^* \|_2^2 + 2 \frac{\delta_k^2}{(1 - \delta_k)^2} \|\mathbf{W}^* \|_2^2 \operatorname{dist}^2(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*)$$

$$\leq 4 \|\mathbf{W}^* \|_2^2$$
(62)

where (62) follows under the assumption that $\delta_k \leq 1/2$. Thus, as long as $\eta \leq 1/(4\bar{\sigma}_{\max,*}^2)$, we have by Weyl's Inequality:

$$\|\mathbf{I}_{k} - \frac{\eta}{rn} (\mathbf{W}^{t+1})^{\top} \mathbf{W}^{t+1} \|_{2}$$

$$\leq 1 - \frac{\eta}{rn} \lambda_{\min} ((\mathbf{W}^{t+1})^{\top} \mathbf{W}^{t+1})$$

$$= 1 - \frac{\eta}{rn} \lambda_{\min} ((\mathbf{W}^{*} \hat{\mathbf{B}}^{*^{\top}} \hat{\mathbf{B}}^{t} - \mathbf{F})^{\top} (\mathbf{W}^{*} \hat{\mathbf{B}}^{*^{\top}} \hat{\mathbf{B}}^{t} - \mathbf{F}))$$
(63)

$$\leq 1 - \frac{\eta}{rn} \sigma_{\min}^{2}(\mathbf{W}^{*}(\hat{\mathbf{B}}^{*})^{\top}\hat{\mathbf{B}}^{t}) + \frac{2\eta}{rn} \sigma_{\max}(\mathbf{F}^{\top}\mathbf{W}^{*}(\hat{\mathbf{B}}^{*})^{\top}\hat{\mathbf{B}}^{t}) - \frac{\eta}{rn} \sigma_{\min}^{2}(\mathbf{F})$$
(64)

$$\leq 1 - \frac{\eta}{rn} \sigma_{\min}^2(\mathbf{W}^*) \sigma_{\min}^2((\hat{\mathbf{B}}^*)^\top \hat{\mathbf{B}}^t) + \frac{2\eta}{rn} \|\mathbf{F}\|_2 \|\mathbf{W}^*(\hat{\mathbf{B}}^*)^\top \hat{\mathbf{B}}^t\|_2$$
 (65)

$$\leq 1 - \frac{\eta}{rn} \sigma_{\min}^2(\mathbf{W}^*) \sigma_{\min}^2((\hat{\mathbf{B}}^*)^{\top} \hat{\mathbf{B}}^t) + \frac{2\eta}{rn} \frac{\delta_k}{1 - \delta_k} \|\mathbf{W}^*\|_2^2$$
(66)

$$= 1 - \eta \bar{\sigma}_{\min,*}^2 \sigma_{\min}^2((\hat{\mathbf{B}}^*)^\top \hat{\mathbf{B}}^t) + 2\eta \frac{\delta_k}{1 - \delta_k} \bar{\sigma}_{\max,*}^2$$

$$(67)$$

where (64) follows by again applying Weyl's inequality, under the condition that $2\sigma_{\max}(\mathbf{F}^{\top}\mathbf{W}^{*}(\hat{\mathbf{B}}^{*})^{\top}\hat{\mathbf{B}}^{t}) \leq \sigma_{\min}^{2}(\mathbf{W}^{*})\sigma_{\min}^{2}((\hat{\mathbf{B}}^{*})^{\top}\hat{\mathbf{B}}^{t})$, which we will enforce to be true (otherwise we would not have contraction). Also, (65) follows by the Cauchy-Schwarz inequality, and we use Lemma 4 to obtain (66). Lastly, (67) follows by the definitions of $\bar{\sigma}_{\min,*}$ and $\bar{\sigma}_{\max,*}$. In order to lower bound $\sigma_{\min}^{2}((\hat{\mathbf{B}}^{*})^{\top}\hat{\mathbf{B}}^{t})$, note that

$$\sigma_{\min}^2((\hat{\mathbf{B}}^*)^{\top}\hat{\mathbf{B}}^t) \ge 1 - \|(\hat{\mathbf{B}}_{\perp}^*)^{\top}\hat{\mathbf{B}}^t\|_2^2 = 1 - \operatorname{dist}^2(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) \ge 1 - \operatorname{dist}^2(\hat{\mathbf{B}}^0, \hat{\mathbf{B}}^*) =: E_0$$
 (68)

