Minimax Learning for Remote Prediction

Cheuk Ting Li*, Xiugang Wu[†], Ayfer Ozgur[†], Abbas El Gamal[†]
*Department of Electrical Engineering and Computer Sciences
University of California, Berkeley
Email: ctli@berkeley.edu

†Department of Electrical Engineering
Stanford University

Email: {x23wu, aozgur}@stanford.edu; abbas@ee.stanford.edu

Abstract

The classical problem of supervised learning is to infer an accurate predictor of a target variable Y from a measured variable X by using a finite number of labeled training samples. Motivated by the increasingly distributed nature of data and decision making, in this paper we consider a variation of this classical problem in which the prediction is performed remotely based on a rate-constrained description M of X. Upon receiving M, the remote node computes an estimate \hat{Y} of Y. We follow the recent minimax approach to study this learning problem and show that it corresponds to a one-shot minimax noisy source coding problem. We then establish information theoretic bounds on the risk-rate Lagrangian cost and a general method to design a near-optimal descriptor-estimator pair, which can be viewed as a rate-constrained analog to the maximum conditional entropy principle used in the classical minimax learning problem. Our results show that a naive estimate-compress scheme for rate-constrained prediction is not in general optimal.

I. Introduction

The classical problem of supervised learning is to infer an accurate predictor of a target variable Y from a measured variable X on the basis of n labeled training samples $\{(X_i,Y_i)\}_{i=1}^n$ 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 predictors and a loss function that quantifies the risk of each predictor, ERM chooses the predictor with minimal risk under the empirical distribution of samples. To avoid overfitting, the set of allowable predictors is restricted to a class with limited complexity.

Recently, an alternative viewpoint has emerged which seeks distributionally robust predictors. Given the labeled training samples, this approach learns a predictor by minimizing its worst-case risk over an ambiguity distribution set centered at the empirical distribution of samples. In other words, instead of restricting the set of allowable predictors, it aims to avoid overfitting by requiring that the learned predictor 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 a predictor 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, and when designing the predictor it is desirable to also take into account the cost of communicating X. This is motivated by the fact that bandwidth and energy limitations on communication in networks and within multiprocessor systems 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 it (or features of it) are communicated over bandwidth-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 prediction in distributed and rate-constrained settings.

In this paper, we study an extension of the classical learning problem in which 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 as 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 predictor 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 predictor under this worst-case distribution. In this paper, we show that an analogous principle approximately holds for rate-constrained prediction. 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. 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 whereby we first find the optimal predictor ignoring the rate constraint, then compress and communicate the predictor 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. We show that when prediction 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 prediction is logarithmic loss, our results provide a new one-shot operational interpretation of the information bottleneck problem. A key technical ingredient in our results is the strong functional representation lemma (SFRL) developed in 4, which we use to design the optimal descriptor-estimator pair for the worst-case distribution.

Notation

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.

II. PROBLEM FORMULATION

We begin by reviewing the minimax approach to the classical learning problem [1].

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 predictor 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 predictors \mathcal{F} that consists of functions $f:\mathcal{X}\to\hat{\mathcal{Y}}$ (where the reconstruction alphabet $\hat{\mathcal{Y}}$ can be in general different from \mathcal{Y}) and a loss function $\ell:\hat{\mathcal{Y}}\times\mathcal{Y}\to\mathbb{R}$. The risk associated with a predictor f when the underlying joint distribution of X and Y is P is

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

ERM simply chooses the predictor $f_n \in \mathcal{F}$ with minimal risk under the empirical distribution P_n of the training samples. Recently, an alternative approach has emerged which seeks distributionally robust predictors. This approach learns a predictor by minimizing its worst-case risk over an ambiguity distribution set $\Gamma(P_n)$, i.e.,

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

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 predictors to a low-complexity class to prevent overfitting, in the minimax approach overfitting is prevented by explicitly requiring that the chosen predictor 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 Π

Learning:
$$\{(X_i,Y_i)\}_{i=1}^n \quad \Rightarrow P_n \quad \Rightarrow \Gamma(P_n) \quad \Rightarrow \text{Learner} \quad \Rightarrow \\ \text{Inference:} \quad X \quad \qquad \Rightarrow \hat{Y} = f_n(X) \\ \Rightarrow \quad \text{Inferrer } f_n \quad \Rightarrow \\ Y = f_n(X) \\ \Rightarrow \quad Y =$$

Fig. 1. Minimax approach to supervised learning.

B. Minimax Learning for Remote Prediction

In this paper, we extend the minimax learning approach to the setting in which the prediction 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 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 $R(e, P) \triangleq \mathsf{E}_P[|e(X)|]$ be the rate of the descriptor e and $L(e, f, P) \triangleq \mathsf{E}_P[\ell(f(e(X)), Y)]$ 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

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

Note that this cost function takes into account both the resultant statistical prediction 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.,

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

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.

Fig. 2. Minimax learning for remote prediction.

