

On Sample Complexity of Learning Shared Representations: The Asymptotic Regime

Xinyi Tong¹, Xiangxiang Xu², and Shao-Lun Huang³

Abstract—Federated learning and multi-task learning have achieved great successes in recent works, which exploit the knowledge contained in different tasks to alleviate the influence of lacking training samples. In order to understand the knowledge utilization mechanism, we study the sample complexity of the federated learning algorithms by computing their expected population risks, where the asymptotic regime is considered that the number of tasks is large. Specifically, we focus on several typical algorithms, including (i) training personalized models merely by local data, (ii) jointly training a single model by all the data of devices, and (iii) jointly training the shared feature representation and learning individual classifiers of different devices, where the latter one is recently widely-adopted. The results reveal that the validity of the shared representation algorithm is affected by the task similarities, the model dimensionality, and the sample size, which quantify the sample size range where such algorithm is preferable. Moreover, we develop a new algorithm in consideration of collaborating the classifiers of different devices under the shared representations, which can have lower sample complexity than conventional methods under the asymptotic regime of large sample size. This provides theoretical guidance for practical algorithm designs.

I. INTRODUCTION

Simultaneously solving multiple machine learning tasks on the same domain of data has gained much attention in federated learning and multi-task learning [1]. For example, it is often considered in federated learning that there are many distributed devices and a server, where each device can collect data for solving its own learning task. The learning tasks across the devices can be correlated, and the devices can interact with the server for jointly solving the learning tasks. In multi-task learning, there are several correlated machine learning tasks with different training data, and the goal is to jointly solve these tasks by using the task correlations and available data. One critical challenge for the abovementioned problems is that the training data size of each task may not be sufficient for training a good machine learning model. As such, the interactions between different tasks and certain kind of data sharing in joint model training are often necessary for achieving better learning performance. A naive approach considered in federated learning for addressing this problem is to learning a single model for all different tasks via the training data for all tasks. When the learning tasks across

the devices are considerably similar, and the size of the training data from all devices are sufficiently large, the overall performance of all the learning tasks can be improved comparing to solving each task with its own data. However, when the learning tasks of different devices are not quite similar, such approach can lead to significant performance degradation [2]. Therefore, instead of sharing the entire machine learning model across different tasks, it is recently proposed in [3] to share only the feature representation of the data with different classifiers. Specifically, this approach applies all training data to train a shared feature representation for all tasks, and then each task trains its own linear classifier according to the shared feature representation and its own training data. Intuitively, such approach can perform well when different tasks are correlated but not considerably similar, and has gained successes in some applications. Nevertheless, the theoretical understanding of why and when learning shared representation can be beneficial appears to be quite behind. Therefore, our goal is to provide a theoretical framework to quantify the positive/negative gains of learning shared representation in different scenarios.

In this paper, we consider the asymptotic regime, such that the number of tasks k is large, which is practical in many federated learning applications [4]. Under such regime, we first investigate three typical approaches in federated learning: (i) Each device learns a model only by its own training data; (ii) All devices jointly learn a single model for all tasks by all training data; (iii) All devices jointly learn a shared representation by all training data, and each device learns its own classifier by its training data. These three approaches are illustrated in Figure 1. The learning model of these approaches are trained by the empirical risk minimization (ERM) through the training data, and the performance of the algorithms are evaluated by the expected population risks (EPR). In particular, we demonstrate analytical characterizations of the EPRs in the asymptotic regime for the three approaches in terms of the underlying distributions and the number of training samples of each task. Then, we compare the EPRs for different approaches, and specify the parameter regimes such that learning shared representation can be beneficial comparing to other approaches. Such regime can be expressed as a comparison between the task similarities between devices and the number of training samples at each device, which quantifies the intuition of learning shared representations.

Moreover, it is interesting to note that the approach (iii) does not exploit the correlations between the classifiers. Hence, we extend the approach (iii) by allowing the classifiers to collaborate via linear filtering. Thus, the classifier at

¹Xinyi Tong is with Tsinghua-Berkeley Shenzhen Institute, Tsinghua University, Shenzhen 518055, China txyl8@mails.tsinghua.edu.cn

²Xiangxiang Xu is with EECS, MIT, Cambridge, MA 02139, USA xuux@mit.edu

³Shao-Lun Huang is with the Data Science and Information Technology Research Center, Tsinghua-Berkeley Shenzhen Institute, Shenzhen 518055, China shaolun.huang@sz.tsinghua.edu.cn

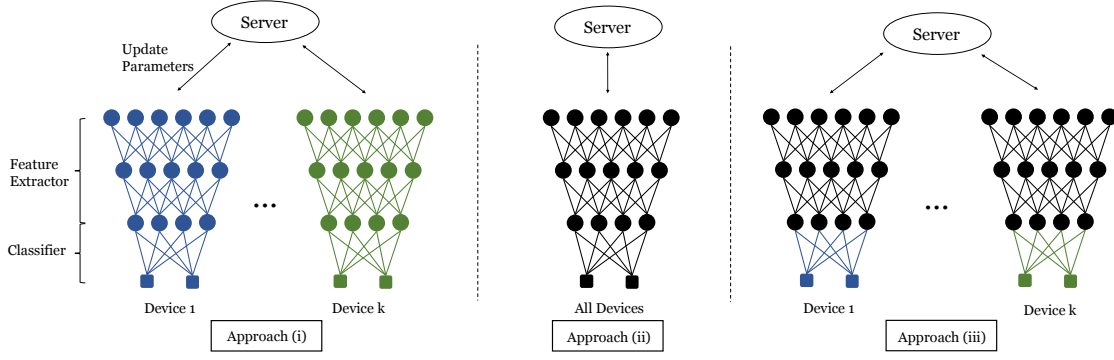


Fig. 1. A diagram to show the three federated learning algorithms, where the learning model can be separated as the feature extractor with many parameters and the classifier with few ones. Particularly, Approach (i) learns an individual model of each device only by its own data; Approach (ii) jointly learn a single model for all tasks; and Approach (iii) learns a shared feature extractor by all training data, where each device learns its own classifier.

each device in this new approach is represented by linear combinations of the classifiers from the approach (iii), where the combining coefficients are subjected to be designed to optimize the EPR. We present the analytical solution of the optimal coefficients in terms of the underlying distributions of the tasks and the number of training samples. This leads to a lower bound for the EPR of this new approach, since such optimal coefficients are not computable without the knowledge of the underlying distributions. Then, we propose a computable algorithm which estimates the coefficients by the training data, and present theoretical analyses of the EPR performance of this algorithm. Our results demonstrate the performance gain of the classifier collaborations in terms of the task similarities and the number of samples in the asymptotic regime. This provides the theoretical guidance for algorithm designs in federated learning and multi-task learning problems.

The contributions of this paper are summarized as follows:

- 1) We provide the sample complexity analyses of the aforementioned typical federated algorithms under the discrete data space by computing their EPRs under the large task number regime. Based on the results, we confirm the sample size range that the shared representation algorithm could outperform the other two algorithms and interpret the influence of the task similarity and the model dimensionality.
- 2) We improve the shared representation algorithm by collaborating the classifiers learned in different devices, where a lower bound of its EPR can be given under the optimal collaboration. We also provide a practical collaboration approach and verify its effectiveness under the regime of large sample sizes.
- 3) We extend the sample complexity analyses to the continuous sample space and provide the EPRs of different algorithms under the large sample size regime with respect to the number of parameters contained in the learning models, where similar interpretations as in the discrete case are provided.

A. Related Works

In consideration of the fact that approaches (i) and (ii) are not particularly intended for federated learning [5] and multi-

task learning [6], we maintain our main focus on approach (iii). Such a shared representation method is the major part of the representation learning area, which focuses on building up joint representations among different tasks [1]. The successes of shared representation methods have been demonstrated when multiple tasks are simultaneously tackled, e.g., in the image classification tasks [7]. Finally, theoretical analyses in this area are in progress recently, which mainly cover the typical Gaussian signal analyses based on the ERM principle [8], and the mechanism of concrete methods, e.g., low-rank method and task clustering method [9]. Compared with previous works, this paper can provide insights of how the sample sizes affect the effectiveness of the shared representation method based on a more general framework. On the other hand, since federated learning is concentrated on privacy and computation capacity, theoretical analyses in this area mainly focus on the convergence rate from the perspective of optimization [10].

II. PROBLEM FORMULATION

A. Notations

Let X and Y be the random variables of the data and label with domains \mathcal{X} and \mathcal{Y} , respectively. For ease of illustration, we assume X to be discrete in Section II and Section III, and will extend our analyses to continuous cases later.

To start with, we consider that the federated learning problem has k tasks. For each task of device $i = 1, \dots, k$, we assume that n training samples $\{(x_\ell^{(i)}, y_\ell^{(i)})\}_{\ell=1}^n$ are i.i.d. generated from some underlying joint distribution $P_{XY}^{(i)}$, and the corresponding empirical distribution of samples is defined as $\hat{P}_{XY}^{(i)}(x, y) \triangleq \frac{1}{n} \sum_{\ell=1}^n \mathbb{1}\{x_\ell^{(i)} = x, y_\ell^{(i)} = y\}$, where $\mathbb{1}\{\cdot\}$ denotes the indicator function [11].

In this paper, we study the sample complexity of federated learning algorithms under the asymptotic regime, where we make the following assumption:

Assumption 2.1 (Asymptotic Regime Assumption): The task number k is arbitrarily large.

To describe the asymptotic behaviors of the related EPRs when $k \rightarrow \infty$, we employ the big O notations [12], which is with respect to k for the whole paper. Specifically, when computing the EPR of some algorithms, we also assume the sample size n to be large [cf. Assumption 3.6]. Under this circumstance, the big O notation is also with respect to n .