As a result, defining $\bar{\delta}_k := \delta_k + \delta'_k$ and combining (56), (59), (60), (67), and (68) yields

$$\operatorname{dist}(\hat{\mathbf{B}}^{t+1}, \hat{\mathbf{B}}^{*}) \leq \|(\mathbf{R}^{t+1})^{-1}\|_{2} \left(1 - \eta \bar{\sigma}_{\min,*}^{2} E_{0} + 2\eta \frac{\delta_{k}}{1 - \delta_{k}} \bar{\sigma}_{\max,*}^{2} + \eta \delta_{k}'\right) \operatorname{dist}(\hat{\mathbf{B}}^{t}, \hat{\mathbf{B}}^{*})$$

$$\leq \|(\mathbf{R}^{t+1})^{-1}\|_{2} \left(1 - \eta \bar{\sigma}_{\min,*}^{2} E_{0} + 2\eta \frac{\bar{\delta}_{k}}{1 - \bar{\delta}_{k}} \bar{\sigma}_{\max,*}^{2}\right) \operatorname{dist}(\hat{\mathbf{B}}^{t}, \hat{\mathbf{B}}^{*})$$
(69)

where (69) follows from the fact that $krn = \|\mathbf{W}^*\|_F^2 \le k\|\mathbf{W}^*\|_2^2 \implies 1 \le \|\mathbf{W}^*\|_2^2/rn \le \bar{\sigma}_{\max,*}^2$. All that remains to bound is $\|(\mathbf{R}^{t+1})^{-1}\|_2$. Define $\mathbf{S}^t := \frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A}(\mathbf{Q}^t)$ and observe that

$$(\mathbf{R}^{t+1})^{\top} \mathbf{R}^{t+1} = (\bar{\mathbf{B}}^{t+1})^{\top} \bar{\mathbf{B}}^{t+1}$$

$$= \hat{\mathbf{B}}^{t^{\top}} \hat{\mathbf{B}}^{t} - \frac{\eta}{rn} (\hat{\mathbf{B}}^{t^{\top}} \mathbf{S}^{t^{\top}} \mathbf{W}^{t+1} + (\mathbf{W}^{t+1})^{\top} \mathbf{S}^{t} \hat{\mathbf{B}}^{t}) + \frac{\eta^{2}}{(rn)^{2}} (\mathbf{W}^{t+1})^{\top} \mathbf{S}^{t} \mathbf{S}^{t^{\top}} \mathbf{W}^{t+1}$$

$$= \mathbf{I}_{k} - \frac{\eta}{rn} (\hat{\mathbf{B}}^{t^{\top}} \mathbf{S}^{t^{\top}} \mathbf{W}^{t+1} + (\mathbf{W}^{t+1})^{\top} \mathbf{S}^{t} \hat{\mathbf{B}}^{t}) + \frac{\eta^{2}}{(rn)^{2}} (\mathbf{W}^{t+1})^{\top} \mathbf{S}^{t} \mathbf{S}^{t^{\top}} \mathbf{W}^{t+1}$$
(70)

thus, by Weyl's Inequality, we have

$$\sigma_{\min}^{2}(\mathbf{R}_{t+1}) \geq 1 - \frac{\eta}{rn} \lambda_{\max}(\hat{\mathbf{B}}^{t^{\top}} \mathbf{S}^{t^{\top}} \mathbf{W}^{t+1} + (\mathbf{W}^{t+1})^{\top} \mathbf{S}^{t} \hat{\mathbf{B}}^{t}) + \frac{\eta^{2}}{(rn)^{2}} \lambda_{\min}((\mathbf{W}^{t+1})^{\top} \mathbf{S}^{t} \mathbf{S}^{t^{\top}} \mathbf{W}^{t+1})$$

$$\geq 1 - \frac{\eta}{rn} \lambda_{\max}(\hat{\mathbf{B}}^{t^{\top}} \mathbf{S}^{t^{\top}} \mathbf{W}^{t+1} + (\mathbf{W}^{t+1})^{\top} \mathbf{S}^{t} \hat{\mathbf{B}}^{t})$$

$$(71)$$

where (71) follows because $(\mathbf{W}^{t+1})^{\top} \mathbf{S}^t \mathbf{S}^{t^{\top}} \mathbf{W}^{t+1}$ is positive semi-definite. Next, note that

$$\frac{\eta}{rn} \lambda_{\max} (\hat{\mathbf{B}}^{t^{\top}} \mathbf{S}^{t^{\top}} \mathbf{W}^{t+1} + (\mathbf{W}^{t+1})^{\top} \mathbf{S}^{t} \hat{\mathbf{B}}^{t})
= \max_{\mathbf{x}: \|\mathbf{x}\|_{2} = 1} \frac{\eta}{rn} \mathbf{x}^{\top} \hat{\mathbf{B}}^{t^{\top}} (\mathbf{S}^{t})^{\top} \mathbf{W}^{t+1} \mathbf{x} + \mathbf{x}^{\top} (\mathbf{W}^{t+1})^{\top} \mathbf{S}^{t} \hat{\mathbf{B}}^{t} \mathbf{x}
= \max_{\mathbf{x}: \|\mathbf{x}\|_{2} = 1} \frac{2\eta}{rn} \mathbf{x}^{\top} (\mathbf{W}^{t+1})^{\top} \mathbf{S}^{t} \hat{\mathbf{B}}^{t} \mathbf{x}
= \max_{\mathbf{x}: \|\mathbf{x}\|_{2} = 1} \frac{2\eta}{rn} \mathbf{x}^{\top} (\mathbf{W}^{t+1})^{\top} \left(\frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A} (\mathbf{Q}^{t}) - \mathbf{Q}^{t} \right) \hat{\mathbf{B}}^{t} \mathbf{x} + \frac{2\eta}{rn} \mathbf{x}^{\top} (\mathbf{W}^{t+1})^{\top} \mathbf{Q}^{t} \hat{\mathbf{B}}^{t} \mathbf{x}$$
(72)

