

# Minimax Learning for Distributed Inference

Cheuk Ting Li<sup>1</sup>, *Member, IEEE*, Xiugang Wu<sup>2</sup>, *Member, IEEE*, Ayfer Özgür, *Senior Member, IEEE*,  
and Abbas El Gamal<sup>3</sup>, *Life Fellow, IEEE*

**Abstract**—The classical problem of supervised learning is to infer an accurate estimate of a target variable  $Y$  from a measured variable  $X$  using a set of labeled training samples. Motivated by the increasingly distributed nature of data and decision making, this paper considers a variation of this classical problem in which the inference is distributed between two nodes, e.g., a mobile device and a cloud, with a rate constraint on the communication between them. The mobile device observes  $X$  and sends a description  $M$  of  $X$  to the cloud, which computes an estimate  $\hat{Y}$  of  $Y$ . We follow the recent minimax learning approach to study this inference problem and show that it corresponds to a one-shot minimax noisy lossy source coding problem. We then establish information theoretic bounds on the risk-rate Lagrangian cost, leading to a general method for designing a near-optimal descriptor-estimator pair. A key ingredient in the proof of our result is a refined version of the strong functional representation lemma previously used to establish several one-shot source coding theorems. Our results show that a naive estimate-compress scheme for rate-constrained inference is not optimal in general. When the distribution of  $(X, Y)$  is known and the error is measured by the logarithmic loss, our bounds on the risk-rate Lagrangian cost provide a new one-shot operational interpretation of the information bottleneck. We also demonstrate a way to bound the excess risk of the descriptor-estimator pair obtained by our method.

**Index Terms**—Minimax learning, distributionally robust learning, information bottleneck, one-shot source coding, functional representation.

## I. INTRODUCTION

THE classical problem of supervised learning is to infer an accurate estimate of a target variable  $Y$  from a measured variable  $X$  on the basis of  $n$  labeled training samples

Manuscript received November 25, 2018; revised January 27, 2020; accepted April 28, 2020. Date of publication October 6, 2020; date of current version November 20, 2020. This work was supported in part by a gift from Huawei Technologies, in part by the Center for Science of Information (CSol), in part by the NSF Science and Technology Center under Grant CCF-0939370, and in part by NSF under Grant CCF-1704624. This article was presented in part at the 2018 IEEE International Symposium on Information Theory. (*Corresponding author: Cheuk Ting Li.*)

Cheuk Ting Li was with the Department of Electrical Engineering, Stanford University, Stanford, CA 94305 USA. He is now with the Department of Information Engineering, The Chinese University of Hong Kong, Hong Kong (e-mail: ctli@ie.cuhk.edu.hk).

Xiugang Wu was with the Department of Electrical Engineering, Stanford University, Stanford, CA 94305 USA. He is now with the Department of Electrical and Computer Engineering, University of Delaware, Newark, DE 19716 USA, and also with the Department of Computer and Information Sciences, University of Delaware, Newark, DE 19716 USA (e-mail: xwu@udel.edu).

Ayfer Özgür and Abbas El Gamal are with the Department of Electrical Engineering, Stanford University, Stanford, CA 94305 USA (e-mail: aozgur@stanford.edu; abbas@ee.stanford.edu).

Communicated by N. Kiyavash, Associate Editor for Statistical Learning. Digital Object Identifier 10.1109/TIT.2020.3029182

$\{(X_i, Y_i)\}_{i=1}^n$  independently drawn from an unknown joint distribution  $P$ . The standard approach for solving this problem in statistical learning theory is empirical risk minimization (ERM). For a given set of allowable estimators and a loss function that quantifies the risk of each estimator, ERM chooses the estimator with minimal risk under the empirical distribution of the samples. To avoid overfitting, the set of allowable estimators is restricted to a class with limited complexity.

Recently, an alternative viewpoint has emerged which seeks distributionally robust estimators. Given the labeled training samples, this approach learns an estimator by minimizing its worst-case risk over an ambiguity set centered at the empirical distribution of the samples. In other words, instead of restricting the set of allowable estimators, it aims to avoid overfitting by requiring that the learned estimator performs well under any distribution in a chosen neighborhood of the empirical distribution. This minimax approach has been investigated under different assumptions on how the ambiguity set is constructed, e.g., by restricting the moments [1], forming the  $f$ -divergence balls [2] and Wasserstein balls [3] (see also references therein).

In these previous works, the learning algorithm finds an estimator that acts directly on a fresh (unlabeled) sample  $X$  to predict the corresponding target variable  $Y$ . Often, however the fresh sample  $X$  may be only remotely available, for example, at a mobile node, and when designing the estimator it is desirable to also take into account the cost of communicating  $X$ . This is motivated by the bandwidth and energy limitations on communication in networks or within multiprocessor systems, which often impose significant bottlenecks on the performance of algorithms. There are also an increasing number of applications in which data is generated in a distributed manner and this data (or features of it) are communicated over rate-limited links to a central processor to perform inference. For instance, applications such as Google Goggles and Siri process the locally collected data on clouds. It is thus important to study inference in distributed and rate-constrained settings.

In this paper, which is a more complete version of [4], we study an extension of the classical learning problem where given a finite set of training samples, the learning algorithm needs to infer a descriptor-estimator pair with a desired communication rate in between them. This is especially relevant when both  $X$  and  $Y$  come from a large alphabet or are continuous random variables, e.g. in regression problems, so neither the sample  $X$  nor its predicted value of  $Y$  can be simply communicated in a lossless fashion. We adopt the minimax framework for learning the descriptor-estimator pair.

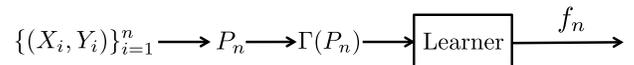
Given a set of labeled training samples, our goal is to find a descriptor-estimator pair by minimizing their resultant worst-case risk over an ambiguity distribution set, where the risk now incorporates both the statistical risk and the communication cost. One of the important conclusions that emerge from the minimax approach to supervised learning in [1] is that the problem of finding the estimator with minimal worst-case risk over an ambiguity set can be broken into two smaller steps: (1) find the worst-case distribution in the ambiguity set that maximizes the (generalized) conditional entropy of  $Y$  given  $X$ , and (2) find the optimal estimator under this worst-case distribution. In this paper, we show that an analogous principle approximately holds for rate-constrained inference. The descriptor-estimator pair with minimal worst-case risk can be found in two steps: (1) find the worst-case distribution in the ambiguity set that maximizes the risk-information Lagrangian cost, and (2) find the optimal descriptor-estimator pair under this worst-case distribution. A key technical ingredient that we use to design the close to optimal descriptor-estimator pair for the worst-case distribution is the strong functional representation lemma (SFRL) used in [5] to establish several one-shot coding theorems, including for lossy source coding and multiple description coding. However, we will need a refined version of this lemma to establish our minimax results.

We demonstrate a way to bound the excess risk of our proposed method, by considering the example in which the ambiguity set is the total variation distance ball around the empirical distribution. We then apply our results to characterize the optimal descriptor-estimator pairs for two applications: rate-constrained linear regression and rate-constrained classification. While a simple scheme in which we first find the optimal estimator ignoring the rate constraint, then compress and communicate the estimator output, is optimal for the linear regression application, we show via the classification application that such an estimate-compress approach is not optimal in general and that when inference is rate-constrained, the optimal descriptor aims to send sufficiently (but not necessarily maximally) informative features of the observed variable, which are at the same time easy to communicate. When applied to the case in which the ambiguity distribution set contains only a single distribution (for example, the true or empirical distribution of  $X, Y$ ) and the loss function for the inference is logarithmic loss, our results provide a new one-shot operational interpretation of the information bottleneck problem.

### Notation

Throughout the paper, we assume that  $\mathcal{X}$ ,  $\mathcal{Y}$  and  $\hat{\mathcal{Y}}$  are Polish spaces, and the loss function  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow \mathbb{R}$  is continuous (note that  $\ell$  is always continuous if  $\mathcal{Y}$  and  $\hat{\mathcal{Y}}$  are discrete). We assume that log is base 2 and the entropy  $H$  is in bits. The length of a variable-length description  $M \in \{0, 1\}^*$  is denoted as  $|M|$ . For random variables  $U, V$ , denote the joint distribution by  $P_{U,V}$  and the conditional distribution of  $U$  given  $V$  by  $P_{U|V}$ . For brevity we denote the distribution of  $(X, Y)$  as  $P$ . We write  $I_P(X; \hat{Y})$  for  $I(X; \hat{Y})$  when  $(X, Y) \sim P$ , and  $P_{\hat{Y}|X}$  is clear from the context.

Learning:



Inference:



Fig. 1. Minimax approach to supervised learning.

## II. PROBLEM FORMULATION

We begin by reviewing the minimax approach to the classical learning problem [1], then extend it to the distributed setting considered in this paper.

### A. Minimax Approach to Supervised Learning

Let  $X \in \mathcal{X}$  and  $Y \in \mathcal{Y}$  be jointly distributed random variables. The problem of statistical learning is to design an accurate estimator of a target variable  $Y$  from a measured variable  $X$  on the basis of a number of independent training samples  $\{(X_i, Y_i)\}_{i=1}^n$  drawn from an unknown joint distribution. The standard approach for solving this problem is to use empirical risk minimization (ERM) in which one defines an admissible class of estimators  $\mathcal{F}$  that consists of functions  $f : \mathcal{X} \rightarrow \hat{\mathcal{Y}}$  (where the reconstruction alphabet  $\hat{\mathcal{Y}}$  can be in general different from  $\mathcal{Y}$ ) and a loss function  $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow \mathbb{R}$ . The risk associated with an estimator  $f$  when the underlying joint distribution of  $X$  and  $Y$  is  $P$  is

$$L(f, P) \triangleq \mathbf{E}_P[\ell(f(X), Y)].$$

ERM simply chooses the estimator  $f_n \in \mathcal{F}$  with minimal risk under the empirical distribution  $P_n$  of the training samples.