In particular, we assume that all the joint distributions $P_{XY}^{(i)}$'s share the same marginal distributions of data and label, denoted as P_X^1 and P_Y , which is natural with all devices sharing the same domains and widely-adopted in practical federated learning implementations of i.i.d. data case [13]. We also make a mild assumption that each label occurs with the same probability, i.e., $P_Y(y) = 1/|\mathcal{Y}|$ for each $y \in \mathcal{Y}$, which is satisfied in most of the popular experimental settings [14]. Moreover, the underlying distribution of each device i is assumed to be decomposed by the feature $\mathbf{f}_i : \mathcal{X} \rightarrow \mathbb{R}^d$ and the classifier $\mathbf{g}_i : \mathcal{Y} \rightarrow \mathbb{R}^d$, where $d \geq |\mathcal{Y}|$, i.e., the (x, y) th entry of $P_{XY}^{(i)}$ is

$$P_{XY}^{(i)}(x, y) = P_X(x)P_Y(y) (1 + \mathbf{f}_i^T(x)\mathbf{g}_i(y)). \quad (1)$$

Furthermore, the feature and classifier are required to be zero-mean under the marginal distributions², i.e., $\mathbb{E}_{P_X}[\mathbf{f}_i(X)] = \mathbb{E}_{P_Y}[\mathbf{g}_i(Y)] = \mathbf{0}$. Then, we define the *canonical dependence matrix* (CDM) [15] $\mathbf{B}(P_{XY}) \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$ of any distribution P_{XY} supported by $\mathcal{X} \times \mathcal{Y}$:

$$[\mathbf{B}(P_{XY})](y, x) \triangleq \frac{P_{XY}(x, y) - P_X(x)P_Y(y)}{\sqrt{P_X(x)P_Y(y)}}, \quad (2)$$

where $[\mathbf{B}(P_{XY})](y, x)$ is the (y, x) th entry of $\mathbf{B}(P_{XY})$.

We also define the *feature matrix* $\Phi(\mathbf{f}) \in \mathbb{R}^{|\mathcal{X}| \times d}$ of any feature function $\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^d$:

$$\Phi(\mathbf{f}) \triangleq \left[\sqrt{P_X(1)}\mathbf{f}^T(1), \dots, \sqrt{P_X(|\mathcal{X}|)}\mathbf{f}^T(|\mathcal{X}|) \right]^T, \quad (3)$$

and the *classifier matrix* $\Psi(\mathbf{g}) \in \mathbb{R}^{|\mathcal{Y}| \times d}$ of any label classifier $\mathbf{g} : \mathcal{Y} \rightarrow \mathbb{R}^d$:

$$\Psi(\mathbf{g}) \triangleq \left[\sqrt{P_Y(1)}\mathbf{g}^T(1), \dots, \sqrt{P_Y(|\mathcal{Y}|)}\mathbf{g}^T(|\mathcal{Y}|) \right]^T. \quad (4)$$

For simplicity, we define the corresponding matrices of each device i as $\Phi_i \triangleq \Phi(\mathbf{f}_i)$, $\Psi_i \triangleq \Psi(\mathbf{g}_i)$, $\forall i = 1, \dots, k$. With such notations, the CDM of the underlying distribution of each device i is $\mathbf{B}(P_{XY}^{(i)}) = \Psi_i \Phi_i^T$. Notice the fact that $\mathbf{B}(P_{XY}^{(i)})$ keeps invariant under the invertible linear mapping in the feature space \mathbb{R}^d , i.e., for any invertible matrix $\mathbf{Q} \in \mathbb{R}^{d \times d}$, $\Psi_i \Phi_i^T = \Psi_i \mathbf{Q}^{-1} \mathbf{Q} \Phi_i^T$. We can without loss of generality assume that for each $i = 1, \dots, k$, $\Phi_i^T \Phi_i = \mathbf{I}_d$, where \mathbf{I}_d denotes the identity matrix in $\mathbb{R}^{d \times d}$.

Finally, we introduce the Frobenius matrix norm [16]: for each matrix $\mathbf{A} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$, $\|\mathbf{A}\|_F^2 \triangleq \text{tr}(\mathbf{A}\mathbf{A}^T) = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} [\mathbf{A}(x, y)]^2$. Such norm is employed for evaluating the performance of the learned features and classifiers. In particular, when $\hat{\mathbf{f}} : \mathcal{X} \rightarrow \mathbb{R}^d$ and $\hat{\mathbf{g}} : \mathcal{Y} \rightarrow \mathbb{R}^d$ are learned from samples for device i , the expected population risk of

the learned features³ is

$$\begin{aligned} R_i(\hat{\mathbf{f}}, \hat{\mathbf{g}}) &\triangleq \mathbb{E} \left[\left\| \mathbf{B}(P_{XY}^{(i)}) - \Psi(\hat{\mathbf{g}})\Phi^T(\hat{\mathbf{f}}) \right\|_F^2 \right] \\ &= \mathbb{E} \left[\left\| \Psi_i \Phi_i^T - \Psi(\hat{\mathbf{g}})\Phi^T(\hat{\mathbf{f}}) \right\|_F^2 \right], \end{aligned} \quad (5)$$

where the expectation is taken over all possible empirical distributions $\hat{P}_{XY}^{(i)}$ following the multinomial distributions [19].

B. Federated learning algorithms

First, we evaluate three typical federated learning algorithms as defined in the following.

Personalized Model Algorithm (Algorithm (i)):

When the samples of each device are sufficient to learn a well-performed model, devices do not require assistance from the others. Under this circumstance, the most trivial way is to learn the personalized model by its own training data, which means for $i = 1, \dots, k$, the feature of data and the classifier of label can be achieved by minimizing the empirical risk

$$\begin{aligned} &(\hat{\mathbf{f}}_{\text{PM}i}, \hat{\mathbf{g}}_{\text{PM}i}) \\ &= \arg \min_{(\mathbf{f}, \mathbf{g}) : \Psi^T(\mathbf{f})\Psi(\mathbf{g}) = \mathbf{I}_d} \left\| \mathbf{B}(\hat{P}_{XY}^{(i)}) - \Psi(\mathbf{g})\Phi^T(\mathbf{f}) \right\|_F^2. \end{aligned} \quad (6)$$

The corresponding EPR is

$$R_{\text{PM}} \triangleq \sum_{i=1}^k R_i(\hat{\mathbf{f}}_{\text{PM}i}, \hat{\mathbf{g}}_{\text{PM}i}). \quad (7)$$

Global Model Algorithm (Algorithm (ii)):

When local samples are extremely insufficient, all devices are forced to learn a global model jointly, otherwise the random noise from the sampling process can lead to a large generalization error. Under this circumstance, the global feature of data and classifier of label can be achieved by

$$\begin{aligned} &(\hat{\mathbf{f}}_{\text{GM}}, \hat{\mathbf{g}}_{\text{GM}}) \\ &= \arg \min_{(\mathbf{f}, \mathbf{g}) : \Psi^T(\mathbf{f})\Psi(\mathbf{g}) = \mathbf{I}_d} \left\| \mathbf{B}\left(\frac{1}{k} \sum_{i=1}^k \hat{P}_{XY}^{(i)}\right) - \Psi(\mathbf{g})\Phi^T(\mathbf{f}) \right\|_F^2. \end{aligned} \quad (8)$$

The corresponding EPR is

$$R_{\text{GM}} \triangleq \sum_{i=1}^k R_i(\hat{\mathbf{f}}_{\text{GM}}, \hat{\mathbf{g}}_{\text{GM}}). \quad (9)$$

Shared Representation Algorithm (Algorithm (iii)):

As introduced in [3], this typical federated learning algorithm shares only the feature representation of the data and learns different classifiers. Under this circumstance, the joint feature representation and individual classifiers of different

¹Without loss of generality, we assume all the entries $P_X(x) > 0$, or we can exclude this element x from the space \mathcal{X} .

²This is not an extra assumption, which is a conclusion of the canonical correlation analysis of joint distribution decomposition. [cf. Proposition 5 [15]]

³Such performance measure based on CDM can be recognized as a good approximation of Kullback-Leibler (K-L) divergence [17] between the underlying distribution $P_{XY}^{(i)}(x, y)$ and the distribution model with learned features $Q_{XY}(x, y) \triangleq P_X(x)P_Y(y)(1 + \hat{\mathbf{f}}^T(x)\hat{\mathbf{g}}(y))$, where $D(P_{XY}^{(i)} \| Q_{XY}) \approx \frac{1}{2} \left\| \mathbf{B}(P_{XY}^{(i)}) - \Psi(\hat{\mathbf{g}})\Phi^T(\hat{\mathbf{f}}) \right\|_F^2$. [18]

devices can be achieved by an iterative algorithm, whose results satisfy

$$\hat{\mathbf{f}}_{\text{SR}} = \arg \min_{\mathbf{f}: \Psi^T(\mathbf{f})\Psi(\mathbf{f})=\mathbf{I}_d} \sum_{i=1}^k \left\| \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) - \Psi(\hat{\mathbf{g}}_{\text{SR}i}) \Phi^T(\mathbf{f}) \right\|_{\text{F}}^2, \quad (10)$$

and

$$\hat{\mathbf{g}}_{\text{SR}i} = \arg \min_{\mathbf{g}} \left\| \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) - \Psi(\mathbf{g}) \Phi^T(\hat{\mathbf{f}}_{\text{SR}}) \right\|_{\text{F}}^2. \quad (11)$$

The corresponding EPR is

$$R_{\text{SR}} \triangleq \sum_{i=1}^k R_i \left(\hat{\mathbf{f}}_{\text{SR}}, \hat{\mathbf{g}}_{\text{SR}i} \right). \quad (12)$$

Classifiers Collaboration Algorithm (Algorithm (iv)):