We first consider the first term. We have

$$\max_{\mathbf{x}:\|\mathbf{x}\|_{2}=1} \frac{2\eta}{rn} \mathbf{x}^{\top} (\mathbf{W}^{t+1})^{\top} \left(\frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A} (\mathbf{Q}^{t}) - \mathbf{Q}^{t} \right) \hat{\mathbf{B}}^{t} \mathbf{x} \leq \frac{2\eta}{rn} \left\| (\mathbf{W}^{t+1})^{\top} \left(\frac{1}{m} \mathcal{A}^{\dagger} \mathcal{A} (\mathbf{Q}^{t}) - \mathbf{Q}^{t} \right) \right\|_{2} \left\| \hat{\mathbf{B}}^{t} \right\|_{2} \leq 2\eta \delta_{k}'$$
(73)

where the last inequality follows with probability at least $1 - e^{-110d} - e^{-110k^2 \log(rn)}$ from Lemma 5. Next we turn to the second term in (72). We have

$$\max_{\mathbf{x}:\|\mathbf{x}\|_{2}=1} \frac{2\eta}{rn} \mathbf{x}^{\top} (\mathbf{W}^{t+1})^{\top} \mathbf{Q}^{t} \hat{\mathbf{B}}^{t} \mathbf{x} = \max_{\mathbf{x}:\|\mathbf{x}\|_{2}=1} \frac{2\eta}{rn} \left\langle \mathbf{Q}^{t}, \mathbf{W}^{t+1} \mathbf{x} \mathbf{x}^{\top} \hat{\mathbf{B}}^{t^{\top}} \right\rangle
= \max_{\mathbf{x}:\|\mathbf{x}\|_{2}=1} \frac{2\eta}{rn} \left\langle \mathbf{Q}^{t}, \mathbf{W}^{*} \hat{\mathbf{B}}^{*^{\top}} \hat{\mathbf{B}}^{t} \mathbf{x} \mathbf{x}^{\top} \hat{\mathbf{B}}^{t^{\top}} \right\rangle - \frac{2\eta}{rn} \left\langle \mathbf{Q}^{t}, \mathbf{F} \mathbf{x} \mathbf{x}^{\top} \hat{\mathbf{B}}^{t^{\top}} \right\rangle$$
(74)

For any $\mathbf{x} \in \mathbb{R}^k : \|\mathbf{x}\|_2 = 1$, we have

$$\frac{2\eta}{rn}\langle\mathbf{Q}^{t},\mathbf{W}^{*}(\hat{\mathbf{B}}^{*})^{\top}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}\rangle
= \frac{2\eta}{rn}\mathrm{tr}((\hat{\mathbf{B}}^{t}(\mathbf{W}^{t+1})^{\top} - \hat{\mathbf{B}}^{*}(\mathbf{W}^{*})^{\top})\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}})
= \frac{2\eta}{rn}\mathrm{tr}((\hat{\mathbf{B}}^{t}\hat{\mathbf{B}}^{t^{\top}}\hat{\mathbf{B}}^{*}\mathbf{W}^{*^{\top}} - \hat{\mathbf{B}}^{t}\mathbf{F}^{\top} - \hat{\mathbf{B}}^{*}\mathbf{W}^{*^{\top}})\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}})
= \frac{2\eta}{rn}\mathrm{tr}((\hat{\mathbf{B}}^{t}\hat{\mathbf{B}}^{t^{\top}} - \mathbf{I})\hat{\mathbf{B}}^{*^{\top}}\mathbf{W}^{*^{\top}}\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}) - \frac{2\eta}{rn}\mathrm{tr}(\hat{\mathbf{B}}^{t}\mathbf{F}^{\top}\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}})
= \frac{2\eta}{rn}\mathrm{tr}(\hat{\mathbf{B}}^{t}_{\perp}\hat{\mathbf{B}}^{*^{\top}}\mathbf{W}^{*^{\top}}\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}) - \frac{2\eta}{rn}\mathrm{tr}(\hat{\mathbf{B}}^{t}\mathbf{F}^{\top}\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}})
= \frac{2\eta}{rn}\mathrm{tr}(\hat{\mathbf{B}}^{*^{\top}}\mathbf{W}^{*^{\top}}\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}\hat{\mathbf{B}}^{t}_{\perp}) - \frac{2\eta}{rn}\mathrm{tr}(\hat{\mathbf{B}}^{t}\mathbf{F}^{\top}\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}})
= -\frac{2\eta}{rn}\mathrm{tr}(\mathbf{F}^{\top}\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}\hat{\mathbf{B}}^{t})
= -\frac{2\eta}{rn}\mathrm{tr}(\mathbf{F}^{\top}\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}\hat{\mathbf{B}}^{t})
= -\frac{2\eta}{rn}\mathrm{tr}(\mathbf{F}^{\top}\mathbf{W}^{*}\hat{\mathbf{B}}^{*^{\top}}\hat{\mathbf{B}}^{t}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}\hat{\mathbf{B}}^{t})$$
(75)