III. MAIN RESULTS

We first consider the case where Γ consists of a single distribution P, which may be the empirical distribution P_n as in ERM. Define the minimax risk-rate cost as

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

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

$$\begin{split} L_{\lambda}^* &\geq \inf_{P_{\hat{Y}|X}} \left(\mathsf{E} \left[\ell(\hat{Y}, Y) \right] + \lambda I(X; \hat{Y}) \right), \\ L_{\lambda}^* &\leq \inf_{P_{\hat{Y}|X}} \left(\mathsf{E} \left[\ell(\hat{Y}, Y) \right] + \lambda \left(I(X; \hat{Y}) + \log(I(X; \hat{Y}) + 1) + 5 \right) \right). \end{split}$$

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 it is large.

To prove Theorem [1] we use the strong functional representation lemma given in [4] (also see [5], [6]): 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

$$H(\hat{Y}|W) \le I(X;\hat{Y}) + \log(I(X;\hat{Y}) + 1) + 4.$$
 (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 an 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, we use the Poisson functional representation construction detailed in \mathbb{A} .

Proof of Theorem 1: Recall that $\hat{Y} = f(e(X,W),W)$. The lower bound follows from the fact that $I_P(X;\hat{Y}) \leq H_P(M) \leq E[|M|]$. 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

$$\mathsf{E}[|M|] \le H(\hat{Y}|W) + 1 \le I(X;\hat{Y}) + \log(I(X;\hat{Y}) + 1) + 5.$$

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

Remark 1. If we consider the logarithmic loss $\ell(\hat{y}, y) = -\log \hat{y}(y)$, where \hat{y} is a distribution over \mathcal{Y} , then the lower bound in Theorem 1 reduces to

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

which is the information bottleneck function [7]. Therefore the setting of remote prediction provides an approximate one-shot operational interpretation of the information bottleneck (up to a logarithmic gap). In [8], [9] it was shown that the asymptotic noisy source coding problem also provides an operational interpretation of the information bottleneck. Our operational interpretation, however, is more satisfying since the feature extraction problem originally considered in [7] is by nature one-shot.

We now extend Theorem 1 to the minimax setting.

Theorem 2. Suppose Γ is convex. Then

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

This result is related to minimax noisy source coding [10]. The main difference is that we consider the one-shot expected length instead of the asymptotic rate.

To prove this theorem, we first invoke a minimax result for relative entropy in [11] (which generalizes the redundancy-capacity theorem [12]). Then we apply the following refined version of the strong functional representation lemma that is proved in the proof of Theorem 1 in [4] (also see [5]).

Lemma 1. For any $P_{\hat{Y}|X}$ and $\tilde{P}_{\hat{Y}}$, there exists random variable W, and functions $k(x,w) \in \{1,2,\ldots\}$ and $\hat{y}(k,w)$ such that $\hat{y}(k(x,W),W) \sim P_{\hat{Y}|X}(\cdot|x)$, and

$$\mathsf{E}\left[\log k(x,W)\right] \le D\left(P_{\hat{Y}|X}(\cdot|x) \,\middle\|\, \tilde{P}_{\hat{Y}}\right) + 1.6. \tag{6}$$

We are now ready to prove Theorem 2.

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

$$\begin{split} L' &= \sup_{P \in \Gamma} \Big(\, \mathsf{E}_P \left[\ell(\hat{Y}, Y) \right] \\ &+ \lambda \left(I_P(X; \hat{Y}) + 2 \log(I_P(X; \hat{Y}) + 1) + 6 \right) \Big). \end{split}$$

Let

$$\begin{split} g(P, \tilde{P}_{\hat{Y}}) &= \mathsf{E}_{P} \left[\ell(\hat{Y}, Y) \right] + \lambda \bigg(\int D \Big(P_{\hat{Y}|X=x} \mathbin{\big\|} \tilde{P}_{\hat{Y}} \Big) dP(x) \\ &+ 2 \log \bigg(\int D \Big(P_{\hat{Y}|X=x} \mathbin{\big\|} \tilde{P}_{\hat{Y}} \Big) dP(x) + 1 \bigg) + 6 \bigg). \end{split}$$

Note that g is concave in P for fixed $\tilde{P}_{\hat{Y}}$ since $\mathsf{E}_P\left[\ell(\hat{Y},Y)\right]$ and $\int D\big(P_{\hat{Y}|X=x} \parallel \tilde{P}_{\hat{Y}}\big)dP(x)$ are linear in P. Also g is quasiconvex in $\tilde{P}_{\hat{Y}}$ for fixed P since $\int D\big(P_{\hat{Y}|X=x} \parallel \tilde{P}_{\hat{Y}}\big)dP(x)$ is convex in $\tilde{P}_{\hat{Y}}$, and is lower semicontinuous in $\tilde{P}_{\hat{Y}}$ since $D\big(P_{\hat{Y}|X=x} \parallel \tilde{P}_{\hat{Y}}\big)$ is lower semicontinuous with respect to the topology of weak convergence [13], and hence $\int D\big(P_{\hat{Y}|X=x} \parallel \tilde{P}_{\hat{Y}}\big)dP(x)$ is lower semicontinuous by Fatou's lemma.