As illustrated in Section I, Algorithms (i)-(iii) are the classical algorithms that have been widely-applied in practical implementations. Moreover, we consider that the correlations between the classifiers are not fully exploited, where the classifiers learned in Algorithm (iii) can be collaborated via linear filtering. We simply consider the new classifier of device i as a linear combination of the classifier $\hat{\mathbf{g}}_{\text{SR}i}$ trained by its own data and the averaged learned classifier $\hat{\mathbf{g}}_{\text{SR}}$, where $\hat{\mathbf{g}}_{\text{SR}} \triangleq \frac{1}{k} \sum_{i=1}^k \hat{\mathbf{g}}_{\text{SR}i}$. Then, we have the following classifier for device i :

$$\hat{\mathbf{g}}_{\text{CC}i} = \alpha_i \hat{\mathbf{g}}_{\text{SR}i} + (1 - \alpha_i) \hat{\mathbf{g}}_{\text{SR}}, \quad (13)$$

where α_i is a function of all the data as given in Section III-B. The corresponding EPR is

$$R_{\text{CC}} \triangleq \sum_{i=1}^k R_{\text{CC}i}, \text{ where } R_{\text{CC}i} \triangleq R_i \left(\hat{\mathbf{f}}_{\text{SR}}, \hat{\mathbf{g}}_{\text{CC}i} \right). \quad (14)$$

III. EPR ANALYSES: THE DISCRETE CASE

A. Analyses for classical algorithms

First, we directly establish the analytical expressions of the EPRs (7), (9), and (12). The proofs are provided in the supplementary material.

1) *Personalized model algorithm*: For the Personalized Model Algorithm (Algorithm (i)), it can be easily derived from the definition (6) that the learned $\hat{\mathbf{f}}_{\text{PM}i}$ and $\hat{\mathbf{g}}_{\text{PM}i}$ correspond to the empirical distribution $\hat{P}_{XY}^{(i)}$ as the objective function being zero, i.e., $\Psi(\hat{\mathbf{g}}_{\text{PM}i}) \Phi^T(\hat{\mathbf{f}}_{\text{PM}i}) = \mathbf{B}(\hat{P}_{XY}^{(i)})$. Then, we have the following characterization of the EPR (7).

Theorem 3.1 (EPR of Personalized Model Algorithm): The expected population risk as defined in (7) is

$$R_{\text{PM}} = \frac{k|\mathcal{X}||\mathcal{Y}| - k + \sum_{i=1}^k C_i}{n}, \quad (15)$$

where $C_i \triangleq \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \mathbf{f}_i^T(x) \mathbf{g}_i(y) - \|\Psi_i\|_{\text{F}}^2$, and Ψ_i is the classifier matrix of \mathbf{g}_i .

Let $C_{\min} \triangleq \min_{x \in \mathcal{X}, y \in \mathcal{Y}, i \in \{1, \dots, k\}} \{\mathbf{f}_i^T(x) \mathbf{g}_i(y)\} - \frac{1+|\mathcal{Y}|}{|\mathcal{X}||\mathcal{Y}|}$ and $C_{\max} \triangleq \max_{x \in \mathcal{X}, y \in \mathcal{Y}, i \in \{1, \dots, k\}} \{\mathbf{f}_i^T(x) \mathbf{g}_i(y)\}$. Then, R_{PM} can be

bounded⁴ by

$$\frac{k(1 + C_{\min})|\mathcal{X}||\mathcal{Y}|}{n} \leq R_{\text{PM}} \leq \frac{k(1 + C_{\max})|\mathcal{X}||\mathcal{Y}|}{n}.$$

Additionally, Theorem 3.1 states that the Personalized Model Algorithm provides an unbiased estimation of the underlying distributions $P_{XY}^{(i)}$'s and all the EPRs come from the variances of the sampling process, which is inversely proportional to the sample size n of each device. When estimating the underlying distributions, there are $|\mathcal{X}||\mathcal{Y}|$ parameters (number of entries) contained in the distribution models, which appear in the EPR (15) and reflect the model dimensionality.

2) *Global Model Algorithm*: As for the Global Model Algorithm (Algorithm (ii)), the feature and classifier are learned with nk samples of all the devices, which lead to an extremely small variance and the main EPR comes from the bias by using the global empirical distribution $\frac{1}{k} \sum_{i=1}^k \hat{P}_{XY}^{(i)}$ to estimate each device's distribution. Let

$$\begin{aligned} &(\mathbf{f}_{\text{GM}}, \mathbf{g}_{\text{GM}}) \\ &\triangleq \arg \min_{(\mathbf{f}, \mathbf{g}): \Phi^T(\mathbf{f})\Phi(\mathbf{f})=\mathbf{I}_d} \left\| \mathbf{B} \left(\frac{1}{k} \sum_{i=1}^k \hat{P}_{XY}^{(i)} \right) - \Psi(\mathbf{g}) \Phi^T(\mathbf{f}) \right\|_{\text{F}}^2. \end{aligned} \quad (16)$$

It is easy to verify that $\Psi(\mathbf{g}_{\text{GM}}) \Phi^T(\mathbf{f}_{\text{GM}}) = \mathbf{B} \left(\frac{1}{k} \sum_{i=1}^k \hat{P}_{XY}^{(i)} \right)$. Then we have the following theorem.

Theorem 3.2 (EPR of Global Model Algorithm): With Assumption 2.1, the expected population risk as defined in (9) is

$$R_{\text{GM}} = kV_{\text{GM}} + \frac{|\mathcal{X}||\mathcal{Y}|}{n} \cdot O(1), \quad (17)$$

where

$$V_{\text{GM}} \triangleq \frac{1}{k} \sum_{i=1}^k \left\| \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) - \mathbf{B} \left(\frac{1}{k} \sum_{j=1}^k \hat{P}_{XY}^{(j)} \right) \right\|_{\text{F}}^2. \quad (18)$$

This quantity V_{GM} can be seen as the variance of different distributions, which does not decay to 0 with k . Overall, Theorem 3.2 states that the Global Model Algorithm provides a biased estimation and the variance term can be neglected under a large k .

3) *Shared Representation Algorithm*: For the Shared Representation Algorithm (Algorithm (iii)), the feature representation is learned with nk samples of all the devices, while each classifier is merely learned with its own n samples. We define

$$\mathbf{f}_{\text{SR}} \triangleq \arg \max_{\mathbf{f}: \Phi^T(\mathbf{f})\Phi(\mathbf{f})=\mathbf{I}_d} \Phi^T(\mathbf{f}) \left[\frac{1}{k} \sum_{i=1}^k (\Phi_i \Psi_i^T \Psi_i \Phi_i^T) \right] \Phi(\mathbf{f}), \quad (19)$$

and the corresponding classifiers $\mathbf{g}_{\text{SR}i}$ as $\Psi(\mathbf{g}_{\text{SR}i}) \triangleq \Psi_i \Phi_i^T \Phi(\mathbf{f}_{\text{SR}})$, for each $i = 1, \dots, k$. Then, we can have the following characterization of the EPR (12).

⁴ $\|\Phi_i\|_{\text{F}}^2$ corresponds to the summation of all the Hirschfeld-Gebelein-Rényi (HGR) correlation coefficients of $P_{XY}^{(i)}$, such that $\|\Phi_i\|_{\text{F}}^2 \leq |\mathcal{Y}|$. [cf. Equation (30) [15]]

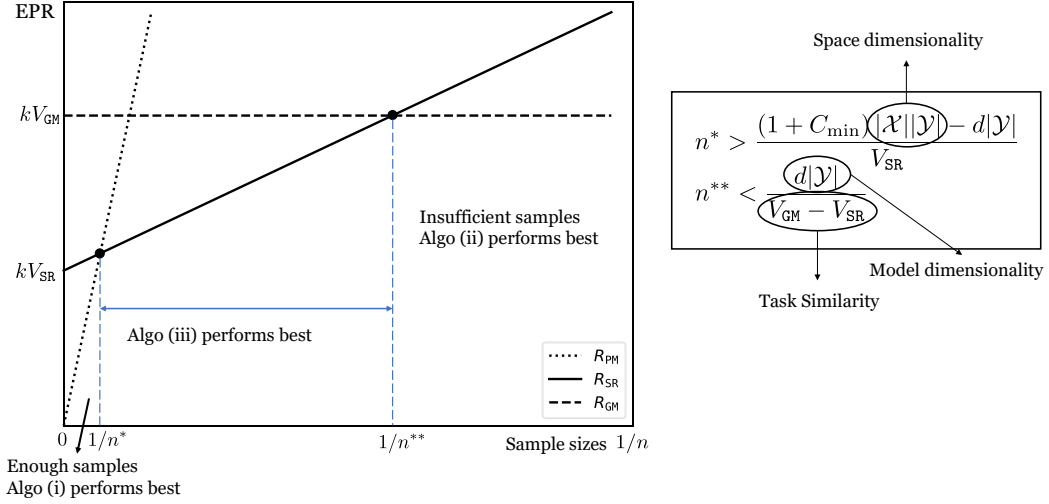


Fig. 2. A Summary of the EPRs as a function of $1/n$, which includes the Personalized Model Algorithm, the Global Model Algorithm, and the Shared Representation Algorithm. Sample sizes n^* and n^{**} represent the transition points that two algorithms reverse the orders in EPR comparisons.

Theorem 3.3 (EPR of Shared Representation Algorithm):

With Assumption 2.1, the expected population risk as defined in (12) is

$$R_{\text{SR}} = \underbrace{kV_{\text{SR}}}_{\text{Bias}} + \underbrace{\frac{kd|\mathcal{Y}| - \sum_{i=1}^k \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2}{n}}_{\text{Variance}} + \frac{|\mathcal{X}||\mathcal{Y}|}{n} \cdot O(1), \quad (20)$$

where

$$V_{\text{SR}} \triangleq \frac{1}{k} \sum_{i=1}^k \left\| \mathbf{B} \left(P_{XY}^{(i)} \right) - \Psi(\mathbf{g}_{\text{SR}i}) \Phi^T(\mathbf{f}_{\text{SR}}) \right\|_{\text{F}}^2. \quad (21)$$

Theorem 3.3 states that the Shared Representation Algorithm has the EPR with both a bias term and a variance term. Specifically, the bias term is smaller than Algorithm (ii), and the variance term is smaller than Algorithm (i), which lead to the results in Section III-A.4.