$$= -\frac{2\eta}{rn} \operatorname{tr}(\mathbf{F}^{\top} \mathbf{W}^* \hat{\mathbf{B}}^{*\top} \hat{\mathbf{B}}^t \mathbf{x} \mathbf{x}^{\top})$$
 (76)

$$\leq \frac{2\eta}{rn} \|\mathbf{F}\|_F \|\mathbf{W}^* \hat{\mathbf{B}}^{*\top} \hat{\mathbf{B}}^t \mathbf{x} \mathbf{x}^\top \|_F \tag{77}$$

$$\leq \frac{2\eta}{rn} \|\mathbf{F}\|_F \|\mathbf{W}^*\|_2 \|\hat{\mathbf{B}}^{*^{\top}}\|_2 \|\hat{\mathbf{B}}^t\|_2 \|\mathbf{x}\mathbf{x}^{\top}\|_F \tag{78}$$

$$\leq \frac{2\eta}{rn} \|\mathbf{F}\|_F \|\mathbf{W}^*\|_2 \tag{79}$$

$$\leq 2\eta \frac{\delta_k}{1 - \delta_k} \bar{\sigma}_{\text{max},*}^2 \tag{80}$$

where (75) follows since $\hat{\mathbf{B}}^{t^{\top}}\hat{\mathbf{B}}_{\perp}^{t} = \mathbf{0}$, (76) follows since $\hat{\mathbf{B}}^{t^{\top}}\hat{\mathbf{B}}^{t} = \mathbf{I}_{k}$, (77) and (78) follows by the Cauchy-Schwarz inequality, (79) follows by the orthonormality of $\hat{\mathbf{B}}^{t}$ and $\hat{\mathbf{B}}^{*}$ and (80) follows by Lemma 4 and the definition of $\bar{\sigma}_{\max,*}$. Next, again for any $\mathbf{x} \in \mathbb{R}^{k} : \|\mathbf{x}\|_{2} = 1$,

$$-\frac{2\eta}{rn}\langle \mathbf{Q}^{t}, \mathbf{F}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}\rangle = -\frac{2\eta}{rn}\mathrm{tr}((\hat{\mathbf{B}}^{t}\hat{\mathbf{B}}^{t^{\top}}\hat{\mathbf{B}}^{*}\mathbf{W}^{*^{\top}} - \hat{\mathbf{B}}^{t}\mathbf{F}^{\top} - \hat{\mathbf{B}}^{*}\mathbf{W}^{*^{\top}})\mathbf{F}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}})$$

$$= -\frac{2\eta}{rn}\mathrm{tr}((\hat{\mathbf{B}}^{t}\hat{\mathbf{B}}^{t^{\top}} - \mathbf{I}_{d})\hat{\mathbf{B}}^{*}\mathbf{W}^{*^{\top}}\mathbf{F}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}) + \frac{2\eta}{rn}\mathrm{tr}(\mathbf{F}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}\hat{\mathbf{B}}^{t}\mathbf{F}^{\top})$$

$$= -\frac{2\eta}{rn}\mathrm{tr}(\hat{\mathbf{B}}^{*}\mathbf{W}^{*^{\top}}\mathbf{F}\mathbf{x}\mathbf{x}^{\top}\hat{\mathbf{B}}^{t^{\top}}\mathbf{B}_{\perp}^{t}) + \frac{2\eta}{rn}\mathbf{x}^{\top}\mathbf{F}^{\top}\mathbf{F}\mathbf{x}$$

$$= \frac{2\eta}{rn}\mathbf{x}^{\top}\mathbf{F}^{\top}\mathbf{F}\mathbf{x}$$

$$\leq \frac{2\eta}{rn}\|\mathbf{F}\|_{2}^{2}$$

$$\leq 2\eta \frac{\delta_{k}^{2}}{(1-\delta_{k})^{2}}\bar{\sigma}_{\max,*}^{2}$$
(81)