Recently, an alternative approach has emerged which seeks distributionally robust estimators. This approach learns an estimator by minimizing its worst-case risk over an ambiguity distribution set  $\Gamma(P_n)$ , i.e.,

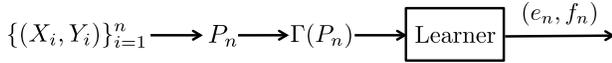
$$f_n = \underset{f}{\operatorname{argmin}} \max_{P \in \Gamma(P_n)} L(f, P), \quad (1)$$

where  $f$  can be any function and  $\Gamma(P_n)$  can be constructed in various ways, e.g., by restricting the moments, forming the  $f$ -divergence balls or Wasserstein balls. While in ERM it is important to restrict the set  $\mathcal{F}$  of admissible estimators to a low-complexity class to prevent overfitting, in the minimax approach overfitting is prevented by explicitly requiring that the chosen estimator is distributionally robust. The learned function  $f_n$  can be then used for predicting  $Y$  when presented with fresh samples of  $X$ . The learning and inference phases are illustrated in Figure 1.

### B. Minimax Learning for Distributed Inference

We extend the minimax learning approach to the setting in which the inference needs to be performed based on a rate-constrained description of  $X$ . In particular, given a set of finite training samples  $\{(X_i, Y_i)\}_{i=1}^n$  independently drawn from an unknown joint distribution  $P$ , our goal is to learn a pair of

Learning:



Inference:

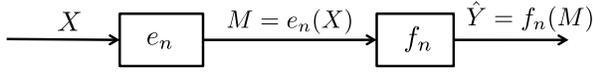


Fig. 2. Minimax learning for distributed inference.

functions  $(e, f)$ , where  $e$  is a descriptor used to compress  $X$  into  $M = e(X) \in \{0, 1\}^*$  (a prefix-free code), and  $f$  is an estimator that takes the compression  $M$  and generates an estimate  $\hat{Y}$  of  $Y$ . See Figure 2.

Let  $R(e, P) \triangleq \mathbf{E}_P[|e(X)|]$  be the rate of the descriptor  $e$  and  $L(e, f, P) \triangleq \mathbf{E}_P[\ell(f(e(X)), Y)]$  be the risk associated with the descriptor-estimator pair  $(e, f)$ , when the underlying distribution of  $(X, Y)$  is  $P$ , and define the risk-rate Lagrangian cost (parametrized by  $\lambda > 0$ ) as

$$L_\lambda(e, f, P) = L(e, f, P) + \lambda R(e, P). \quad (2)$$

Note that this cost function takes into account both the resultant statistical inference risk of  $(e, f)$ , as well as the communication rate they require. The task of a minimax learner is to find an  $(e_n, f_n)$  pair that minimizes the worst-case  $L_\lambda(e, f, P)$  over the ambiguity distribution set  $\Gamma(P_n)$ , i.e.,

$$(e_n, f_n) = \underset{(e, f)}{\operatorname{argmin}} \max_{P \in \Gamma(P_n)} L_\lambda(e, f, P), \quad (3)$$

for an appropriately chosen  $\Gamma(P_n)$  centered at the empirical distribution of samples  $P_n$ . Note that we allow here all possible  $(e, f)$  pairs. We also assume that the descriptor and the estimator can use unlimited common randomness  $W$  which is independent of the data, i.e.,  $e$  and  $f$  can be expressed as functions of  $(X, W)$  and  $(M, W)$ , respectively, and the prefix-free codebook for  $M$  can depend on  $W$ . The availability of such common randomness can be justified by the fact that in practice, although the inference scheme is one-shot, it is used many times (by the same user and by different users), hence the descriptor and the estimator can share a common randomness seed before communication commences without impacting the communication rate.

Note that the main difference between our minimax learning setup and the minimax noisy source coding problem studied in [6] is that here we are considering the one-shot variable-length setting instead of the asymptotic setting. As such, the proof of our result is quite different from that of the asymptotic setting. In [5], a subset of the authors used the *strong functional representation lemma* (SFRL) to derive an upper bound on the (average) rate-distortion function for the one-shot lossy source coding in terms of the rate distortion function for the asymptotic case. While this proof extends easily to the noisy one-shot lossy compression setting (see Theorem 1 and its proof), we need a refined version of SFRL and several other arguments to establish the corresponding result for the minimax setting in Theorem 2.

### III. MAIN RESULTS

We first consider the case in which  $\Gamma$  consists of a single distribution  $P$ , for example, the empirical distribution  $P_n$  as in ERM. Define the minimax risk-rate cost as

$$L_\lambda^*(\Gamma) = \inf_{(e, f)} \sup_{P \in \Gamma} L_\lambda(e, f, P). \quad (4)$$

While it is difficult to minimize the risk-rate cost (2) directly, the minimax risk-rate cost can be bounded in terms of the mutual information between  $X$  and  $\hat{Y}$ .

*Theorem 1:* Let  $\Gamma = \{P\}$ . Then

$$L_\lambda^* \geq \inf_{P_{\hat{Y}|X}} (\mathbf{E}[\ell(\hat{Y}, Y)] + \lambda I(X; \hat{Y})),$$

$$L_\lambda^* \leq \inf_{P_{\hat{Y}|X}} (\mathbf{E}[\ell(\hat{Y}, Y)] + \lambda(I(X; \hat{Y}) + \log(I(X; \hat{Y}) + 1) + 5)).$$

As in other one-shot compression results (e.g., zero-error compression), there is a gap between the upper and lower bound. While the logarithmic gap in Theorem 1 is not as small as the 1-bit gap in the zero-error compression, it is dominated by the linear term  $I(X; \hat{Y})$  when the alphabet of  $X, Y, \hat{Y}$  are very large, for example, if  $Y$  is a large feature set of an image and  $X$  is the image itself.

To prove Theorem 1, we use the strong functional representation (also see [7], [8]).

*Lemma 1* (see [5]): For any random variables  $X, \hat{Y}$ , there exists random variable  $W$  independent of  $X$ , such that  $\hat{Y}$  is a function of  $(X, W)$ , and

$$H(\hat{Y}|W) \leq I(X; \hat{Y}) + \log(I(X; \hat{Y}) + 1) + 4. \quad (5)$$

Here,  $W$  can be intuitively viewed as the part of  $\hat{Y}$  which is not contained in  $X$ . Note that for any  $W$ , such that  $\hat{Y}$  is a function of  $(X, W)$  and  $W$  is independent of  $X$ ,  $H(\hat{Y}|W) \geq I(X; \hat{Y})$ . The statement (5) ensures the existence of such  $W$ , independent of  $X$ , which comes close to this lower bound, and in this sense it is most informative about  $\hat{Y}$ . This is critical for the proof of Theorem 1 as we will see next. Identifying the part of  $\hat{Y}$  which is not contained in  $X$  allows us to generate and share this part between the descriptor and the estimator ahead of time, eliminating the need to communicate it during the course of inference. To find  $W$  and the function that generates  $\hat{Y}$ , we use the Poisson functional representation construction detailed in [5].

*Proof of Theorem 1:* Recall that  $\hat{Y} = f(e(X, W), W)$ . The lower bound follows from the fact that  $I_P(X; \hat{Y}) \leq H_P(M) \leq \mathbf{E}[|M|]$ . To establish the upper bound, fix any  $P_{\hat{Y}|X}$ . Let  $W$  be obtained from (5). Note that  $W$  is independent of  $X$  and can be generated from a random seed shared between the descriptor and the estimator ahead of time. For a given  $w$ , take  $m = e(x, w)$  to be the Huffman codeword of  $\hat{Y}(x, w)$  according to the distribution  $P_{\hat{Y}|W}(\cdot|w)$  (recall that  $\hat{Y}$  is a function of  $(X, W)$ ), and take  $f(m, w)$  to be the decoding function of the Huffman code. The expected codeword length  $\mathbf{E}[|M|] \leq H(\hat{Y}|W) + 1 \leq I(X; \hat{Y}) + \log(I(X; \hat{Y}) + 1) + 5$ .

Taking an infimum over all  $P_{\hat{Y}|X}$  completes the proof. ■

*Remark 1 (Relationship to the Information Bottleneck):* If we consider the logarithmic loss  $\ell(\hat{y}, y) = -\log \hat{y}(y)$ , where

$\hat{y}$  is a distribution over the discrete space  $\mathcal{Y}$ , then the lower bound in Theorem 1 reduces to

$$\inf_{P_{U|X}} (H(Y|U) + \lambda I(X;U)) = H(Y) + \inf_{P_{U|X}} (\lambda I(X;U) - I(Y;U)),$$

which is the information bottleneck (IB) function [9]. Therefore the setting of distributed inference provides an approximate one-shot operational interpretation of the IB (up to a logarithmic gap). In [10], [11] it was shown that the problem of asymptotic lossy source coding with noisy observations also provides an operational interpretation of the IB. Our operational interpretation, however, is more satisfying since the feature extraction problem is by nature one-shot.