4) Comparisons between the EPRs of Classical Algorithms:

Based on the results above, we can have a summary of the mentioned EPRs, as shown in Figure 2. When neglecting the $O(1)$ terms, the expressions can be seen as functions of the sample size n . In view of the EPRs as linear functions of $1/n$, we set the abscissa as the reciprocal of the sample size.

Specifically, when the sample size n is large (as the left part of Figure 2), the data of the device allow training a well-performed model, where the Personalized Model Algorithm is preferable. When the sample size n is small (as the right part of Figure 2), the data of each device lead to a model with high error, where therefore the Global Model Algorithm is preferable. Moreover, under the following sample size range, the Shared Representation Algorithm can outperform the other two algorithms.

Corollary 3.4: When $n^{**} < n < n^*$, $R_{\text{SR}} < R_{\text{PM}}$ and $R_{\text{SR}} < R_{\text{GM}}$ are both satisfied, i.e., the Shared Representation Algorithm can perform best among the classical algorithms,

where

$$n^* = \frac{k|\mathcal{X}||\mathcal{Y}| - k - kd|\mathcal{Y}| + \sum_{i=1}^k (C_i + \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2)}{kV_{\text{SR}}} > \frac{(1 + C_{\min})|\mathcal{X}||\mathcal{Y}| - d|\mathcal{Y}|}{V_{\text{SR}}} \quad (22)$$

and

$$n^{**} = \frac{kd|\mathcal{Y}| - \sum_{i=1}^k \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2}{kV_{\text{GM}} - kV_{\text{SR}}} < \frac{d|\mathcal{Y}|}{V_{\text{GM}} - V_{\text{SR}}}. \quad (23)$$

Based on Corollary 3.4, we have the following interpretations of the transition points n^* and n^{**} . The transition point n^* is proportional to the *space dimensionality* $|\mathcal{X}||\mathcal{Y}|$ (considering the practical scene that $|\mathcal{X}|$ is much larger than d), and inversely proportional to the averaged distance between the underlying distributions and the learned models V_{SR} . It indicates that a high-dimensional sample space and well-adapted model can lead to the preference of the Shared Representation Algorithm compared with the Personalized Model Algorithm.

On the other hand, the transition point n^{**} is proportional to $d|\mathcal{Y}|$, where it can be recognized as the *model dimensionality* of the classifiers with $d|\mathcal{Y}|$ training parameters. Meanwhile, the transition point is inversely proportional to $V_{\text{GM}} - V_{\text{SR}}$. When $\mathbf{f}_{\text{SR}} = \mathbf{f}_{\text{GM}}$, the numerator becomes $V_{\text{GM}} - V_{\text{SR}} = \text{Var}(\{\mathbf{g}_{\text{SR}i}\}) \triangleq \frac{1}{k} \sum_{i=1}^k \|\Psi(\mathbf{g}_{\text{SR}i}) - \Psi(\frac{1}{k} \sum_{i=1}^k \mathbf{g}_{\text{SR}i})\|_{\text{F}}^2$, which can be seen as the variance of the learned classifiers and quantifies the *task similarity*. It explains that the low-dimensional model and the large task similarities lead to the preference of the Shared Representation Algorithm, in comparison with the Global Model Algorithm.

B. Classifiers Collaboration Algorithm

As introduced in Section II, the Classifiers Collaboration Algorithm (Algorithm (iv)) jointly employs the classifiers

learned by different devices. We prove in the following that its EPR can be declined through selecting proper combining coefficient α_i as defined in (13). As this algorithm is an operation based on Algorithm (iii), we do not provide the comparisons with Algorithm (i) or (ii) in this part.

First, for each device i , the EPR of the Shared Representation Algorithm as defined in (12) is

$$R_{\text{SR}i} \triangleq R_i(\hat{\mathbf{f}}_{\text{SR}}, \hat{\mathbf{g}}_{\text{SR}i}) = \frac{d|\mathcal{Y}|}{n} - \frac{1}{n} \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2 + \|\Psi_i^{\text{T}} \Phi_i - \Psi^{\text{T}}(\mathbf{g}_{\text{SR}i}) \Phi(\mathbf{f}_{\text{SR}})\|_{\text{F}}^2 + O\left(\frac{1}{nk}\right). \quad (24)$$

When α_i is a constant, we can compute the EPR (5) of each device as a function of α_i

$$\begin{aligned} \tilde{R}_i(\alpha_i) &\triangleq R_i(\hat{\mathbf{f}}_{\text{SR}}, \alpha_i \hat{\mathbf{g}}_{\text{SR}i} + (1 - \alpha_i) \hat{\mathbf{g}}_{\text{SR}}) \\ &= \alpha_i^2 \left(\frac{d|\mathcal{Y}|}{n} - \frac{1}{n} \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2 \right) \\ &\quad + (1 - \alpha_i)^2 \|\Psi(\mathbf{g}_{\text{SR}i}) - \Psi(\mathbf{g}_{\text{SR}})\|_{\text{F}}^2 \\ &\quad + \|\Psi_i^{\text{T}} \Phi_i - \Psi^{\text{T}}(\mathbf{g}_{\text{SR}i}) \Phi(\mathbf{f}_{\text{SR}})\|_{\text{F}}^2 + O\left(\frac{1}{nk}\right), \end{aligned} \quad (25)$$

where $\mathbf{g}_{\text{SR}} = \frac{1}{k} \sum_{i=1}^k \mathbf{g}_{\text{SR}i}$ and $\tilde{R}_i(1) = R_{\text{SR}i}$. Thus, the optimal α_i^* is

$$\begin{aligned} \alpha_i^* &= \arg \min_{\alpha_i} \tilde{R}_i(\alpha_i) \\ &= \frac{\|\Psi(\mathbf{g}_{\text{SR}i}) - \Psi(\mathbf{g}_{\text{SR}})\|_{\text{F}}^2}{\|\Psi(\mathbf{g}_{\text{SR}i}) - \Psi(\mathbf{g}_{\text{SR}})\|_{\text{F}}^2 + \frac{d|\mathcal{Y}|}{n} - \frac{1}{n} \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2}. \end{aligned} \quad (26)$$

Then we have the following characterization of the optimal coefficient α_i^* .

Proposition 3.5: The EPR as defined in (25) with respect to the parameter α_i^* as defined in (26) is

$$\begin{aligned} \tilde{R}_i(\alpha_i^*) &= \frac{\left(\frac{d|\mathcal{Y}|}{n} - \frac{1}{n} \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2 \right) \|\Psi(\mathbf{g}_{\text{SR}i}) - \Psi(\mathbf{g}_{\text{SR}})\|_{\text{F}}^2}{\frac{d|\mathcal{Y}|}{n} - \frac{1}{n} \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2 + \|\Psi(\mathbf{g}_{\text{SR}i}) - \Psi(\mathbf{g}_{\text{SR}})\|_{\text{F}}^2} \\ &\quad + \|\Psi_i^{\text{T}} \Phi_i - \Psi^{\text{T}}(\mathbf{g}_{\text{SR}i}) \Phi(\mathbf{f}_{\text{SR}})\|_{\text{F}}^2 + O\left(\frac{1}{nk}\right) < R_{\text{SR}i}. \end{aligned} \quad (27)$$

Therefore, the Classifiers Collaboration Algorithm can provide a better performance than the Shared Representation Algorithm by selecting proper coefficient. In practice, the coefficient α_i^* also needs estimating. As the simplest way, we substitute the empirical classifiers for the true classifiers, i.e., we select

$$\alpha_i = \frac{\|\Psi(\hat{\mathbf{g}}_{\text{SR}i}) - \Psi(\hat{\mathbf{g}}_{\text{SR}})\|_{\text{F}}^2}{\|\Psi(\hat{\mathbf{g}}_{\text{SR}i}) - \Psi(\hat{\mathbf{g}}_{\text{SR}})\|_{\text{F}}^2 + \frac{d|\mathcal{Y}|}{n} - \frac{1}{n} \|\Psi(\hat{\mathbf{g}}_{\text{SR}i})\|_{\text{F}}^2}. \quad (28)$$

However, due to the complicated expression of the undecided parameter α_i , the expectation of the population risk (14) is difficult to derive. Specifically, we consider the large sample size regime:

Assumption 3.6: The sample size n of each task satisfies $k^{\frac{1}{2}} > n \gg 1$.

As introduced in Section III-A.4, the right bound of the sample size interval that Algorithm (iii) performs best is proportional to the space dimensionality, where such assumption is often satisfied. Then we have the following characterization of $R_{\text{CC}i}$ as defined in (14).

Theorem 3.7 (EPR of Shared Representation Algorithm): With Assumption 2.1 and 3.6, the expected population risk as defined in (14) satisfies

$$R_{\text{SR}i} - R_{\text{CC}i} = (1 - M_i) \frac{(d|\mathcal{Y}| - \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2)^2}{n^2 \|\Psi(\mathbf{g}_{\text{SR}i}) - \Psi(\mathbf{g}_{\text{SR}})\|_{\text{F}}^2} + O(n^{-3}), \quad (29)$$

where M_i is defined in (48) of Appendix E.