Thus, we have the following bound on the second term of (72):

$$\max_{\mathbf{x}:\|\mathbf{x}\|_2=1} \frac{2\eta}{rn} \langle \mathbf{Q}^t, \mathbf{W}^{t+1} \mathbf{x} \mathbf{x}^\top \hat{\mathbf{B}}^{t^\top} \rangle \le 2\eta \bar{\sigma}_{\max,*}^2 \left(\frac{\delta_k}{1-\delta_k} + \frac{\delta_k^2}{(1-\delta_k)^2} \right) \le 4\eta \frac{\delta_k}{(1-\delta_k)^2} \bar{\sigma}_{\max,*}^2$$
(82)

since $0 \le \delta_k \le 1 \implies \delta_k^2 \le \delta_k$. Therefore, using (71), (72), (73) and (82), we have

$$\sigma_{\min}^{2}(\mathbf{R}_{t+1}) \ge 1 - 2\eta \delta_{k}' - 4\eta \frac{\delta_{k}}{(1 - \delta_{k})^{2}} \bar{\sigma}_{\max,*}^{2} \ge 1 - 4\eta \frac{\bar{\delta}_{k}}{(1 - \bar{\delta}_{k})^{2}} \bar{\sigma}_{\max,*}^{2}$$
(83)

where $\bar{\delta}_k = \delta_k' + \delta_k$. This means that

$$\|(\mathbf{R}^{t+1})^{-1}\|_{2} \le \left(1 - 4\eta \frac{\bar{\delta}_{k}}{(1 - \bar{\delta}_{k})^{2}} \bar{\sigma}_{\max,*}^{2}\right)^{-1/2}$$
(84)

Note that $1 - 4\eta \frac{\bar{\delta}_k}{(1 - \bar{\delta}_k)^2} \bar{\sigma}_{\max,*}^2$ is strictly positive as long as $\frac{\bar{\delta}_k}{(1 - \bar{\delta}_k)^2} < 1$, which we will verify shortly, due to our earlier assumption that $\eta \leq 1/(4\bar{\sigma}_{\max,*}^2)$. Therefore, from (69), we have

$$\operatorname{dist}(\hat{\mathbf{B}}^{t+1}, \hat{\mathbf{B}}^*) \leq \frac{1}{\sqrt{1 - 4\eta \frac{\bar{\delta}_k}{(1 - \bar{\delta}_k)^2} \bar{\sigma}_{\max,*}^2}} \left(1 - \eta \bar{\sigma}_{\min,*}^2 E_0 + 2\eta \frac{\bar{\delta}_k}{(1 - \bar{\delta}_k)^2} \bar{\sigma}_{\max,*}^2\right) \operatorname{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*)$$

Next, let $\bar{\delta}_k < 16E_0/(25 \cdot 5\kappa^2)$. This implies that $\bar{\delta}_k < 1/5$. Then $\bar{\delta}_k/(1-\bar{\delta}_k)^2 < 25\bar{\delta}_k/16 \le E_0/(5\kappa^2) \le 1$, validating (84). Further, it is easily seen that

$$1 - \eta E_0 \bar{\sigma}_{\min,*}^2 + \eta \frac{\delta_k}{(1 - \bar{\delta}_k)^2} \bar{\sigma}_{\max,*}^2 \le 1 - 4\eta \frac{\delta_k}{(1 - \bar{\delta}_k)^2} \bar{\sigma}_{\max,*}^2$$

$$\le 1 - \eta E_0 \bar{\sigma}_{\min,*}^2 / 2$$
(85)

Thus

$$\operatorname{dist}(\hat{\mathbf{B}}^{t+1}, \hat{\mathbf{B}}^*) \le (1 - \eta E_0 \bar{\sigma}_{\min,*}^2 / 2)^{1/2} \operatorname{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*).$$

Finally, recall that $\bar{\delta}_k = \delta_k + \delta_k' = c \left(\frac{k^{3/2} \sqrt{\log(rn)}}{\sqrt{m}} + \frac{k\sqrt{d}}{\sqrt{rnm}} \right)$ for some absolute constant c. Choosing $m \geq c'(\kappa^4 k^3 \log(rn)/E_0^2 + \kappa^4 k^2 d/(E_0^2 rn))$ for another absolute constant c' satisfies $\bar{\delta}_k \leq 16E_0/(25 \cdot 5\kappa^2)$. Also, we have conditioned on two events, described in Lemmas 4 and 5, which occur with probability at least $1 - e^{-110d} - e^{-110k^2 \log(rn)} \geq 1 - e^{-100 \min(k^2 \log(rn), d)}$, completing the proof.

Finally, Theorem 1 follows by recursively applying Lemma 6 and taking a union bound over all $t \in [T]$.

B.4 Initialization

As mentioned in the main body, our interpretation of Theorem 1 assumes that the initial distance is bounded above by a constant less than one, i.e., E_0 is bounded below by a constant greater than zero. We can achieve such an initialization without increasing the overall sample complexity via the Method-of-Moments algorithm, ignoring log factors. To show this, we adapt a result from Tripuraneni et al. [2020a].