The optimal risk-rate cost  $L_\lambda^*(\{P\})$  also shares some similarities with the deterministic information bottleneck (DIB) [12]  $\inf_{P_{T|X}} (\lambda H(T) - I(Y;T))$  (stated in a slightly different form), which, compared to the IB, considers the entropy  $H(T)$  instead of the mutual information  $I(X;U)$ . It is demonstrated in [12] that the DIB produces a compression  $T$  with significantly smaller entropy than the compression  $U$  produced by the IB, making the DIB more suitable for clustering. The rate term  $\mathbf{E}_P[\ell(e(X))]$  in  $L_\lambda^*(\{P\})$  is close to the entropy  $H(e(X))$  (or its conditional entropy given the common randomness  $W$ ). Therefore, the optimal risk-rate cost approach shares the same advantage as the DIB (low entropy of the compression), while having a provable logarithmic gap from the original IB as shown in Theorem 1.

We now extend Theorem 1 to the minimax setting.

**Theorem 2:** Suppose  $\Gamma$  is convex and closed. Then

$$\begin{aligned} L_\lambda^* &\geq \inf_{P_{\hat{Y}|X}} \sup_{P \in \Gamma} (\mathbf{E}_P[\ell(\hat{Y}, Y)] + \lambda I_P(X; \hat{Y})) \\ L_\lambda^* &\leq \inf_{P_{\hat{Y}|X}} \sup_{P \in \Gamma} (\mathbf{E}_P[\ell(\hat{Y}, Y)] + \lambda (I_P(X; \hat{Y}) \\ &\quad + 2 \log(I_P(X; \hat{Y}) + 1) + 6)). \end{aligned}$$

The proof of this theorem uses the following refined version of the strong functional representation lemma which is established in the course of proving Theorem 1 on page 6976 of [5]. We substitute  $P_Y \leftarrow \tilde{P}_{\hat{Y}}$  and  $P_{Y|X} \leftarrow P_{\hat{Y}|X}$  in the Poisson functional representation (note that while the Poisson functional representation in [5] is stated in terms of  $P_{XY}$ , it only depends on  $P_Y$  and  $P_{Y|X}$ ). Also see a similar bound in [7].

**Lemma 2:** For any  $P_{\hat{Y}|X}$  and  $\tilde{P}_{\hat{Y}}$ , there exists random variable  $W$ , and functions  $k(x, w) \in \{1, 2, \dots\}$  and  $\hat{Y}(k, w)$  such that for any  $x$ , we have  $\hat{Y}(k(x, W), W) \sim P_{\hat{Y}|X}(\cdot|x)$ , and

$$\mathbf{E} [\log k(x, W)] \leq D(P_{\hat{Y}|X}(\cdot|x) \parallel \tilde{P}_{\hat{Y}}) + 1.6. \quad (6)$$

The rest of the proof is given in Section V-A.

Theorem 2 suggests that we can simplify the analysis of the risk-rate cost (2)  $L_\lambda = \mathbf{E}_P[\ell(\hat{Y}, Y)] + \lambda \mathbf{E}_P[|M|]$  by replacing the rate  $\mathbf{E}_P[|M|]$  with the mutual information  $I_P(X; \hat{Y})$ . Define the *risk-information cost* as

$$\tilde{L}_\lambda(P_{\hat{Y}|X}, P) = \mathbf{E}_P[\ell(\hat{Y}, Y)] + \lambda I_P(X; \hat{Y}). \quad (7)$$

Theorem 2 implies that the minimax risk-rate cost  $L_\lambda^*$  can be approximated by the *minimax risk-information cost*

$$\tilde{L}_\lambda^*(\Gamma) = \inf_{P_{\hat{Y}|X}} \sup_{P \in \Gamma} \tilde{L}_\lambda(P_{\hat{Y}|X}, P), \quad (8)$$

within a logarithmic gap. Theorem 2 can also be stated in the following slightly weaker form

$$\tilde{L}_\lambda^* \leq L_\lambda^* \leq \tilde{L}_\lambda^* + 2\lambda \log(\lambda^{-1} \tilde{L}_\lambda^* + 1) + 6\lambda.$$

The risk-information cost has more desirable properties than the risk-rate cost. For example, it is convex in  $P_{\hat{Y}|X}$  for fixed  $P$ , and concave in  $P$  for fixed  $P_{\hat{Y}|X}$ . This allows us to exchange the infimum and supremum in Theorem 2 by Sion's minimax Theorem [13], which gives the following proposition.

**Proposition 1:** Suppose  $\mathcal{X}$ ,  $\mathcal{Y}$  and  $\hat{\mathcal{Y}}$  are finite,  $\Gamma$  is convex and closed, and  $\lambda \geq 0$ , then

$$\tilde{L}_\lambda^*(\Gamma) = \inf_{P_{\hat{Y}|X}} \sup_{P \in \Gamma} \tilde{L}_\lambda(P_{\hat{Y}|X}, P) = \sup_{P \in \Gamma} \inf_{P_{\hat{Y}|X}} \tilde{L}_\lambda(P_{\hat{Y}|X}, P).$$

Moreover, there exists  $P_{\hat{Y}|X}^*$  attaining the infimum in the left hand side, which also attains the infimum on the right hand side when  $P$  is fixed to  $P^*$ , the distribution that attains the supremum on the right hand side.

Proposition 1 means that in order to design a robust descriptor-estimator pair that works for any  $P \in \Gamma$ , we only need to design them according to the worst-case distribution  $P^*$  as follows.

**Principle of maximum risk-information cost:** Given a convex and closed  $\Gamma$ , we design the descriptor-estimator pair based on the worst-case distribution

$$P^* = \arg \max_{P \in \Gamma} \inf_{P_{\hat{Y}|X}} \tilde{L}_\lambda(P_{\hat{Y}|X}, P).$$

We then find  $P_{\hat{Y}|X}$  that minimizes  $\tilde{L}_\lambda(P_{\hat{Y}|X}, P^*)$  and design the descriptor-estimator pair accordingly, e.g. using the construction in Theorem 2 according to  $P_{\hat{Y}|X}$ .

We now demonstrate the use of the principle of maximum risk-information cost in a supervised learning setting by considering the case in which  $\Gamma$  is the a total variation distance ball around the empirical distribution. The following proposition bounds the excess risk, which is the gap between  $L_\lambda(e, f, P)$  (the risk-rate cost over the true distribution  $P$ , of the descriptor-estimator pair  $e, f$  obtained using the principle of maximum risk-information cost) and  $L_\lambda^*(\{P\})$  (the optimal risk-rate cost when the true distribution  $P$  is known). The proof is given in Section V-B. The main idea is that the principle of maximum risk-information cost prescribes that the descriptor-estimator pair should be designed according to one of the distributions in  $\Gamma$  (the worst-case distribution). Therefore, if the optimal risk-information cost  $\tilde{L}_\lambda^*(\{\tilde{P}\})$  for different  $\tilde{P} \in \Gamma$  are close to each other, then the optimal risk-information cost of the worst-case distribution will be close to that of the true distribution if the true distribution is in  $\Gamma$ .

**Proposition 2:** Suppose  $\mathcal{X}$ ,  $\mathcal{Y}$  and  $\hat{\mathcal{Y}}$  are finite, and  $\lambda, \epsilon > 0$ . Let  $\alpha = \sup_{y, \hat{y}} \ell(y, \hat{y}) - \inf_{y, \hat{y}} \ell(y, \hat{y})$ . Let  $P_n$  be the empirical distribution of  $n$  i.i.d. samples drawn from  $P$ , and let  $\Gamma = \{\tilde{P} : d_{\text{TV}}(\tilde{P}, P_n) \leq \epsilon\}$ , where  $d_{\text{TV}}(\tilde{P}, P_n) = \sup_{A \subseteq \mathcal{X} \times \mathcal{Y}} |\tilde{P}(A) - P_n(A)|$  is the total variation distance. Let  $P_{\hat{Y}|X}^*$  be attaining

the infimum of  $\inf_{P_{\hat{Y}|X}} \sup_{\tilde{P} \in \Gamma} \tilde{L}_\lambda(P_{\hat{Y}|X}, \tilde{P})$ , and  $e, f$  be constructed as in the proof of Theorem 2 according to  $P_{\hat{Y}|X}^*$ . Then

$$\begin{aligned} & \mathbb{P} \left\{ L_\lambda(e, f, P) > L_\lambda^*({P}) + 2\lambda \log(\lambda^{-1} L_\lambda^*({P})) \right. \\ & \quad \left. + 2\alpha\epsilon + 1 \right\} + 6\lambda + 2\alpha\epsilon \\ & \leq \frac{1}{2\epsilon} \sqrt{\frac{|\mathcal{X}| \cdot |\mathcal{Y}|}{n}}. \end{aligned}$$

#### IV. APPLICATIONS

We present two applications of the minimax results discussed in the previous section. The first application shows that with a proper choice of  $l, \Gamma$ , we obtain a rate-constrained linear regression scheme in which the mobile performs linear regression, then communicates a compressed version of it to the cloud. This straightforward estimate-compress scheme is shown not to be optimal in general via a simple classification example.