The right-hand side of (29) decides whether the Classifiers Collaboration Algorithm is more preferable than the Shared Representation Algorithm. It is estimated in Appendix E that $M_i \lesssim \frac{4}{d-1}$. Thus, when $d > 5$, we have $R_{\text{SR}i} > R_{\text{CC}i}$, which means in the most practical scenes, the Classifiers Collaboration Algorithm can perform better than the Shared Representation Algorithm.

As a brief summary of this section, we present the EPRs of all the four algorithms and the assumptions of k and n in Table I.

IV. EPR ANALYSES: THE CONTINUOUS CASE

Before this section, all the analyses are based on the sample space with a limited cardinality $|\mathcal{X}|$. However, for practical applications, such an assumption could lead to problems. Therefore, we extend the analyses above to the continuous case. The changes of our theoretical analyses mainly cover: (i) the representations of the feature extractors and classifiers are defined with parameterized models, (ii) the EPRs as functions of CDMs are transformed into integrals, and (iii) for analyses, the EPRs are computed under the large sample size regime.

For point (i), suppose that $\mathbf{w} = (\mathbf{w}_f, \mathbf{w}_g)$ represents all the parameters to construct the feature extractor and classifier, where \mathbf{w}_f and \mathbf{w}_g represent the parameters contained in the feature extractor and classifier, respectively. In detail, the feature representation $\mathbf{f}(\mathbf{w}_f)$ expresses the function of data x with values $\mathbf{f}(x; \mathbf{w}_f)$, and similarly $\mathbf{g}(\mathbf{w}_g)$ expresses the function of label y with values $\mathbf{g}(y; \mathbf{w}_g)$. Specifically, the parameter vector \mathbf{w} contains D parameters, i.e., $\mathbf{w} \in \mathbb{R}^D$.

For change (ii), the risk function with respect to CDMs can be computed in following way:

$$\begin{aligned} &\left\| \mathbf{B}(\hat{P}_{XY}^{(i)}) - \Psi(\mathbf{g}(\mathbf{w}_g)) \Phi^{\text{T}}(\mathbf{f}(\mathbf{w}_f)) \right\|_{\text{F}}^2 \\ &= \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} \frac{\left(\hat{P}_{XY}^{(i)}(x, y) - Q_{XY}^{(\mathbf{w}_f, \mathbf{w}_g)}(x, y) \right)^2}{P_X(x) P_Y(y)} dx, \end{aligned} \quad (30)$$

where

$$Q_{XY}^{(\mathbf{w}_f, \mathbf{w}_g)}(x, y) \triangleq P_X(x) P_Y(y) (1 + \mathbf{f}^{\text{T}}(x; \mathbf{w}_f) \mathbf{g}(y; \mathbf{w}_g)).$$

For change (iii), we employ Assumption 3.6. Under these changes, the difference of our main results lies in the EPR of the Personalized Model Algorithm. The parameters achieved

TABLE I
SUMMARY OF THE EPRS OF THE MENTIONED ALGORITHMS

EPR	Task number k	Sample size n	Expression
R_{PM}	Large	No request	$\frac{k \mathcal{X} \mathcal{Y} - k + \sum_{i=1}^k C_i}{n}$
R_{GM}	Large	No request	$kV_{\text{GM}} + \frac{ \mathcal{X} \mathcal{Y} }{n} \cdot O(1)$
R_{SR}	Large	No request	$kV_{\text{SR}} + \frac{kd \mathcal{Y} - \sum_{i=1}^k \ \Psi(\mathbf{g}_{\text{SR}i})\ _{\text{F}}^2}{n} + \frac{ \mathcal{X} \mathcal{Y} }{n} \cdot O(1)$
R_{CC}	Large	$k^{\frac{1}{2}} > n \gg 1$	$kV_{\text{SR}} + \frac{kd \mathcal{Y} - \sum_{i=1}^k \ \Psi(\mathbf{g}_{\text{SR}i})\ _{\text{F}}^2}{n}$ $- \sum_{i=1}^k (1 - M_i) \frac{(d \mathcal{Y} - \ \Psi(\mathbf{g}_{\text{SR}i})\ _{\text{F}}^2)^2}{n^2 \ \Psi(\mathbf{g}_{\text{SR}i}) - \Psi(\mathbf{g}_{\text{SR}})\ _{\text{F}}^2} + O\left(\frac{k}{n^3}\right)$

TABLE II
SUMMARY OF THE EPRS OF THE MENTIONED ALGORITHMS UNDER THE CONTINUOUS DATA SPACE

EPR	Task number k	Sample size n	Expression
R_{PM}	Large	$k^{\frac{1}{2}} > n \gg 1$	$\sum_{i=1}^k \frac{D_i}{n} + O\left(\frac{k}{n^2}\right)$
R_{GM}	Large	$k^{\frac{1}{2}} > n \gg 1$	$kV_{\text{GM}} + O\left(\frac{1}{n}\right)$
R_{SR}	Large	$k^{\frac{1}{2}} > n \gg 1$	$kV_{\text{SR}} + \frac{kd \mathcal{Y} - \sum_{i=1}^k \ \Psi(\mathbf{g}_{\text{SR}i})\ _{\text{F}}^2}{n} + O\left(\frac{1}{n}\right)$
R_{CC}	Large	$k^{\frac{1}{2}} > n \gg 1$	$kV_{\text{SR}} + \frac{kd \mathcal{Y} - \sum_{i=1}^k \ \Psi(\mathbf{g}_{\text{SR}i})\ _{\text{F}}^2}{n}$ $- \sum_{i=1}^k (1 - M_i) \frac{(d \mathcal{Y} - \ \Psi(\mathbf{g}_{\text{SR}i})\ _{\text{F}}^2)^2}{n^2 \ \Psi(\mathbf{g}_{\text{SR}i}) - \Psi(\mathbf{g}_{\text{SR}})\ _{\text{F}}^2} + O\left(\frac{k}{n^3}\right)$

can be defined with the feature and classifiers as in (6) such that

$$\hat{\mathbf{f}}_{\text{PM}i} = \mathbf{f}(\hat{\mathbf{w}}_{\text{fPM}i}), \text{ and } \hat{\mathbf{g}}_{\text{PM}i} = \mathbf{g}(\hat{\mathbf{w}}_{\text{gPM}i}). \quad (31)$$

We define the risk

$$L_i(\mathbf{w}) \triangleq \left\| \mathbf{B}(P_{XY}^{(i)}) - \Phi(\mathbf{g}(\mathbf{w}_g))\Psi^T(\mathbf{f}(\mathbf{w}_f)) \right\|_{\text{F}}^2.$$

Then, we have the following characterization of its EPR.

Theorem 4.1: With Assumption 2.1 and 3.6, suppose that there exist parameters $\mathbf{w}_{\text{PM}i} = (\mathbf{w}_{\text{fPM}i}, \mathbf{w}_{\text{gPM}i})$ such that $\mathbf{f}_i^T(x)\mathbf{g}_i(y) = \mathbf{f}^T(x; \mathbf{w}_{\text{fPM}i})\mathbf{g}(y; \mathbf{w}_{\text{gPM}i})$, for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$ and let $D_i \triangleq \text{rank}\left(\nabla_{\mathbf{w}}^2 L_i(\mathbf{w})|_{\mathbf{w}=\mathbf{w}_{\text{PM}i}}\right)$. The expected population risk as defined in (7) is

$$R_{\text{PM}} = \sum_{i=1}^k \frac{D_i}{n} + O\left(\frac{k}{n^2}\right) \leq \frac{kD}{n}. \quad (32)$$

The other ERPs of Algorithm (ii)-(iv) need slight modifications with respect to Assumption 3.6. For simplicity, we summarize all the results in Table II.

The transition point as defined in (22) can be estimated as $n^* \approx \frac{D}{V_{\text{SR}}}$, and the expression of n^{**} does not change. Similar interpretations can be given as in Section III-A.4, where n^* indicates that more parameters contained in the

feature extractors can lead to the preference of the Shared Representation Algorithm compared with the Personalized Model Algorithm.

V. CONCLUSION

This paper introduces a mathematical framework for federate learning analyses under the large task number regime. We provide the closed-form EPR expressions of typical algorithms and characterize the roles of the task similarity and the model dimensionality, which explain the sample size range for employing the Shared Representation Algorithm. In addition, we propose an improved algorithm by collaborating the classifiers learned in different devices, whose effectiveness is verified by the EPR analyses.

APPENDIX

A. Proof of Theorem 3.1

According to the definition (7) and the fact $\Psi(\hat{\mathbf{g}}_{\text{PM}i})\Phi^T(\hat{\mathbf{f}}_{\text{PM}i}) = \mathbf{B}(\hat{P}_{XY}^{(i)})$, the EPR can be computed as

$$R_{\text{PM}} = \sum_{i=1}^k \mathbb{E} \left[\left\| \mathbf{B}\left(P_{XY}^{(i)}\right), \mathbf{B}\left(\hat{P}_{XY}^{(i)}\right) \right\|_{\text{F}}^2 \right]$$

$$\begin{aligned}
&= \sum_{i=1}^k \mathbb{E} \left[\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left(P_{XY}^{(i)}(x, y) - \hat{P}_{XY}^{(i)}(x, y) \right)^2}{P_X(x) P_Y(y)} \right] \\
&= \sum_{i=1}^k \frac{1}{n} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{P_{XY}^{(i)}(x, y) - [P_{XY}^{(i)}(x, y)]^2}{P_X(x) P_Y(y)} \\
&= \sum_{i=1}^k \frac{1}{n} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \left[(1 + \mathbf{f}_i^T(x) \mathbf{g}_i(y)) \right. \\
&\quad \left. - P_X(x) P_Y(y) (1 + \mathbf{f}_i^T(x) \mathbf{g}_i(y))^2 \right]. \quad (33)
\end{aligned}$$

With $\mathbb{E}_{P_X}[\mathbf{f}_i(X)] = \mathbb{E}_{P_Y}[\mathbf{g}_i(Y)] = \mathbf{0}$ and $\Phi_i^T \Phi_i = \mathbf{I}_d$, we can derive the theorem.