Theorem 2 (Theorem 3, Tripuraneni et al. [2020a]). In addition to Assumptions 2 and 3, suppose that $\mathbf{x}_i^{0,j} \sim \mathcal{N}(0,\mathbf{I}_d)$ independently for all $i \in [n], j \in [m]$. If each client $i \in [n]$ sends the server $\mathbf{Z}_i := \frac{1}{m} \sum_{j=1}^m (y_i^{0,j})^2 \mathbf{x}_i^{0,j} (\mathbf{x}_i^{0,j})^{\top}$ and the server computes $\hat{\mathbf{U}} \mathbf{D} \hat{\mathbf{U}}^{\top} \leftarrow rank-k$ $SVD(\frac{1}{n} \sum_{i=1}^n \mathbf{Z}_i)$ and sets $\mathbf{B}^0 = \hat{\mathbf{U}}$. Then, if $m \geq c \operatorname{polylog}(d, mn) \tilde{\kappa}^2 kd/(\sigma_{\min,*}^2 n)$,

$$\operatorname{dist}\left(\mathbf{B}^{0}, \mathbf{B}^{*}\right) \leq \tilde{O}\left(\frac{\tilde{\kappa}^{2}kd}{\sigma_{\min,*}^{2}mn}\right) \tag{86}$$

with probability at least $1-O((mn)^{-100})$ for some absolute constant c, where $\sigma_{\min,*} := \sigma_{\min}\left(\frac{1}{\sqrt{kn}}\mathbf{W}^*\right)$, $\tilde{\kappa} := \frac{\sigma_{\max}\left(\frac{1}{\sqrt{kn}}\mathbf{W}^*\right)}{\sigma_{\min,*}}$ and $\tilde{O}(\cdot)$ hides log factors.

The above result is a direct adaptation of Theorem 3 in Tripuraneni et al. [2020a] so we omit the proof. Note that the $\frac{1}{\sqrt{k}}$ factor in the definition of $\sigma_{\min,*}$ is a scaling factor to enforce consistency with the assumption that $\|\mathbf{w}_i^*\| = \Theta(1)$ in Tripuraneni et al. [2020a] (since we have assumed $\|\mathbf{w}_i^*\| = \sqrt{k}$). This result shows that $m_{\text{init}} = \tilde{\Omega}(\frac{\tilde{\kappa}^2 k d}{\sigma_{\min,*}^2 n})$ samples are required for proper initialization. Since $\frac{1}{\sigma_{\min,*}^2} \leq k\kappa^2$ (as $\sigma_{\max,*}^2 \geq 1/k$, see (69)), the overall sample complexity does not increase by more than log factors.

B.5 Proof Challenges

We next discuss two analytical challenges involved in proving Theorem 1.

(i) Row-wise sparse measurements. Recall that the measurement matrices $\mathbf{A}_{i,j}^t$ have non-zero elements only in the *i*-th row. This property is beneficial in the sense that it allows for distributing

computation across the *n* clients. However, it also means that the operators $\{\frac{1}{\sqrt{m}}\mathcal{A}^t\}_t$ do not satisfy Restricted Isometry Property (RIP), which therefore prevents us from using standard RIP-based analysis. The RIP is defined as follows:

Definition 4 (Restricted Isometry Property). An operator $\mathcal{B}: \mathbb{R}^{n \times d} \to \mathbb{R}^{nm}$ satisfies the k-RIP with parameter $\delta_k \in [0,1)$ if and only if

$$(1 - \delta_k) \|\mathbf{M}\|_F^2 \le \|\mathcal{B}(\mathbf{M})\|_2^2 \le (1 + \delta_k) \|\mathbf{M}\|_F^2$$
(87)

holds simultaneously for all $\mathbf{M} \in \mathbb{R}^{n \times d}$ of rank at most k.

Claim 1. Let $\mathcal{A}: \mathbb{R}^{rn \times d} \to \mathbb{R}^{rnm}$ such that $\mathcal{A}(\mathbf{M}) = [\langle \mathbf{e}_i(\mathbf{x}_i^j)^\top, \mathbf{M} \rangle]_{1 \leq i \leq rn, 1 \leq j \leq m}$, and let the samples \mathbf{x}_i^j be i.i.d. sub-gaussian random vectors with mean $\mathbf{0}_d$ and covariance \mathbf{I}_d . Then if $m \leq d/2$, with probability at least $1 - e^{-cd}$ for some absolute constant $c, \frac{1}{\sqrt{m}} \mathcal{A}$ does not satisfy 1-RIP for any constant $\delta_1 \in [0,1)$.

Proof. Let $\mathbf{M} = \mathbf{e}_1(\mathbf{x}_1^1)^{\top}$. Then

$$\|\frac{1}{\sqrt{m}}\mathcal{A}(\mathbf{M})\|_{2}^{2} = \frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{m} \langle \mathbf{e}_{i}(\mathbf{x}_{i}^{j})^{\top}, \mathbf{e}_{1}(\mathbf{x}_{1}^{1})^{\top} \rangle^{2}$$

$$= \frac{1}{m} \|\mathbf{x}_{1}^{1}\|_{2}^{4} + \frac{1}{m} \sum_{j=2}^{m} \langle \mathbf{x}_{1}^{j}, \mathbf{x}_{1}^{1} \rangle^{2}$$