##### A. Rate-Constrained Minimax Linear Regression

Suppose  $\mathbf{X} \in \mathbb{R}^d$ ,  $Y \in \mathbb{R}$ ,  $\ell(\hat{y}, y) = (y - \hat{y})^2$  is the mean-squared loss, and we observe the data  $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ . Let  $\boldsymbol{\mu}_{\mathbf{X},n}$ ,  $\mu_{Y,n}$ ,  $\Sigma_{\mathbf{X},n}$ , and  $C_{\mathbf{X}Y,n}$ , respectively, be the empirical means, covariance matrix, and cross covariance matrix estimated from the data. Take  $\Gamma$  to be the set of distributions with these first and second moments, i.e.,

$$\Gamma = \left\{ P_{\mathbf{X}Y} : \boldsymbol{\mu}_{\mathbf{X}} = \boldsymbol{\mu}_{\mathbf{X},n}, \mu_Y = \mu_{Y,n}, \right. \\ \left. \Sigma_{\mathbf{X}} = \Sigma_{\mathbf{X},n}, \sigma_Y^2 = \sigma_{Y,n}^2, C_{\mathbf{X}Y} = C_{\mathbf{X}Y,n} \right\}, \quad (9)$$

The following proposition shows that for this natural choice of the ambiguity set, the distributions  $P^*$  and  $P_{\hat{Y}|X}^*$  that achieve the minimax risk-rate cost are both Gaussian.

*Proposition 3 (Linear regression with rate constraint):* Consider mean-squared loss and define  $\Gamma$  as in (9). Then the distribution that achieves the supremum in Proposition 1,  $P_{\mathbf{X}Y}^*$ , is Gaussian with its mean and covariance matrix specified in (9), and the optimal estimate and minimax risk-information cost (8) are as follows:

$$\hat{Y} = \begin{cases} a \cdot C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} (\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}}) + \mu_Y + Z & \text{if } a > 0 \\ \mu_Y & \text{otherwise} \end{cases} \quad (10)$$

$$\tilde{L}_\lambda^* = \begin{cases} \sigma_Y^2 - C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y} + \frac{\lambda}{2} \log \frac{2\epsilon C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y}}{\lambda \log e} & \text{if } a > 0 \\ \sigma_Y^2 & \text{otherwise,} \end{cases} \quad (11)$$

where

$$a = 1 - \frac{\lambda \log e}{2 C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y}},$$

and  $Z \sim \mathcal{N}(0, \sigma_Z^2)$  is independent of  $\mathbf{X}$  with  $\sigma_Z^2 = \lambda a \log e / 2$ .

Note that this setting does not satisfy the conditions in Proposition 1. Hence, we analyze (8) directly to obtain the optimal  $P_{\mathbf{X}Y}^*$ . Given the optimal  $P_{\mathbf{X}Y}^*$ , Theorem 2 and Lemma 2 can be used to construct the scheme. Operationally,

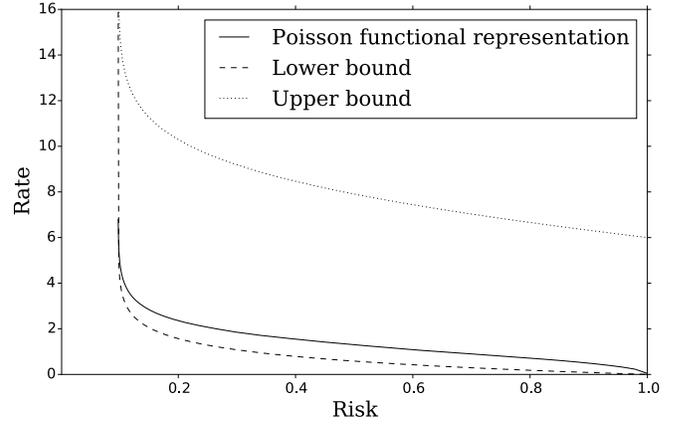


Fig. 3. Tradeoff between the rate and risk in rate-constrained minimax linear regression.

$e_n(x, w)$  is a random quantizer of  $a \cdot C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} (\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}})$  such that the quantization noise follows  $\mathcal{N}(0, \sigma_Z^2)$ . Hence, with this natural choice of the ambiguity set, our formulation recovers a compressed version of the familiar MMSE estimator.

Figure 3 plots the tradeoff between the rate and the risk when  $d = 1$ ,  $\mu_{\mathbf{X}} = \mu_Y = 0$ ,  $\sigma_{\mathbf{X}}^2 = \sigma_Y^2 = 1$ ,  $\sigma_{\mathbf{X}Y} = 0.95$  for the scheme constructed using the Poisson functional representation in [5], with the lower bound given by the minimax risk-information cost  $\tilde{L}_\lambda^*$ , and the upper bound given in Theorem 2. The proof of this proposition is in Section V-C.

The optimal scheme in the above example corresponds to compressing and communicating the minimax optimal rate-unconstrained estimate  $\bar{Y} = C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} (\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}}) + \mu_Y$ , since the optimal  $\hat{Y}$  can be obtained from  $\bar{Y}$  by shifting, scaling and adding noise. This estimate-compress approach can be thought as a *separation* scheme, since we first optimally estimate  $\bar{Y}$ , then optimally communicate it while satisfying the rate constraint. In the next application, we show that such separation is not optimal in general.

##### B. Rate-Constrained Minimax Classification

Let  $\mathcal{Y} = \hat{\mathcal{Y}} = \{1, \dots, k\}$  and  $\mathcal{X}$  be finite,  $\ell(\hat{y}, y) = \mathbf{1}_{\{\hat{y} \neq y\}}$ , and  $\Gamma$  be closed and convex. The following gives the minimax risk-information cost and the optimal estimator for this classification setup.

*Proposition 4:* Consider the setting described above. The minimax risk-information cost is

$$\tilde{L}_\lambda^* = \sup_{P \in \Gamma} \left( 1 + \lambda \inf_{\tilde{P}_{\hat{Y}}} \mathbf{E}_P \left( -\log \sum_y 2^{\lambda^{-1} P_{Y|X}(y|X)} \tilde{P}_{\hat{Y}}(y) \right) \right),$$

the worst-case distribution  $P^*$  is the one attaining the supremum, and the optimal estimator is  $P_{\hat{Y}|X}^*(\hat{Y}|x) \propto 2^{\lambda^{-1} P_{Y|X}^*(\hat{Y}|x)} \tilde{P}_{\hat{Y}}^*(\hat{Y})$ , where  $\tilde{P}_{\hat{Y}}^*$  attains the infimum (when  $P = P^*$ ) and  $P_{\hat{Y}|X}^*$  is obtained from  $P^*$ . In particular, if  $\Gamma$  is symmetric for different values of  $Y$  (i.e., for any  $y_1, y_2 \in \mathcal{Y}$ , there exists permutation  $\pi$  of  $\mathcal{Y}$ ,  $\tau$  of  $\mathcal{X}$  such that  $\pi(y_1) = y_2$  and  $P_{X,Y} \in \Gamma \Leftrightarrow P_{\tau(X),\pi(Y)} \in \Gamma$ ), then

$$\tilde{L}_\lambda^* = \sup_{P \in \Gamma} \left( 1 + \lambda \log k - \lambda \mathbf{E}_P \left( \log \sum_y 2^{\lambda^{-1} P_{Y|X}(y|X)} \right) \right).$$

We can see that when  $\lambda \rightarrow 0$ ,  $P_{\hat{Y}|X}^*$  tends to the MAP estimator (under  $\bar{P}^*$ , the worst-case distribution when  $\lambda = 0$ ). The proof of this proposition is in Section V-D.

To show that the estimate-compress approach is not always optimal, we consider the following.

*Example 1 (Estimate-Compress Not Optimal):* Considering the above classification application, Let  $\mathcal{Y} = \mathcal{Y}_1 \cup \mathcal{Y}_2$ , where  $\mathcal{Y}_1 \cap \mathcal{Y}_2 = \emptyset$  and  $|\mathcal{Y}_i| = k_i$ ,  $i = 1, 2$ . Let  $\Gamma = \{P\}$ , where  $P$  is such that  $(X_1, X_2) \sim \text{Unif}(\mathcal{Y}_1 \times \mathcal{Y}_2)$ , and  $Y = X_i$  with probability  $q_i$  for  $i = 1, 2$  ( $q_2 = 1 - q_1$ ). By Proposition 4, the optimal risk-information cost is

$$\bar{L}_\lambda^* = 1 - \lambda \log \max \left\{ \frac{a_1}{k_1}, \frac{a_2}{k_2} \right\}, \text{ where} \\ a_1 = 2^{\lambda^{-1}q_1} + k_1 - 1, \quad a_2 = 2^{\lambda^{-1}q_2} + k_2 - 1. \quad (12)$$

The optimal estimator is as follows. If  $a_1/k_1 > a_2/k_2$ , then

$$\hat{Y} = \begin{cases} X_1 & \text{w.p. } a_1^{-1}2^{\lambda^{-1}q_1}, \\ \hat{Y}_1 \sim \text{Unif}(\mathcal{Y}_1 \setminus \{X_1\}) & \text{w.p. } a_1^{-1}, \end{cases} \quad (13)$$

and, if  $a_1/k_1 \leq a_2/k_2$ , then simply exchange the subscripts 1 and 2 in the above.