B. Proof of Theorem 3.2

According to the definition (9) and the fact $\Psi(\hat{\mathbf{g}}_{\text{GM}}) \Phi^T(\hat{\mathbf{f}}_{\text{GM}}) = \mathbf{B} \left(\frac{1}{k} \sum_{i=1}^k \hat{P}_{XY}^{(i)} \right)$, the EPR can be computed as

$$\begin{aligned}
R_{\text{GM}} &= \sum_{i=1}^k \mathbb{E} \left[\left\| \mathbf{B} \left(P_{XY}^{(i)} \right), \mathbf{B} \left(\frac{1}{k} \sum_{i=1}^k \hat{P}_{XY}^{(i)} \right) \right\|_{\text{F}}^2 \right] \\
&= \sum_{i=1}^k \mathbb{E} \left[\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left(P_{XY}^{(i)}(x, y) - \frac{1}{k} \sum_{j=1}^k \hat{P}_{XY}^{(j)}(x, y) \right)^2}{P_X(x) P_Y(y)} \right] \\
&= \sum_{i=1}^k \mathbb{E} \left[\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left(P_{XY}^{(i)}(x, y) - \frac{1}{k} \sum_{j=1}^k P_{XY}^{(j)}(x, y) \right)^2}{P_X(x) P_Y(y)} \right] \\
&\quad + \sum_{i=1}^k \mathbb{E} \left[\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left(\sum_{j=1}^k \hat{P}_{XY}^{(j)}(x, y) - \sum_{j=1}^k P_{XY}^{(j)}(x, y) \right)^2}{k^2 P_X(x) P_Y(y)} \right], \quad (34)
\end{aligned}$$

where the last equation comes from the fact that $\mathbb{E}[\hat{P}_{XY}^{(j)}] = P_{XY}^{(j)}$. With the result in Section A, we have

$$\begin{aligned}
&\sum_{i=1}^k \mathbb{E} \left[\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left(\sum_{j=1}^k \hat{P}_{XY}^{(j)}(x, y) - \sum_{j=1}^k P_{XY}^{(j)}(x, y) \right)^2}{k^2 P_X(x) P_Y(y)} \right] \\
&= \frac{|\mathcal{X}| |\mathcal{Y}| \cdot O(1)}{n}.
\end{aligned}$$

Then, Theorem 3.2 is proved.

C. Proof of Theorem 3.3

Note that

$$\hat{\mathbf{g}}_{\text{SR}i} = \arg \min_{\mathbf{g}} \left\| \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) - \Psi(\mathbf{g}) \Phi^T(\hat{\mathbf{f}}_{\text{SR}}) \right\|_{\text{F}}^2. \quad (35)$$

We can compute by the derivatives that

$$\Psi(\hat{\mathbf{g}}_{\text{SR}i}) = \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) \Phi(\hat{\mathbf{f}}_{\text{SR}}). \quad (36)$$

With

$$\hat{\mathbf{f}}_{\text{SR}} = \arg \min_{\mathbf{f}: \Phi^T(\mathbf{f}) \Phi(\mathbf{f}) = \mathbf{I}_d} \sum_{i=1}^k \left\| \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) - \Psi(\hat{\mathbf{g}}_{\text{SR}i}) \Phi^T(\mathbf{f}) \right\|_{\text{F}}^2, \quad (37)$$

we have

$$\begin{aligned}
\hat{\mathbf{f}}_{\text{SR}} &= \arg \max_{\mathbf{f}: \Phi^T(\mathbf{f}) \Phi(\mathbf{f}) = \mathbf{I}_d} \\
&\Phi^T(\mathbf{f}) \left[\frac{1}{k} \sum_{i=1}^k \mathbf{B}^T \left(\hat{P}_{XY}^{(i)} \right) \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) \right] \Phi(\mathbf{f}).
\end{aligned}$$

Note that (19) defines

$$\begin{aligned}
\mathbf{f}_{\text{SR}} &\triangleq \arg \max_{\mathbf{f}: \Phi^T(\mathbf{f}) \Phi(\mathbf{f}) = \mathbf{I}_d} \\
&\Phi^T(\mathbf{f}) \left[\frac{1}{k} \sum_{i=1}^k \mathbf{B}^T \left(P_{XY}^{(i)} \right) \mathbf{B} \left(P_{XY}^{(i)} \right) \right] \Phi(\mathbf{f}).
\end{aligned}$$

For simplicity, let $\mathbf{A} \triangleq \frac{1}{k} \sum_{i=1}^k \mathbf{B}^T \left(P_{XY}^{(i)} \right) \mathbf{B} \left(P_{XY}^{(i)} \right)$ and $\hat{\mathbf{A}} \triangleq \frac{1}{k} \sum_{i=1}^k \mathbf{B}^T \left(\hat{P}_{XY}^{(i)} \right) \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right)$.

(i) When $\mathbf{A} = \hat{\mathbf{A}}$, $\Phi^T(\mathbf{f}_{\text{SR}}) = \Phi^T(\hat{\mathbf{f}}_{\text{SR}})$.

(ii) When $\mathbf{A} \neq \hat{\mathbf{A}}$, according to Davis-Kahan's Theorem [21], we have

$$\begin{aligned}
\frac{1}{2} \left\| \Phi(\mathbf{f}_{\text{SR}}) - \Phi(\hat{\mathbf{f}}_{\text{SR}}) \right\|_{\text{F}}^2 &\leq \left\| \sin \Theta \left(\Phi(\mathbf{f}_{\text{SR}}), \Phi(\hat{\mathbf{f}}_{\text{SR}}) \right) \right\|_{\text{F}}^2 \\
&\leq \frac{\left\| \mathbf{A} - \hat{\mathbf{A}} \right\|_{\text{F}}^2}{\delta}, \quad (38)
\end{aligned}$$

where the eigenvalues of \mathbf{A} and $\hat{\mathbf{A}}$ are $\lambda_1 \geq \dots \geq \lambda_{|\mathcal{X}|}$ and $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_{|\mathcal{X}|}$, and

$$\delta \triangleq \inf_{\mathbf{A}} \inf_{1 \leq i \leq d, d \leq j \leq |\mathcal{X}|} |\lambda_i - \hat{\lambda}_j|. \quad (39)$$

It is easy to verify that

$$\mathbb{E} \left[\left\| \mathbf{A} - \hat{\mathbf{A}} \right\|_{\text{F}}^2 \right] = \frac{|\mathcal{X}| |\mathcal{Y}|}{n} \cdot O(k^{-1}). \quad (40)$$

Considering the cases (i) and (ii), we have

$$\mathbb{E} \left[\left\| \Phi(\mathbf{f}_{\text{SR}}) - \Phi(\hat{\mathbf{f}}_{\text{SR}}) \right\|_{\text{F}}^2 \right] = \frac{|\mathcal{X}| |\mathcal{Y}|}{n} \cdot O(k^{-1}). \quad (41)$$

Now, we can compute the EPR (12):

$$\begin{aligned}
R_{\text{SR}} &= \sum_{i=1}^k \mathbb{E} \left[\left\| \mathbf{B} \left(P_{XY}^{(i)} \right) - \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) \Phi(\hat{\mathbf{f}}_{\text{SR}}) \Phi^T(\hat{\mathbf{f}}_{\text{SR}}) \right\|_{\text{F}}^2 \right] \\
&= \sum_{i=1}^k \mathbb{E} \left[\left\| \mathbf{B} \left(P_{XY}^{(i)} \right) - \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) \Phi(\mathbf{f}_{\text{SR}}) \Phi^T(\mathbf{f}_{\text{SR}}) \right\|_{\text{F}}^2 \right] \\
&\quad + \frac{|\mathcal{X}||\mathcal{Y}|}{n} \cdot O(1) \\
&= \sum_{i=1}^k \mathbb{E} \left[\left\| \mathbf{B} \left(P_{XY}^{(i)} \right) - \mathbf{B} \left(P_{XY}^{(i)} \right) \Phi(\mathbf{f}_{\text{SR}}) \Phi^T(\mathbf{f}_{\text{SR}}) \right\|_{\text{F}}^2 \right] \\
&\quad + \sum_{i=1}^k \mathbb{E} \left[\left\| \mathbf{B} \left(P_{XY}^{(i)} - \hat{P}_{XY}^{(i)} \right) \Phi(\mathbf{f}_{\text{SR}}) \Phi^T(\mathbf{f}_{\text{SR}}) \right\|_{\text{F}}^2 \right] \\
&\quad + \frac{|\mathcal{X}||\mathcal{Y}|}{n} \cdot O(1) \\
&= kV_{\text{SR}} + \sum_{i=1}^k \frac{1}{n} \sum_{y=1}^{|\mathcal{Y}|} \text{tr} \left(\mathbb{E}_{P_{X|Y}^{(i)}} [\mathbf{f}_{\text{SR}}(X) \mathbf{f}_{\text{SR}}^T(X)] \right) \\
&\quad - \sum_{i=1}^k \frac{1}{n} \sum_{y=1}^{|\mathcal{Y}|} P_Y(y) \left\| \mathbb{E}_{P_{X|Y}^{(i)}} [\mathbf{f}_{\text{SR}}(X)] \right\|^2 \\
&\quad + \frac{|\mathcal{X}||\mathcal{Y}|}{n} \cdot O(1) \\
&= kV_{\text{SR}} + \frac{kd|\mathcal{Y}| - \sum_{i=1}^k \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2}{n} + \frac{|\mathcal{X}||\mathcal{Y}|}{n} \cdot O(1). \tag{42}
\end{aligned}$$

Then, Theorem 3.3 is proved.