$$\geq \frac{1}{m} \|\mathbf{x}_{1}^{1}\|_{2}^{4}$$

$$(88)$$

Also observe that $\|\mathbf{M}\|_F^2 = \|\mathbf{x}_1^1\|_2^2$. Therefore, we have

$$\mathbb{P}\left(\frac{\left\|\frac{1}{\sqrt{m}}\mathcal{A}(\mathbf{M})\right\|_{2}^{2}}{\left\|\mathbf{M}\right\|_{F}^{2}} \geq \frac{d}{2m}\right) \geq \mathbb{P}\left(\frac{\frac{1}{m}\left\|\mathbf{x}_{1}^{1}\right\|_{2}^{4}}{\left\|\mathbf{x}_{1}^{1}\right\|_{2}^{2}} \geq \frac{d}{2m}\right)$$

$$= \mathbb{P}\left(\left\|\mathbf{x}_{1}^{1}\right\|_{2}^{2} \geq \frac{d}{2}\right)$$

$$= 1 - \mathbb{P}\left(\left\|\mathbf{x}_{1}^{1}\right\|_{2}^{2} - d \leq \frac{-d}{2}\right)$$

$$\geq 1 - e^{-cd}$$
(89)

where the last inequality follows for some absolute constant c by the sub-exponential property of $\|\mathbf{x}_1^1\|_2^2$ and the fact that $\mathbb{E}[\|\mathbf{x}_1^1\|_2^2] = d$. Thus, with probability at least $1 - e^{-cd}$, $\left\|\frac{1}{\sqrt{m}}\mathcal{A}(\mathbf{M})\right\|_2^2 \ge \frac{d}{2m}\|\mathbf{M}\|_2^2$, meaning that $\frac{1}{\sqrt{m}}\mathcal{A}$ does not satisfy 1-RIP with high probability if $m \le \frac{d}{2}$.

Claim 1 shows that we cannot use the RIP to show $\mathcal{O}(d/(rn))$ sample complexity for m - instead, this approach would require $m = \Omega(d)$. Fortunately, we do not need concentration of the measurements for all rank-k matrices \mathbf{M} , but only a particular class of rank-k matrices that are row-wise incoherent, due to the row-wise incoherence of \mathbf{W}^* (see Assumption 3 and Definition 3). Leveraging the row-wise incoherence of the matrices being measured allows us to show that we only require $m = \Omega(k^3 \log(rn) + k^2 d/(rn))$ samples per user (ignoring dimension-independent constants).

(ii) Non-symmetric updates. Existing analyses for nonconvex matrix sensing study algorithms with symmetric update schemes for the factors **W** and **B**, either alternating minimization, e.g. [Jain et al., 2013], or alternating gradient descent, e.g. [Tu et al., 2016]. Here we show contraction due to the gradient descent step in principal angle distance, differing from the standard result for gradient descent using Procrustes distance [Park et al., 2018, Tu et al., 2016, Zheng and Lafferty, 2016]. We combine aspects of both types of analysis in our proof.

References

- Rie Kubota Ando, Tong Zhang, and Peter Bartlett. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6(11), 2005.
- Manoj Ghuhan Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Federated learning with personalization layers. arXiv preprint arXiv:1912.00818, 2019.
- Maria-Florina Balcan, Avrim Blum, and Santosh Vempala. Efficient representations for lifelong learning and autoencoding. In *Conference on Learning Theory*, pages 191–210. PMLR, 2015.
- Jonathan Baxter. A model of inductive bias learning. *Journal of artificial intelligence research*, 12: 149–198, 2000.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- Brian Bullins, Elad Hazan, Adam Kalai, and Roi Livni. Generalize across tasks: Efficient algorithms for linear representation learning. In *Algorithmic Learning Theory*, pages 235–246. PMLR, 2019.
- Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings. arXiv preprint arXiv:1812.01097, 2018.
- Fei Chen, Mi Luo, Zhenhua Dong, Zhenguo Li, and Xiuqiang He. Federated meta-learning with fast convergence and efficient communication. arXiv preprint arXiv:1802.07876, 2018.
- Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In 2017 International Joint Conference on Neural Networks (IJCNN), pages 2921–2926. IEEE, 2017.
- Giulia Denevi, Carlo Ciliberto, Dimitris Stamos, and Massimiliano Pontil. Learning to learn around a common mean. In *ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 31* (NIPS 2018), volume 31. NIPS Proceedings, 2018.
- Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive personalized federated learning. arXiv preprint arXiv:2003.13461, 2020.
- Simon S. Du, Wei Hu, Sham M. Kakade, Jason D. Lee, and Qi Lei. Few-shot learning via learning the representation, provably, 2020.

- Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning: A meta-learning approach, 2020.
- Farzin Haddadpour, Mohammad Mahdi Kamani, Aryan Mokhtari, and Mehrdad Mahdavi. Federated learning with compression: Unified analysis and sharp guarantees. arXiv preprint arXiv:2007.01154, 2020.
- Filip Hanzely and Peter Richtárik. Federated learning of a mixture of global and local models. arXiv preprint arXiv:2002.05516, 2020.
- Daniel Hsu, Sham M Kakade, and Tong Zhang. Random design analysis of ridge regression. In *Conference on learning theory*, pages 9–1. JMLR Workshop and Conference Proceedings, 2012.
- Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating minimization. *Proceedings of the 45th annual ACM symposium on Symposium on theory of computing STOC '13*, 2013.
- Yihan Jiang, Jakub Konečný, Keith Rush, and Sreeram Kannan. Improving federated learning personalization via model agnostic meta learning. arXiv preprint arXiv:1909.12488, 2019.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- Mikhail Khodak, Maria-Florina F Balcan, and Ameet S Talwalkar. Adaptive gradient-based metalearning methods. In *Advances in Neural Information Processing Systems*, pages 5915–5926, 2019.
- Weihao Kong, Raghav Somani, Zhao Song, Sham Kakade, and Sewoong Oh. Meta-learning for mixed linear regression. In *International Conference on Machine Learning*, pages 5394–5404. PMLR, 2020.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. nature, 521(7553):436–444, 2015.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. arXiv preprint arXiv:1812.06127, 2018.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Feddane: A federated newton-type method. In 2019 53rd Asilomar Conference on Signals, Systems, and Computers, pages 1227–1231. IEEE, 2019.
- Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. arXiv: 2012.04221, 2020.
- Paul Pu Liang, Terrance Liu, Liu Ziyin, Ruslan Salakhutdinov, and Louis-Philippe Morency. Think locally, act globally: Federated learning with local and global representations. arXiv preprint arXiv:2001.01523, 2020.
- Yishay Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three approaches for personalization with applications to federated learning. arXiv preprint arXiv:2002.10619, 2020.

- Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning. *Journal of Machine Learning Research*, 17(81):1–32, 2016.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- Aritra Mitra, Rayana Jaafar, George J Pappas, and Hamed Hassani. Achieving linear convergence in federated learning under objective and systems heterogeneity. arXiv preprint arXiv:2102.07053, 2021.
- Dohyung Park, Anastasios Kyrillidis, Constantine Caramanis, and Sujay Sanghavi. Finding low-rank solutions via nonconvex matrix factorization, efficiently and provably. *SIAM Journal on Imaging Sciences*, 11(4):2165–2204, 2018.
- Reese Pathak and Martin J Wainwright. Fedsplit: An algorithmic framework for fast federated optimization. arXiv preprint arXiv:2005.05238, 2020.
- Massimiliano Pontil and Andreas Maurer. Excess risk bounds for multitask learning with trace norm regularization. In *Conference on Learning Theory*, pages 55–76. PMLR, 2013.
- Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse? towards understanding the effectiveness of maml. arXiv preprint arXiv:1909.09157, 2019.
- Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. arXiv preprint arXiv:2003.00295, 2020.
- Amirhossein Reisizadeh, Isidoros Tziotis, Hamed Hassani, Aryan Mokhtari, and Ramtin Pedarsani. Straggler-resilient federated learning: Leveraging the interplay between statistical accuracy and system heterogeneity. arXiv preprint arXiv:2012.14453, 2020.
- Irina Rish, Genady Grabarnik, Guillermo Cecchi, Francisco Pereira, and Geoffrey J Gordon. Closed-form supervised dimensionality reduction with generalized linear models. In *Proceedings of the 25th international conference on Machine learning*, pages 832–839, 2008.
- Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar. Federated multi-task learning. In *Advances in neural information processing systems*, pages 4424–4434, 2017.
- Nilesh Tripuraneni, Chi Jin, and Michael I. Jordan. Provable meta-learning of linear representations, 2020a.
- Nilesh Tripuraneni, Michael I Jordan, and Chi Jin. On the theory of transfer learning: The importance of task diversity. arXiv preprint arXiv:2006.11650, 2020b.
- Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *International Conference on Machine Learning*, pages 964–973. PMLR, 2016.
- Roman Vershynin. High-dimensional probability: An introduction with applications in data science, volume 47. Cambridge university press, 2018.

- Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. arXiv preprint arXiv:2007.07481, 2020.
- Kangkang Wang, Rajiv Mathews, Chloé Kiddon, Hubert Eichner, Françoise Beaufays, and Daniel Ramage. Federated evaluation of on-device personalization. arXiv preprint arXiv:1910.10252, 2019.
- Tao Yu, Eugene Bagdasaryan, and Vitaly Shmatikov. Salvaging federated learning by local adaptation, 2020.
- Qinqing Zheng and John Lafferty. Convergence analysis for rectangular matrix completion using burer-monteiro factorization and gradient descent. arXiv preprint arXiv:1605.07051, 2016.
- Kai Zhong, Prateek Jain, and Inderjit S Dhillon. Efficient matrix sensing using rank-1 gaussian measurements. In *International conference on algorithmic learning theory*, pages 3–18. Springer, 2015.