If  $q_1 > q_2$ , then the MAP estimator gives  $\bar{Y} = X_1$ . An estimate-compress approach would either communicate a compressed version of  $\bar{Y} = X_1$  as in (13), or randomly select an element in  $\mathcal{Y}_2$  (giving a risk of  $1 - q_2/k_2$ ). The risk-information cost achieved by this approach is

$$\bar{L}_\lambda = 1 - \lambda \log \max \left\{ \frac{a_1}{k_1}, 2^{\lambda^{-1}q_2}k_2^{-1} \right\}. \quad (14)$$

Now, if  $k_1 \gg k_2$ , the optimal rate-constrained descriptor (13) communicates a compressed version of  $X_2$ , and the risk of estimate-compress in (14) is larger than (12). Moreover, the gap between the rates needed by the two approaches for a fixed risk can be unbounded. Let  $q_1 = 1 - q_2 = 2/3$ ,  $k_2 = 2$ ,  $k_1 \geq 15$ . The minimum rate needed to achieve a risk  $2/3$  is 1 (by  $\hat{Y} = X_2$ ). For the estimate-compress approach, since  $\hat{Y} \sim \text{Unif}(\mathcal{Y}_2)$  gives a risk  $5/6$ , we have to compress  $X_1$  (by passing it through a symmetric channel with  $P\{\hat{Y} = X_1\} = 1/2$ ) to achieve a risk  $2/3$ , which as  $k_1$  increases, requires an unbounded rate

$$I(X; \hat{Y}) = H(\hat{Y}) - H(\hat{Y}|X_1) = \log k_1 - \frac{1}{2} \log(k_1 - 1) - \frac{1}{2}.$$

Figure 4 compares risk-rate tradeoff for the optimal scheme, the lower bound obtained from the optimal risk-information tradeoff (12), the upper bound in Theorem 1, and the risk-information cost of the estimate-compress approach (14) for  $q_1 = 1 - q_2 = 2/3$ ,  $k_1 = 2^{32}$ ,  $k_2 = 2$ . Note that the optimal scheme (attaining the optimal risk-rate tradeoff) performs time-sharing between encoding  $X_1$  using 32 bits with risk  $1/3$ , encoding  $X_2$  using 1 bit with risk  $2/3$ , and fixing the output at one value of  $X_2$  with zero rate needed and risk  $5/6$ . The mutual information needed by the estimate-compress approach (which is a lower bound on the actual rate needed by this approach) is strictly greater than the optimal rate (except when the risk is at its minimum  $1/3$  or maximum  $5/6$ ).

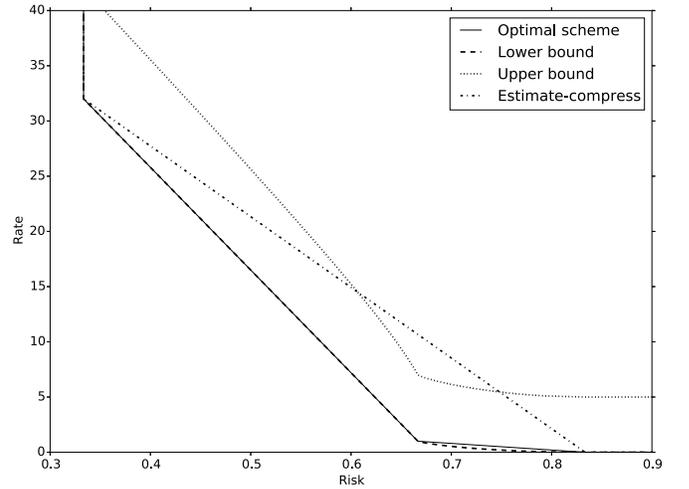


Fig. 4. Tradeoff between the rate and risk in rate-constrained minimax linear classification for the optimal scheme, lower bound (12), upper bound by Theorem 1, and estimate-compress approach (14).

## V. PROOFS OF THE RESULTS

### A. Proof of Theorem 2

To prove this theorem, we first invoke a minimax result for relative entropy in [14], which generalizes the redundancy-capacity theorem in [15] and use Lemma 2.

Instead of using the Huffman code as in Theorem 1, we apply a code over positive integers (e.g. Elias delta code [16]) on  $k(X, W)$  to produce  $M$ .

*Proof:* The lower bound follows from  $\mathbf{E}_P[|M|] \geq H_P(M) \geq I_P(X; \hat{Y})$ . To prove the upper bound, we fix any  $P_{\hat{Y}|X}$ , and show that the following risk-rate cost is achievable:

$$L' = \sup_{P \in \Gamma} \left( \mathbf{E}_P[\ell(\hat{Y}, Y)] + \lambda (I_P(X; \hat{Y}) + 2 \log(I_P(X; \hat{Y}) + 1) + 5.5) \right).$$

Let

$$g(P, \tilde{P}_{\hat{Y}}) = \mathbf{E}_P[\ell(\hat{Y}, Y)] + \lambda \left( \int D(P_{\hat{Y}|X=x} \| \tilde{P}_{\hat{Y}}) dP(x) + 2 \log \left( \int D(P_{\hat{Y}|X=x} \| \tilde{P}_{\hat{Y}}) dP(x) + 1 \right) + 5.4 \right).$$

Note that  $g$  is concave in  $P$  for fixed  $\tilde{P}_{\hat{Y}}$ , since  $\mathbf{E}_P[\ell(\hat{Y}, Y)]$  and  $\int D(P_{\hat{Y}|X=x} \| \tilde{P}_{\hat{Y}}) dP(x)$  are linear in  $P$ . Also  $g$  is quasiconvex in  $\tilde{P}_{\hat{Y}}$  for fixed  $P$ , since  $\int D(P_{\hat{Y}|X=x} \| \tilde{P}_{\hat{Y}}) dP(x)$  is convex in  $\tilde{P}_{\hat{Y}}$ , and is lower semicontinuous in  $\tilde{P}_{\hat{Y}}$ , since  $D(P_{\hat{Y}|X=x} \| \tilde{P}_{\hat{Y}})$  is lower semicontinuous with respect to the topology of weak convergence [17]. Hence by Fatou's lemma,  $\int D(P_{\hat{Y}|X=x} \| \tilde{P}_{\hat{Y}}) dP(x)$  is lower semicontinuous.

We write  $P_{\hat{Y}|X} \circ P$  for the distribution of  $\hat{Y}$  when  $(X, Y) \sim P$  and  $\hat{Y}|X=x \sim P_{\hat{Y}|X}(\cdot|x)$ . Let  $\Gamma_{\hat{Y}} = \{P_{\hat{Y}|X} \circ P : P \in \Gamma\}$  and  $\overline{\Gamma_{\hat{Y}}}$  be the closure of  $\Gamma_{\hat{Y}}$  in the topology of weak convergence. It can be shown using the same arguments as in [14] (on  $g$  instead of relative entropy, and using Sion's

minimax theorem [13] instead of Lemma 2 in [14]) that if  $\Gamma_{\hat{Y}}$  is uniformly tight, then

$$\inf_{\tilde{P}_{\hat{Y}} \in \overline{\Gamma_{\hat{Y}}}} \sup_{P \in \Gamma} g(P, \tilde{P}_{\hat{Y}}) \leq \sup_{P \in \Gamma} \inf_{\tilde{P}_{\hat{Y}} \in \overline{\Gamma_{\hat{Y}}}} g(P, \tilde{P}_{\hat{Y}}) = L' - 0.1\lambda,$$

where the second infimum is over all distributions  $\tilde{P}_{\hat{Y}}$  on  $\hat{\mathcal{Y}}$ . Therefore there exists  $P_{\hat{Y}}^* \in \overline{\Gamma_{\hat{Y}}}$  such that

$$\sup_{P \in \Gamma} g(P, P_{\hat{Y}}^*) \leq L'. \quad (15)$$

If  $\Gamma_{\hat{Y}}$  is not uniformly tight, then by Lemma 4 in [14],

$$\sup_{P \in \Gamma} \inf_{\tilde{P}_{\hat{Y}}} \int D(P_{\hat{Y}|X=x} \| \tilde{P}_{\hat{Y}}) dP(x) = \infty.$$

Hence  $L' = \sup_{P \in \Gamma} \inf_{\tilde{P}_{\hat{Y}}} g(P, \tilde{P}_{\hat{Y}}) + 0.1\lambda = \infty$ .

Applying Lemma 2 to  $P_{\hat{Y}|X}$  and  $P_{\hat{Y}}^*$ , we obtain  $W$  independent of  $X$ , and  $K = k(X, W) \in \{1, 2, \dots\}$  and  $\hat{Y} = \hat{Y}(K, W)$  following the conditional distribution  $P_{\hat{Y}|X}(\cdot|x)$  such that,

$$\mathbf{E} [\log K | X = x] \leq D(P_{\hat{Y}|X} \| P_{\hat{Y}}^* | X = x) + 1.6.$$

We then use Elias delta code [16] for  $K$  to produce  $M$ . To encode an integer  $k \geq 1$ , let  $n$  be the number of digits in the binary representation of  $k$ , and let  $l$  be the number of digits in the binary representation of  $n$ . The codeword is formed by appending  $l - 1 = \lfloor \log(\lfloor \log k \rfloor + 1) \rfloor$  zeroes, the binary representation of  $n$  ( $\lfloor \log(\lfloor \log k \rfloor + 1) \rfloor + 1$  bits), and the binary representation of  $k$  except the leading digit ( $\lfloor \log k \rfloor$  bits). Since the length of the Elias delta codeword for an integer  $k$  is upper bounded by  $\log k + 2 \log(\log k + 1) + 1$ , by Jensen's inequality,

$$\begin{aligned} \mathbf{E}_P[|M|] &\leq \mathbf{E}_P[\log K] + 2 \log \left( \mathbf{E}_P[\log K] + 1 \right) + 1 \\ &\leq \int D(P_{\hat{Y}|X=x} \| P_{\hat{Y}}^*) dP(x) \\ &\quad + 2 \log \left( \int D(P_{\hat{Y}|X=x} \| P_{\hat{Y}}^*) dP(x) + 1 \right) + 5.4. \end{aligned}$$