D. Local Approximation of the multinomial distribution

According to the multinomial distribution,

$$\mathbb{P} \left(\hat{P}_{XY}^{(i)}; P_{XY}^{(i)} \right) \doteq \exp \left(-nD \left(\hat{P}_{XY}^{(i)} \| P_{XY}^{(i)} \right) \right). \tag{43}$$

Moreover, we introduce the following lemma:

Lemma 5.1: Suppose that $|P(z) - Q(z)| < \epsilon$ for each $z \in \mathcal{Z}$, where $\epsilon/|\mathcal{Z}| \ll 1$, K-L divergence between $P(z)$ and $Q(z)$ can have the following approximation

$$D(P\|Q) = \frac{1}{2} \chi^2(P, Q) + O(\epsilon^3). \tag{44}$$

With (43) and Lemma 5.1, we let

$$\mathbb{P} \left(\hat{P}_{XY}^{(i)}; P_{XY}^{(i)} \right) \propto \exp \left(-\frac{n}{2} \chi^2 \left(P_{XY}^{(i)}, \hat{P}_{XY}^{(i)} \right) \right). \tag{45}$$

E. Proof of Theorem 3.7

First, we prove the EPR with constant α_i as define in (25).

$$\begin{aligned}
\tilde{R}_i(\alpha_i) &\triangleq R_i \left(\hat{\mathbf{f}}_{\text{SR}}, \alpha_i \hat{\mathbf{g}}_{\text{SR}i} + (1 - \alpha_i) \hat{\mathbf{g}}_{\text{SR}} \right) \\
&= \mathbb{E} \left[\left\| \mathbf{B} \left(P_{XY}^{(i)} \right) - \left(\alpha_i \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) \Phi(\hat{\mathbf{f}}_{\text{SR}}) \right. \right. \right. \\
&\quad \left. \left. + (1 - \alpha_i) \mathbf{B} \left(\frac{1}{k} \sum_{j=1}^k \hat{P}_{XY}^{(j)} \right) \Phi(\hat{\mathbf{f}}_{\text{SR}}) \right) \Phi^T(\hat{\mathbf{f}}_{\text{SR}}) \right\|_{\text{F}}^2 \right] \\
&= \mathbb{E} \left[\left\| \mathbf{B} \left(P_{XY}^{(i)} \right) - \left(\alpha_i \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) \Phi(\hat{\mathbf{f}}_{\text{SR}}) \right. \right. \right. \\
&\quad \left. \left. + (1 - \alpha_i) \mathbf{B} \left(\frac{1}{k} \sum_{j=1}^k P_{XY}^{(j)} \right) \Phi(\mathbf{f}_{\text{SR}}) \right) \Phi^T(\mathbf{f}_{\text{SR}}) \right\|_{\text{F}}^2 \right] + O \left(\frac{1}{nk} \right) \\
&= \alpha_i^2 \mathbb{E} \left[\left\| \mathbf{B} \left(P_{XY}^{(i)} - \hat{P}_{XY}^{(i)} \right) \Phi(\mathbf{f}_{\text{SR}}) \Phi^T(\mathbf{f}_{\text{SR}}) \right\|_{\text{F}}^2 \right] \\
&\quad + (1 - \alpha_i)^2 \left\| \mathbf{B} \left(P_{XY}^{(i)} - \frac{1}{k} \sum_{j=1}^k P_{XY}^{(j)} \right) \Phi(\mathbf{f}_{\text{SR}}) \Phi^T(\mathbf{f}_{\text{SR}}) \right\|_{\text{F}}^2 \\
&\quad + \left\| \mathbf{B} \left(P_{XY}^{(i)} \right) - \mathbf{B} \left(P_{XY}^{(i)} \right) \Phi(\mathbf{f}_{\text{SR}}) \Phi^T(\mathbf{f}_{\text{SR}}) \right\|_{\text{F}}^2 + O \left(\frac{1}{nk} \right) \\
&= \alpha_i^2 \left(\frac{d|\mathcal{Y}|}{n} - \frac{1}{n} \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2 \right) \\
&\quad + (1 - \alpha_i)^2 \|\Psi(\mathbf{g}_{\text{SR}i}) - \Psi(\mathbf{g}_{\text{SR}})\|_{\text{F}}^2 \\
&\quad + \|\Psi_i^T \Phi_i - \Psi^T(\mathbf{g}_{\text{SR}i}) \Phi(\mathbf{f}_{\text{SR}})\|_{\text{F}}^2 + O \left(\frac{1}{nk} \right). \tag{46}
\end{aligned}$$

Accordingly, we can derive α_i^* . With Assumption 3.6 and $\Delta \mathbf{g}_{\text{SR}i} \triangleq \hat{\mathbf{g}}_{\text{SR}i} - \mathbf{g}_{\text{SR}i}$,

$$\begin{aligned}
\alpha_i &= \frac{\|\Psi(\hat{\mathbf{g}}_{\text{SR}i}) - \Psi(\hat{\mathbf{g}}_{\text{SR}})\|_{\text{F}}^2}{\|\Psi(\hat{\mathbf{g}}_{\text{SR}i}) - \Psi(\hat{\mathbf{g}}_{\text{SR}})\|_{\text{F}}^2 + \frac{d|\mathcal{Y}|}{n} - \frac{1}{n} \|\Psi(\hat{\mathbf{g}}_{\text{SR}i})\|_{\text{F}}^2} \\
&= \alpha_i^* + O(\|\Psi(\Delta \mathbf{g}_{\text{SR}i})\|_{\text{F}}^2) \\
&\quad + 2 \frac{\langle \theta_i \Psi(\mathbf{g}_{\text{SR}i} - \mathbf{g}_{\text{SR}}) + \theta \Psi(\mathbf{g}_{\text{SR}i}), \Psi(\Delta \mathbf{g}_{\text{SR}i}) \rangle_{\text{F}}}{n\theta^2},
\end{aligned}$$

where

$$\theta = \|\Psi(\mathbf{g}_{\text{SR}i} - \mathbf{g}_{\text{SR}})\|_{\text{F}}^2,$$

and

$$\theta_i = d|y| - \|\Psi(\mathbf{g}_{\text{SR}i})\|_{\text{F}}^2$$

Then, we can have Theorem 3.7 as follows

$$\begin{aligned}
R_i \left(\hat{\mathbf{f}}_{\text{SR}}, \hat{\mathbf{g}}_{\text{SR}i} \right) &= \mathbb{E} \left[\left\| \mathbf{B} \left(P_{XY}^{(i)} \right) - \left(\alpha_i \mathbf{B} \left(\hat{P}_{XY}^{(i)} \right) \Phi(\hat{\mathbf{f}}_{\text{SR}}) \right. \right. \right. \\
&\quad \left. \left. + (1 - \alpha_i) \mathbf{B} \left(\frac{1}{k} \sum_{j=1}^k \hat{P}_{XY}^{(j)} \right) \Phi(\hat{\mathbf{f}}_{\text{SR}}) \right) \Phi^T(\hat{\mathbf{f}}_{\text{SR}}) \right\|_{\text{F}}^2 \right] \\
&= \tilde{R}_i(\alpha_i^*) + O \left(\frac{1}{n^3} \right) \\
&\quad + 4 \mathbb{E} \left[\alpha_i^* \frac{\langle \theta_i \Psi(\mathbf{g}_{\text{SR}i} - \mathbf{g}_{\text{SR}}) + \theta \Psi(\mathbf{g}_{\text{SR}i}), \Psi(\Delta \mathbf{g}_{\text{SR}i}) \rangle_{\text{F}}}{n\theta^2} \right]
\end{aligned}$$

$$\begin{aligned} & \left[\Psi(\mathbf{g}_{\text{SR}i} - \mathbf{g}_{\text{SR}}), \Psi(\Delta \mathbf{g}_{\text{SR}i}) \right]_{\text{F}} \\ &= \tilde{R}_i(\alpha_i^*) + M_i \frac{\theta_i^2}{n^2 \theta} + O\left(\frac{1}{n^3}\right), \end{aligned} \quad (47)$$

where

$$\begin{aligned} M_i &= \frac{4}{\theta \theta_i} \left(\sum_{x,y} P_{XY}^{(i)}(x,y) \mathbf{f}_{\text{SR}}^{\text{T}}(x) (\mathbf{h}_1(y) + \mathbf{h}_2(y)) \mathbf{f}_{\text{SR}}^{\text{T}}(x) \mathbf{h}_1(y) \right. \\ &\quad \left. - \sum_y P_Y(y) \mathbf{h}_1^{\text{T}}(y) \mathbf{g}_{\text{SR}i}(y) \cdot \sum_y P_Y(y) \mathbf{h}_2^{\text{T}}(y) \mathbf{g}_{\text{SR}i}(y) \right) \end{aligned} \quad (48)$$

where $\mathbf{h}_1 \triangleq \mathbf{g}_{\text{SR}i} - \mathbf{g}_{\text{SR}}$ and $\mathbf{h}_2 \triangleq \frac{\theta}{\theta_i} \mathbf{g}_{\text{SR}i}$.