Thus,

$$\tilde{L}_\lambda^* \leq \sup_{P \in \Gamma} (\mathbf{E}_P[\ell(\hat{Y}, Y) + \lambda |M|]) \leq \sup_{P \in \Gamma} g(P, P_{\hat{Y}}^*) \leq L'. \quad \blacksquare$$

### B. Proof of Proposition 2

Since  $\mathbf{E}[d_{\text{TV}}(P, P_n)] \leq (1/2)\sqrt{(|\mathcal{X}| \cdot |\mathcal{Y}|)/n}$  (e.g. see [18]), by Markov's inequality, we have

$$\mathbf{P}\{P \notin \Gamma\} \leq \frac{1}{2\epsilon} \sqrt{\frac{|\mathcal{X}| \cdot |\mathcal{Y}|}{n}}.$$

Let  $\hat{Y}|\{X=x\} \sim P_{\hat{Y}|X}^*(\cdot|x)$ . If  $P \in \Gamma$ , then

$$\begin{aligned} &L_\lambda(e, f, P) \\ &= \mathbf{E}_P[\ell(\hat{Y}, Y) + \lambda|M|] \\ &\leq \mathbf{E}_P[\ell(\hat{Y}, Y)] + \mathbf{E}_P[\log K] + 2 \log \left( \mathbf{E}_P[\log K] + 1 \right) + 1 \\ &\leq g(P, \tilde{P}_{\hat{Y}}^*) \end{aligned}$$

$$\begin{aligned} &\leq \sup_{\tilde{P} \in \Gamma} g(\tilde{P}, \tilde{P}_{\hat{Y}}^*) \\ &\stackrel{(a)}{\leq} \sup_{\tilde{P} \in \Gamma} \inf_{\tilde{P}_{\hat{Y}} \in \overline{\Gamma_{\hat{Y}}}} g(\tilde{P}, \tilde{P}_{\hat{Y}}) + 0.1\lambda \\ &= \sup_{\tilde{P} \in \Gamma} (\mathbf{E}_{\tilde{P}}[\ell(\hat{Y}, Y)] + \lambda(I_{\tilde{P}}(X; \hat{Y}) \\ &\quad + 2 \log(I_{\tilde{P}}(X; \hat{Y}) + 1) + 5.5)) \\ &\leq \sup_{\tilde{P} \in \Gamma} (\tilde{L}_\lambda(P_{\hat{Y}|X}^*, \tilde{P}) + 2\lambda \log(\lambda^{-1} \tilde{L}_\lambda(P_{\hat{Y}|X}^*, \tilde{P}) + 1) + 5.5\lambda) \\ &= \sup_{\tilde{P} \in \Gamma} (\tilde{L}_\lambda^*(\{\tilde{P}\}) + 2\lambda \log(\lambda^{-1} \tilde{L}_\lambda^*(\{\tilde{P}\}) + 1) + 5.5\lambda) \\ &\stackrel{(b)}{\leq} \tilde{L}_\lambda^*(\{P\}) + 2\alpha\epsilon + 2\lambda \log(\lambda^{-1} \tilde{L}_\lambda^*(\{P\}) + 2\alpha\epsilon + 1) + 5.5\lambda \\ &\leq L_\lambda^*(\{P\}) + 2\alpha\epsilon + 2\lambda \log(\lambda^{-1} L_\lambda^*(\{P\}) + 2\alpha\epsilon + 1) + 5.5\lambda, \end{aligned}$$

where  $K$ ,  $g$  and  $\tilde{P}_{\hat{Y}}^*$  are defined as in the proof of Theorem 2, (a) is due to (15). To show (b), fix any  $\tilde{P} \in \Gamma$  and let  $((X, Y), (\tilde{X}, \tilde{Y}))$  be a coupling of  $(P, \tilde{P})$  (i.e.,  $(X, Y) \sim P$ ,  $(\tilde{X}, \tilde{Y}) \sim \tilde{P}$ ) such that  $\mathbf{P}\{(X, Y) \neq (\tilde{X}, \tilde{Y})\} \leq \epsilon$ . Let  $\check{X}|\{\tilde{X} = \tilde{x}\} \sim P_{\check{X}|\tilde{X}}(\cdot|\tilde{x})$  be conditionally independent of  $(X, Y, \tilde{Y})$  given  $\tilde{X}$ . For any  $P_{\hat{Y}|X}$ , let  $\hat{Y}|\{X=x\} \sim P_{\hat{Y}|X}(\cdot|x)$  be conditionally independent of  $(Y, \tilde{X}, \tilde{Y}, \check{X})$  given  $X$ , and let  $\check{Y}|\{\check{X} = \check{x}\} \sim P_{\check{Y}|\check{X}}(\cdot|\check{x})$  be conditionally independent of  $(X, Y, \tilde{X}, \tilde{Y}, \hat{Y})$  given  $\check{X}$ . We have

$$\begin{aligned} &\tilde{L}_\lambda^*(\{\tilde{P}\}) \\ &\leq \tilde{L}_\lambda(P_{\hat{Y}|\tilde{X}}, \tilde{P}) \\ &= \mathbf{E}[\ell(\tilde{Y}, \tilde{Y})] + \lambda I(\tilde{X}; \tilde{Y}) \\ &\stackrel{(c)}{\leq} \mathbf{E}[\ell(\hat{Y}, Y)] + \alpha d_{\text{TV}}(P_{\hat{Y}, \tilde{Y}}, P_{Y, \tilde{Y}}) + \lambda I(\check{X}; \check{Y}) \\ &\stackrel{(d)}{\leq} \mathbf{E}[\ell(\hat{Y}, Y)] + \alpha d_{\text{TV}}(P_{\hat{Y}, \tilde{X}}, P_{Y, \tilde{X}}) + \lambda I(X; \hat{Y}) \\ &\leq \mathbf{E}[\ell(\hat{Y}, Y)] + \alpha \mathbf{P}\{(\tilde{Y}, \tilde{X}) \neq (Y, X)\} + \lambda I(X; \hat{Y}) \\ &\leq \mathbf{E}[\ell(\hat{Y}, Y)] + \alpha \mathbf{P}\{(\tilde{Y}, \tilde{X}) \neq (Y, X)\} \\ &\quad + \alpha \mathbf{P}\{\tilde{X} \neq \check{X}\} + \lambda I(X; \hat{Y}) \\ &\leq \mathbf{E}[\ell(\hat{Y}, Y)] + 2\alpha\epsilon + \lambda I(X; \hat{Y}) \\ &= \tilde{L}_\lambda(P_{\hat{Y}|X}, P) + 2\alpha\epsilon, \end{aligned}$$

where (c) is by  $\alpha = \sup_{y, \hat{y}} \ell(y, \hat{y}) - \inf_{y, \hat{y}} \ell(y, \hat{y})$  and the data processing inequality, and (d) is because  $P_{\hat{Y}, \tilde{X}, \tilde{Y}} = P_{\hat{Y}, \tilde{X}} P_{\hat{Y}|X}$  and  $P_{Y, \tilde{X}, \tilde{Y}} = P_{Y, \tilde{X}} P_{\hat{Y}|X}$ . Taking infimum over  $P_{\hat{Y}|X}$ , we have  $\tilde{L}_\lambda^*(\{\tilde{P}\}) \leq \tilde{L}_\lambda^*(\{P\}) + 2\alpha\epsilon$ , and thus (b) holds.

### C. Proof of Proposition 3

Without loss of generality, assume  $\mu_{\mathbf{X}} = \mathbf{0}$  and  $\mu_Y = 0$ . We first prove “ $\leq$ ” in (11). For this, fix  $P_{\hat{Y}|X}$  as given in the Proposition and consider any  $P \in \Gamma$ . When  $\frac{\lambda \log e}{2} < C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y}$ , we have

$$\begin{aligned} \mathbf{E}_P[\ell(\hat{Y}, Y)] &= \mathbf{E}_P[(\hat{Y} - Y)^2] \\ &\leq \sigma_Y^2 + \frac{\lambda \log e}{2} - C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y}, \text{ and} \end{aligned}$$

$$\begin{aligned} I_P(\mathbf{X}; \hat{Y}) &= h(\hat{Y}) - h(\hat{Y}|\mathbf{X}) \\ &\leq \frac{1}{2} \log \left( \frac{2 C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}\hat{Y}}}{\lambda \log e} \right). \end{aligned}$$

Therefore,

$$\inf_{P_{\hat{Y}|\mathbf{X}}} \sup_{P \in \Gamma} \left( \mathbf{E}_P [\ell(\hat{Y}, Y)] + \lambda I_P(\mathbf{X}; \hat{Y}) \right) \leq \text{R.H.S. of (11)}.$$

It can also be checked that the above relation holds when  $\frac{\lambda \log e}{2} \geq C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}\hat{Y}}$ , and thus we have proved “ $\leq$ ” in (11).

To prove “ $\geq$ ” in (11), fix a Gaussian  $P_{\mathbf{X}\hat{Y}}$  with its mean and covariance matrix specified in (9) and consider an arbitrary  $P_{\hat{Y}|\mathbf{X}}$ . We have

$$\begin{aligned} \mathbf{E}_P [\ell(\hat{Y}, Y)] &= \mathbf{E}_P [(Y - \hat{Y})^2] \\ &= \sigma_Y^2 - C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}\hat{Y}} + \mathbf{E}_P [(\hat{Y} - C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{X})^2], \end{aligned}$$

and

$$\begin{aligned} I_P(\mathbf{X}; \hat{Y}) &= I_P(C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{X}; \hat{Y}) \\ &\geq h(C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{X}) - h(C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{X} - \hat{Y}) \\ &\geq \frac{1}{2} \log C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}\hat{Y}} - \frac{1}{2} \log \mathbf{E}_P [(\hat{Y} - C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{X})^2]. \end{aligned}$$

Letting  $\gamma = \mathbf{E}_P [(\hat{Y} - C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{X})^2]$ , we have

$$\begin{aligned} \mathbf{E}_P [\ell(\hat{Y}, Y)] + \lambda I_P(\mathbf{X}; \hat{Y}) &\geq \sigma_Y^2 - C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}\hat{Y}} + \frac{\lambda}{2} \log C_{\mathbf{X}\hat{Y}}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}\hat{Y}} + \gamma - \frac{\lambda \log \gamma}{2} \\ &\geq \text{R.H.S. of (11)}, \end{aligned}$$

where the second inequality follows by evaluating the minimum value of  $\gamma - \frac{\lambda \log \gamma}{2}$ . Combing this with the above completes the proof of Proposition 3.

#### D. Proof of Proposition 4

Assume  $\Gamma$  is closed and convex. By Proposition 1, the minimax risk-information cost is  $\tilde{L}_\lambda^* = \sup_{P \in \Gamma} \inf_{P_{\hat{Y}|\mathbf{X}}} \tilde{L}_\lambda(P_{\hat{Y}|\mathbf{X}}, P)$ , where

$$\begin{aligned} \inf_{P_{\hat{Y}|\mathbf{X}}} \tilde{L}_\lambda(P_{\hat{Y}|\mathbf{X}}, P) &= \inf_{P_{\hat{Y}|\mathbf{X}}} \left( \mathbf{E}_P [\ell(\hat{Y}, Y)] + \lambda I_P(\mathbf{X}; \hat{Y}) \right) \\ &= \inf_{P_{\hat{Y}|\mathbf{X}}} \left( P\{\hat{Y} \neq Y\} + \lambda \inf_{\tilde{P}_{\hat{Y}}} \int D(P_{\hat{Y}|\mathbf{X}=x} \| \tilde{P}_{\hat{Y}}) dP(x) \right) \\ &= \inf_{\tilde{P}_{\hat{Y}}, P_{\hat{Y}|\mathbf{X}}} \left( P\{\hat{Y} \neq Y\} + \lambda \int D(P_{\hat{Y}|\mathbf{X}=x} \| \tilde{P}_{\hat{Y}}) dP(x) \right) \\ &= 1 + \lambda \inf_{\tilde{P}_{\hat{Y}}, P_{\hat{Y}|\mathbf{X}}} \mathbf{E}_P \left( \sum_y P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X}) \right. \\ &\quad \cdot \left. \left( \log \frac{P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X})}{\tilde{P}_{\hat{Y}}(y)} - \lambda^{-1} P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X}) \right) \right) \end{aligned}$$

$$\begin{aligned} &= 1 + \lambda \inf_{\tilde{P}_{\hat{Y}}} \inf_{P_{\hat{Y}|\mathbf{X}}} \mathbf{E}_P \left( \sum_y P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X}) \right. \\ &\quad \cdot \left. \left( \log \frac{P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X})}{2^{\lambda^{-1} P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X})} \tilde{P}_{\hat{Y}}(y) / \sum_{y'} 2^{\lambda^{-1} P_{\hat{Y}|\mathbf{X}}(y'|\mathbf{X})} \tilde{P}_{\hat{Y}}(y')} \right) \right. \\ &\quad \left. - \log \sum_y 2^{\lambda^{-1} P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X})} \tilde{P}_{\hat{Y}}(y) \right) \\ &\stackrel{(a)}{=} 1 + \lambda \inf_{\tilde{P}_{\hat{Y}}} \mathbf{E}_P \left( - \log \sum_y 2^{\lambda^{-1} P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X})} \tilde{P}_{\hat{Y}}(y) \right), \end{aligned}$$

where (a) is due to the fact that relative entropy is nonnegative and equality is attained when  $P_{\hat{Y}|\mathbf{X}}(y|x) \propto 2^{\lambda^{-1} P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X})} \tilde{P}_{\hat{Y}}(y)$ .

Next we consider the case in which  $\Gamma$  is symmetric. Consider the minimax risk-information cost

$$\begin{aligned} \tilde{L}_\lambda^* &= \inf_{P_{\hat{Y}|\mathbf{X}}} \sup_{P \in \Gamma} \tilde{L}_\lambda(P_{\hat{Y}|\mathbf{X}}, P) \\ &= \inf_{P_{\hat{Y}|\mathbf{X}}} \sup_{P \in \Gamma} \left( \mathbf{E}_P [\ell(\hat{Y}, Y)] + \lambda I_P(\mathbf{X}; \hat{Y}) \right). \end{aligned}$$

For any  $i, j \in \mathcal{Y} = \{1, \dots, k\}$ , let  $\pi_{ij}$  be the permutation over  $\mathcal{Y}$  such that  $\pi_{ij}(i) = j$  and let  $\tau_{ij}$  be the corresponding permutation over  $\mathcal{X}$  in the symmetry assumption. Since the function

$$P_{\hat{Y}|\mathbf{X}} \mapsto \sup_{P \in \Gamma} \tilde{L}_\lambda(P_{\hat{Y}|\mathbf{X}}, P)$$

is convex and symmetric about  $\pi_{ij}$  and  $\tau_{ij}$  (i.e.,  $\sup_{P \in \Gamma} \tilde{L}_\lambda(P_{\hat{Y}|\mathbf{X}}, P) = \sup_{P \in \Gamma} \tilde{L}_\lambda(P_{\pi_{ij} \hat{Y} | \tau_{ij} \mathbf{X}}, P)$ ), to find its infimum, we only need to consider  $P_{\hat{Y}|\mathbf{X}}$ 's satisfying  $P_{\hat{Y}|\mathbf{X}} = P_{\pi_{ij} \hat{Y} | \tau_{ij} \mathbf{X}}$  for all  $i, j$  (if not, we can instead consider the average of  $P_{\pi_{ij} \hat{Y} | \tau_{ij} \mathbf{X}}$  for  $a$  from 1 up to the product of the periods of  $\pi_{ij}$  and  $\tau_{ij}$ , which gives a value of the function not larger than that of  $P_{\hat{Y}|\mathbf{X}}$ ). For brevity we say  $P_{\hat{Y}|\mathbf{X}}$  is symmetric if it satisfies this condition.

Fix any symmetric  $P_{\hat{Y}|\mathbf{X}}$ . Since the function

$$P \mapsto \tilde{L}_\lambda(P_{\hat{Y}|\mathbf{X}}, P)$$

is concave and symmetric about  $\pi_{ij}$  and  $\tau_{ij}$  (i.e.,  $\tilde{L}_\lambda(P_{\hat{Y}|\mathbf{X}}, P_{\mathbf{X}, Y}) = \tilde{L}_\lambda(P_{\hat{Y}|\mathbf{X}}, P_{\tau_{ij} \mathbf{X}, \pi_{ij} Y})$ ), to find its supremum, we only need to consider symmetric  $P$ 's. Hence,

$$\begin{aligned} \tilde{L}_\lambda^* &= \inf_{P_{\hat{Y}|\mathbf{X}}} \sup_{P \in \Gamma_{\text{symm}}} \tilde{L}_\lambda(P_{\hat{Y}|\mathbf{X}}, P) \\ &= \inf_{P_{\hat{Y}|\mathbf{X}}} \sup_{P \in \Gamma_{\text{symm}}} \left( \mathbf{E}_P [\ell(\hat{Y}, Y)] + \lambda I_P(\mathbf{X}; \hat{Y}) \right) \\ &= \inf_{P_{\hat{Y}|\mathbf{X}}} \sup_{P \in \Gamma_{\text{symm}}} \left( P\{\hat{Y} \neq Y\} + \lambda (\log k - H_P(\hat{Y}|\mathbf{X})) \right) \\ &= 1 + \lambda \log k + \lambda \inf_{P_{\hat{Y}|\mathbf{X}}} \sup_{P \in \Gamma_{\text{symm}}} \mathbf{E}_P \left( \sum_y P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X}) \left( \log P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X}) - \lambda^{-1} P_{\hat{Y}|\mathbf{X}}(y|\mathbf{X}) \right) \right) \end{aligned}$$

$$\begin{aligned}
 &= 1 + \lambda \log k + \lambda \inf_{P_{\hat{Y}|X} \text{ symm}} \sup_{P \in \Gamma \text{ symm}} \mathbf{E}_P \left( \right. \\
 &\quad \left. \sum_y P_{\hat{Y}|X}(y|X) \log \frac{P_{\hat{Y}|X}(y|X)}{2^{\lambda^{-1}P_{Y|X}(y|X)} / \sum_{y'} 2^{\lambda^{-1}P_{Y|X}(y'|X)}} \right. \\
 &\quad \left. - \log \sum_y 2^{\lambda^{-1}P_{Y|X}(y|X)} \right) \\
 &\geq 1 + \lambda \log k + \lambda \inf_{P_{\hat{Y}|X} \text{ symm}} \sup_{P \in \Gamma \text{ symm}} \mathbf{E}_P \left( \right. \\
 &\quad \left. - \log \sum_y 2^{\lambda^{-1}P_{Y|X}(y|X)} \right) \\
 &= \sup_{P \in \Gamma \text{ symm}} \left( 1 + \lambda \log k - \lambda \mathbf{E}_P \log \sum_y 2^{\lambda^{-1}P_{Y|X}(y|X)} \right),
 \end{aligned}$$

where the inequality is because relative entropy is nonnegative (and equality is attained when  $P_{\hat{Y}|X}(y|x) \propto 2^{\lambda^{-1}P_{Y|X}(y|x)}$ ). Note that

$$\begin{aligned}
 &1 + \lambda \log k - \lambda \mathbf{E}_P \log \sum_y 2^{\lambda^{-1}P_{Y|X}(y|X)} \\
 &= \inf_{P_{\hat{Y}|X}} \left( P\{\hat{Y} \neq Y\} + \lambda(\log k - H_P(\hat{Y}|X)) \right)
 \end{aligned}$$

is an infimum of affine functions of  $P$ , hence it is concave in  $P$ . Also it is symmetric about  $\pi$  and  $\tau$ , hence

$$\begin{aligned}
 \tilde{L}_\lambda^* &\geq \sup_{P \in \Gamma \text{ symm.}} \left( 1 + \lambda \log k - \lambda \mathbf{E}_P \log \sum_y 2^{\lambda^{-1}P_{Y|X}(y|X)} \right) \\
 &= \sup_{P \in \Gamma} \left( 1 + \lambda \log k - \lambda \mathbf{E}_P \log \sum_y 2^{\lambda^{-1}P_{Y|X}(y|X)} \right).
 \end{aligned}$$