Consider that $\mathbf{g}_{\text{SR}i} - \mathbf{g}_{\text{SR}}$ and \mathbf{g}_i can have a positive inner product or negative inner product. We apply that in average $\mathbb{E}[\mathbf{h}_1^{\text{T}}(Y) \mathbf{g}_{\text{SR}i}(Y)] \approx 0$. Note that $\mathbb{E}[\mathbf{f}_{\text{SR}}(X) \mathbf{f}_{\text{SR}}^{\text{T}}(X)] = \mathbf{I}$. Then we can estimate M_i as

$$\begin{aligned} M_i &\approx \frac{4}{\theta \theta_i} \sum_{x,y} P_{XY}^{(i)}(x,y) \mathbf{f}_{\text{SR}}^{\text{T}}(x) (\mathbf{h}_1(y) + \mathbf{h}_2(y)) \mathbf{f}_{\text{SR}}^{\text{T}}(x) \mathbf{h}_1(y) \\ &\leq \frac{4}{\theta \theta_i} |\mathcal{Y}| \sum_y P_Y(y) (\mathbf{h}_1(y) + \mathbf{h}_2(y))^{\text{T}} \mathbf{h}_1(y) \\ &\approx \frac{4|\mathcal{Y}|}{\theta_i} \leq \frac{4|\mathcal{Y}|}{d|\mathcal{Y}| - |\mathcal{Y}|} = \frac{4}{d-1} \end{aligned} \quad (49)$$

Note that

$$\tilde{R}_i(\alpha_i^*) - \tilde{R}_i(\alpha_i^*) = \frac{\theta_i^2}{n^2 \theta} + O\left(\frac{1}{n^3}\right). \quad (50)$$

We can easily derive Theorem 3.7.

F. Proof of Theorem 4.1

Note that the risk function is $L_i(\mathbf{w}) = \|\mathbf{B}(P_{XY}^{(i)}) - \Phi(\mathbf{g}(\mathbf{w}_g)) \Psi^{\text{T}}(\mathbf{f}(\mathbf{w}_f))\|_{\text{F}}^2$. We also define the empirical risk

$$l_i(\mathbf{p}, \mathbf{w}) = \|\mathbf{B}(\mathbf{p}) - \Phi(\mathbf{g}(\mathbf{w}_g)) \Psi^{\text{T}}(\mathbf{f}(\mathbf{w}_f))\|_{\text{F}}^2, \quad (51)$$

and $L_i(\mathbf{w}) = l_i(P_{XY}^{(i)}, \mathbf{w})$. With $\mathbf{w}_{\text{PM}i} = (\mathbf{w}_{\text{fPM}i}, \mathbf{w}_{\text{gPM}i})$ and $\hat{\mathbf{w}}_{\text{PM}i} = (\hat{\mathbf{w}}_{\text{fPM}i}, \hat{\mathbf{w}}_{\text{gPM}i})$

$$\nabla_{\mathbf{w}} l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) = \nabla_{\mathbf{w}} l_i(\hat{P}_{XY}^{(i)}, \hat{\mathbf{w}}_{\text{PM}i}) = 0. \quad (52)$$

By the Taylor expansions of $\nabla_{\mathbf{w}} l_i(\mathbf{p}, \mathbf{w})$, let $\Delta \mathbf{p}_i \triangleq \hat{P}_{XY}^{(i)} - P_{XY}^{(i)} \in \mathbb{R}^{|\mathcal{X}||\mathcal{Y}|}$, $\Delta \mathbf{w}_i \triangleq \hat{\mathbf{w}}_{\text{PM}i} - \mathbf{w}_{\text{PM}i}$, and then we have

$$\begin{aligned} \Delta \mathbf{w}_i &= \left[\nabla_{\mathbf{w}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right]^{-1} \\ &\quad \cdot \left[\nabla_{\mathbf{w}, \mathbf{p}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right] \Delta \mathbf{p}_i + O(\|\Delta \mathbf{p}_i\|^2). \end{aligned} \quad (53)$$

Then, the EPR is

$$\begin{aligned} R_i(\hat{\mathbf{f}}_{\text{PM}i}, \hat{\mathbf{g}}_{\text{PM}i}) &= \mathbb{E}[L_i(\hat{\mathbf{w}}_{\text{PM}i})] \\ &= \frac{1}{2} \mathbb{E} \left[\Delta \mathbf{p}_i^{\text{T}} \left[\nabla_{\mathbf{w}, \mathbf{p}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right]^{\text{T}} \right. \\ &\quad \left[\nabla_{\mathbf{w}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right]^{-1} \left[\nabla_{\mathbf{w}, \mathbf{p}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right] \Delta \mathbf{p}_i \Big] \\ &\quad + \mathbb{E}[O(\|\Delta \mathbf{p}_i\|^3)] \\ &= \frac{1}{2} \left[\nabla_{\mathbf{w}, \mathbf{p}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right] \mathbb{E}[\Delta \mathbf{p}_i \Delta \mathbf{p}_i^{\text{T}}] \\ &\quad \left[\nabla_{\mathbf{w}, \mathbf{p}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right]^{\text{T}} \left[\nabla_{\mathbf{w}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right]^{-1} \\ &\quad + O\left(\frac{1}{n^2}\right) \\ &= \frac{1}{2n} \text{tr} \left(\left[\nabla_{\mathbf{w}, \mathbf{p}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right] \text{diag} \{P_{XY}^{(i)}\} \right. \\ &\quad \left. \left[\nabla_{\mathbf{w}, \mathbf{p}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right]^{\text{T}} \left[\nabla_{\mathbf{w}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right]^{-1} \right) \\ &\quad + O\left(\frac{1}{n^2}\right). \end{aligned} \quad (54)$$

When $l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) = 0$, let $\mathbf{q}(x, y; \mathbf{w}) = \mathbf{f}^{\text{T}}(x; \mathbf{w}_{\text{f}}) \mathbf{g}(y; \mathbf{w}_{\text{gPM}i})$, and then we have

$$\begin{aligned} &\nabla_{\mathbf{w}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \\ &= 2 \mathbb{E}_{P_{XY}^{(i)}} [\nabla_{\mathbf{w}} \mathbf{q}(X, Y; \mathbf{w}_{\text{PM}i}) \nabla_{\mathbf{w}} \mathbf{q}(X, Y; \mathbf{w}_{\text{PM}i})^{\text{T}}], \end{aligned} \quad (55)$$

$$\begin{aligned} &\left[\nabla_{\mathbf{w}, \mathbf{p}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right] \text{diag} \{P_{XY}^{(i)}\} \\ &\left[\nabla_{\mathbf{w}, \mathbf{p}}^2 l_i(P_{XY}^{(i)}, \mathbf{w}_{\text{PM}i}) \right]^{\text{T}} \\ &= 4 \mathbb{E}_{P_{XY}^{(i)}} [\nabla_{\mathbf{w}} \mathbf{q}(X, Y; \mathbf{w}_{\text{PM}i}) \nabla_{\mathbf{w}} \mathbf{q}(X, Y; \mathbf{w}_{\text{PM}i})^{\text{T}}]. \end{aligned} \quad (56)$$

Finally, we have

$$\begin{aligned} R_i(\hat{\mathbf{f}}_{\text{PM}i}, \hat{\mathbf{g}}_{\text{PM}i}) &= \frac{1}{n} \text{rank} \left(\nabla_{\mathbf{w}}^2 L_i(\mathbf{w})|_{\mathbf{w}=\mathbf{w}_{\text{PM}i}} \right) \\ &\quad + O\left(\frac{1}{n^2}\right). \end{aligned} \quad (57)$$

Thus Theorem 4.1 is proved.

REFERENCES

- [1] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [2] P. P. Liang, T. Liu, L. Ziyin, N. B. Allen, R. P. Auerbach, D. Brent, R. Salakhutdinov, and L.-P. Morency, "Think locally, act globally: Federated learning with local and global representations," *arXiv preprint arXiv:2001.01523*, 2020.
- [3] L. Collins, H. Hassani, A. Mokhtari, and S. Shakkottai, "Exploiting shared representations for personalized federated learning," in *Proceedings of the 38th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 139. PMLR, 18–24 Jul 2021, pp. 2089–2099.
- [4] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.

- [5] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [6] R. Caruana, "Multitask learning," *Machine learning*, vol. 28, no. 1, pp. 41–75, 1997.
- [7] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [8] F. Grimberg, M.-A. Hartley, S. P. Karimireddy, and M. Jaggi, "Optimal model averaging: Towards personalized collaborative learning," *arXiv preprint arXiv:2110.12946*, 2021.
- [9] Y. Zhang and Q. Yang, "A survey on multi-task learning," *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–1, 2021.
- [10] Z. Charles and J. Konečný, "Convergence and accuracy trade-offs in federated learning and meta-learning," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 2575–2583.
- [11] W. Feller, *An introduction to probability theory and its applications, vol 2*. John Wiley & Sons, 2008.
- [12] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to algorithms*. MIT press, 2009.
- [13] A. Nilsson, S. Smith, G. Ulm, E. Gustavsson, and M. Jirstrand, "A performance evaluation of federated learning algorithms," in *Proceedings of the Second Workshop on Distributed Infrastructures for Deep Learning*, 2018, pp. 1–8.
- [14] K. J. Rahman, F. Ahmed, N. Akhter, M. Hasan, R. Amin, K. E. Aziz, A. M. Islam, M. S. H. Mukta, and A. N. Islam, "Challenges, applications and design aspects of federated learning: A survey," *IEEE Access*, vol. 9, pp. 124 682–124 700, 2021.
- [15] S.-L. Huang, A. Makur, G. W. Wornell, and L. Zheng, "On universal features for high-dimensional learning and inference," *arXiv preprint arXiv:1911.09105*, 2019.
- [16] R. A. Horn and C. R. Johnson, *Matrix analysis*. Cambridge university press, 2012.
- [17] M. Thomas and A. T. Joy, "Elements of information theory," 2006.
- [18] S.-L. Huang, X. Xu, L. Zheng, and G. W. Wornell, "An information theoretic interpretation to deep neural networks," in *2019 IEEE international symposium on information theory (ISIT)*. IEEE, 2019, pp. 1984–1988.
- [19] I. Csiszár, "The method of types [information theory]," *IEEE Transactions on Information Theory*, vol. 44, no. 6, pp. 2505–2523, 1998.
- [20] I. Csiszár and P. C. Shields, "Information theory and statistics: A tutorial," 2004.
- [21] G. W. Stewart, "Matrix perturbation theory," 1990.