The other direction follows from setting  $\tilde{P}_{\hat{Y}}(y) = 1/k$ .

## VI. CONCLUSION

We introduced the minimax learning problem in which the inference is distributed between two nodes (e.g., a mobile device and a cloud) with a constraint on the communication rate between them. We showed that the minimax risk-rate cost can be approximated by the minimax risk-information cost, which is significantly easier to evaluate and leads to a general method for the design a near-optimal descriptor-estimator pair. We showed that the naive estimate-compress scheme for rate-constrained inference is not in general optimal. Our results also provide a new one-shot operational interpretation of the information bottleneck and extends it to the minimax robust setting. Designing efficient algorithms for practical applications is left for future research. Extending the work to the case in which the data is also distributed, hence learning has a communication constraint, would also be of great interest to real world applications such as federated learning [19], [20].

## REFERENCES

[1] F. Farnia and D. Tse, "A minimax approach to supervised learning," in *Proc. Adv. Neural Inf. Process. Syst.*, 2016, pp. 4240–4248.

[2] H. Namkoong and J. C. Duchi, "Variance-based regularization with convex objectives," in *Proc. Adv. Neural Inf. Process. Syst.*, 2017, pp. 2975–2984.

[3] J. Lee and M. Raginsky, "Minimax statistical learning with Wasserstein distances," 2017, *arXiv:1705.07815*. [Online]. Available: <http://arxiv.org/abs/1705.07815>

[4] C. T. Li, X. Wu, A. Özgür, and A. El Gamal, "Minimax learning for remote prediction," in *Proc. IEEE Int. Symp. Inf. Theory*, Jun. 2018, pp. 541–545.

[5] C. T. Li and A. E. Gamal, "Strong functional representation lemma and applications to coding theorems," *IEEE Trans. Inf. Theory*, vol. 64, no. 11, pp. 6967–6978, Nov. 2018.

[6] A. Dembo and T. Weissman, "The minimax distortion redundancy in noisy source coding," *IEEE Trans. Inf. Theory*, vol. 49, no. 11, pp. 3020–3030, Nov. 2003.

[7] P. Harsha, R. Jain, D. McAllester, and J. Radhakrishnan, "The communication complexity of correlation," *IEEE Trans. Inf. Theory*, vol. 56, no. 1, pp. 438–449, Jan. 2010.

[8] M. Braverman and A. Garg, "Public vs private coin in bounded-round information," in *Proc. Int. Colloq. Automata, Lang., Program.* Berlin, Germany: Springer, 2014, pp. 502–513.

[9] N. Tishby, F. C. Pereira, and W. Bialek, "The information bottleneck method," 2000, *arXiv:physics/0004057*. [Online]. Available: <https://arxiv.org/abs/physics/0004057>

[10] P. Harremoës and N. Tishby, "The information bottleneck revisited or how to choose a good distortion measure," in *Proc. IEEE Int. Symp. Inf. Theory*, Jun. 2007, pp. 566–570.

[11] T. A. Courtade and T. Weissman, "Multiterminal source coding under logarithmic loss," *IEEE Trans. Inf. Theory*, vol. 60, no. 1, pp. 740–761, Jan. 2014.

[12] D. Strouse and D. J. Schwab, "The deterministic information bottleneck," *Neural Comput.*, vol. 29, no. 6, pp. 1611–1630, Jun. 2017.

[13] M. Sion, "On general minimax theorems," *Pacific J. Math.*, vol. 8, no. 1, pp. 171–176, Mar. 1958.

[14] D. Haussler, "A general minimax result for relative entropy," *IEEE Trans. Inf. Theory*, vol. 43, no. 4, pp. 1276–1280, Jul. 1997.

[15] R. G. Gallager, "Source coding with side information and universal coding," Lab. Inf. Decis. Syst., MIT, Cambridge, MA, USA, Tech. Rep. LIDS-P-937, 1979.

[16] P. Elias, "Universal codeword sets and representations of the integers," *IEEE Trans. Inf. Theory*, vol. IT-21, no. 2, pp. 194–203, Mar. 1975.

[17] E. Posner, "Random coding strategies for minimum entropy," *IEEE Trans. Inf. Theory*, vol. IT-21, no. 4, pp. 388–391, Jul. 1975.

[18] S. Kamath, A. Orlitsky, D. Pichapati, and A. T. Suresh, "On learning distributions from their samples," in *Proc. Conf. Learn. Theory*, 2015, pp. 1066–1100.

[19] J. Konečný, H. Brendan McMahan, D. Ramage, and P. Richtárik, "Federated optimization: Distributed machine learning for on-device intelligence," 2016, *arXiv:1610.02527*. [Online]. Available: <http://arxiv.org/abs/1610.02527>

[20] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. Y. Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2017, pp. 1273–1282.

**Cheuk Ting Li** (Member, IEEE) received the B.Sc. degree in mathematics and the B.Eng. degree in information engineering from The Chinese University of Hong Kong in 2012, and the M.S. and Ph.D. degrees in electrical engineering from Stanford University in 2014 and 2018, respectively. He was a Postdoctoral Scholar with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley. He joined the Department of Information Engineering, The Chinese University of Hong Kong, in January 2020. His research interests include generation of random variables, one-shot schemes in information theory, wireless communications, and information-theoretic secrecy.

**Xiugang Wu** (Member, IEEE) received the B.Eng. degree (Hons.) in electronics and information engineering from Tongji University, Shanghai, China, in 2007, and the M.A.Sc. and Ph.D. degrees in electrical and computer engineering from the University of Waterloo, Waterloo, ON, Canada, in 2009 and 2014, respectively. He was a Postdoctoral Fellow with the Department of Electrical Engineering, Stanford University, Stanford, CA, from 2015 to 2018. He is currently an Assistant Professor with the University of Delaware, Newark, DE, where he is jointly appointed with the Department of Electrical and Computer Engineering and the Department of Computer and Information Sciences. His research interests are in information theory, networks, data science, and the interplay between them. He was a recipient of the 2017 NSF Center for Science of Information (CSoI) Postdoctoral Fellowship.

**Ayfer Özgür** (Senior Member, IEEE) received the B.Sc. degrees in electrical engineering and physics and the M.Sc. degree in electrical engineering from Middle East Technical University, Turkey, in 2001 and 2004, respectively, and the Ph.D. degree from the Information Processing Group, EPFL, in 2009. She is currently an Associate Professor of electrical engineering and the Chambers Faculty Scholar with the School of Engineering, Stanford University. Her current research interests include distributed communication and learning, wireless systems, and information theory.

Dr. Özgür received the EPFL Best Ph.D. Thesis Award in 2010, the NSF CAREER Award in 2013, the Okawa Foundation Research Grant and the IEEE Communication Theory Technical Committee (CTTC) Early Achievement Award in 2018, the Google Faculty Research Award in 2019, and was selected as the inaugural Goldsmith Lecturer of the IEEE IT Society in 2020.

**Abbas El Gamal** (Life Fellow, IEEE) received the B.Sc. degree (Hons.) from Cairo University in 1972, and the M.S. degree in statistics and the Ph.D. degree in electrical engineering from Stanford University in 1977 and 1978, respectively. He is currently the Hitachi America Professor with the School of Engineering, Stanford University. From 1978 to 1980, he was an Assistant Professor of electrical engineering with USC. He has been on the faculty of the Department of Electrical Engineering, Stanford University, since 1981. From 2003 to 2012, he was the Director of the Information Systems Laboratory, Stanford University. From 2012 to 2017, he was the Chair of the Department of Electrical Engineering, Stanford University. His research contributions have been in network information theory, FPGAs, digital imaging devices and systems, and smart grid modeling and control. He has authored or coauthored over 230 articles and holds 35 patents in these areas. He is the coauthor of the book *Network Information Theory* (Cambridge Press, 2011). He is a member of the U.S. National Academy of Engineering. He received several honors and awards for his research contributions, including the 2016 IEEE Richard Hamming Medal, the 2012 Claude E. Shannon Award, and the 2004 INFOCOM Paper Award. He served on the Board of Governors of the Information Theory Society from 2009 to 2016 and was the President in 2014.