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 11 (on g instead of relative entropy, and using Sion's minimax theorem 14 instead of Lemma 2 in 11) that if $\Gamma_{\hat{Y}}$ is uniformly tight, then there exists $P_{\hat{Y}}^* \in \overline{\Gamma_{\hat{Y}}}$ such that

$$\sup_{P\in\Gamma}g(P,\tilde{P}_{\hat{Y}}^*)=\sup_{P\in\Gamma}\inf_{\tilde{P}_{\hat{Y}}}g(P,\tilde{P}_{\hat{Y}})=L'.$$

If $\Gamma_{\hat{Y}}$ is not uniformly tight, then by Lemma 4 in [11], $\sup_{P\in\Gamma}\inf_{\tilde{P}_{\hat{Y}}}\int D\big(P_{\hat{Y}|X=x}\,\big\|\,\tilde{P}_{\hat{Y}}\big)dP(x)=\infty$, and hence $L'=\sup_{P\in\Gamma}\inf_{\tilde{P}_{\hat{Y}}}g(P,\tilde{P}_{\hat{Y}})=\infty$.

 $\sup_{P\in\Gamma}\inf_{\tilde{P}_{\hat{Y}}}g(P,\tilde{P}_{\hat{Y}})=\infty.$ Applying Lemma $\boxed{1}$ to $P_{\hat{Y}|X}$, $P_{\hat{Y}}^*$ we obtain W independent of X, random variable $K=k(X,W)\in\{1,2,\ldots\}$, and $\hat{Y}=\hat{y}(K,W)$ following the conditional distribution $P_{\hat{Y}|X}$, and

$$\mathsf{E}\left[\log K \,|\, X = x\right] \le D\left(P_{\hat{Y}|X} \,\|\, P_{\hat{Y}}^* \,|\, X = x\right) + 1.6$$

for any x. Then we use Elias delta code 15 for K to produce M. Note that the average length of the Elias delta code is upper bounded by $\log K + 2 \log (\log K + 1) + 1$. Hence, we have

$$\begin{split} \mathsf{E}_{P} \left[|M| \right] & \leq \mathsf{E}_{P} \left[\log K \right] + 2 \log \left(\mathsf{E}_{P} \left[\log K \right] + 1 \right) + 1 \\ & \leq \int D \left(P_{\hat{Y}|X=x} \, \middle\| \, P_{\hat{Y}}^{*} \right) \! dP(x) \\ & + 2 \log \left(\int D \left(P_{\hat{Y}|X=x} \, \middle\| \, P_{\hat{Y}}^{*} \right) \! dP(x) + 1 \right) + 6. \end{split}$$

Hence

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

Theorem 2 suggest that we can simplify the analysis of the risk-rate cost (2) $L_{\lambda} = \mathsf{E}_{P}\left[\ell(\hat{Y},Y)\right] + \lambda \,\mathsf{E}_{P}\left[|M|\right]$ by replacing the rate $\mathsf{E}_{P}\left[|M|\right]$ with the mutual information $I_{P}(X;\hat{Y})$. Define the *risk-information cost* as

$$\tilde{L}_{\lambda}(P_{\hat{Y}|X}, P) = \mathsf{E}_{P}\left[\ell(\hat{Y}, Y)\right] + \lambda I_{P}(X; \hat{Y}). \tag{7}$$

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

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

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

$$\tilde{L}_{\lambda}^* \le L_{\lambda}^* \le \tilde{L}_{\lambda}^* + 2\lambda \log(\lambda^{-1}\tilde{L}_{\lambda}^* + 1) + 7\lambda.$$

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 14, which gives the following proposition.

Proposition 1. Suppose \mathcal{X} , \mathcal{Y} and $\hat{\mathcal{Y}}$ are finite, Γ 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 \blacksquare means that in order to design a robust descriptor-estimator pair that work 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 Γ , we design the descriptor-estimator pair based on the worst-case distribution

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

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 Lemma 1 on $P_{\hat{Y}|X}$ and the induced distribution $P_{\hat{Y}}^*$ from $P_{\hat{Y}|X}$ and P^* .

IV. APPLICATIONS

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 $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$. Take Γ to be the set of distributions with the same first and second moments as given by the empirical distribution, i.e.,

$$\Gamma = \left\{ P_{\mathbf{X}Y} : \mathsf{E}[\mathbf{X}] = \boldsymbol{\mu}_{\mathbf{X}}, \, \mathsf{E}[Y] = \boldsymbol{\mu}_{Y}, \, \mathsf{Var}[\mathbf{X}] = \boldsymbol{\Sigma}_{\mathbf{X}}, \\ \mathsf{Var}[Y] = \sigma_{Y}^{2}, \, \mathsf{Cov}[\mathbf{X}, Y] = C_{\mathbf{X}Y} \right\}, \tag{9}$$

where $\mu_{\mathbf{X}}, \mu_{Y}, \Sigma_{\mathbf{X}}, \sigma_{Y}^{2}, C_{\mathbf{X}Y}$ are the corresponding statistics of the empirical distribution. The following proposition shows that P^{*} is Gaussian.

Proposition 2 (Linear regression with rate constraint). Consider mean-squared loss and define Γ as in [9]. Then the minimax risk-information cost [8] is

$$\tilde{L}_{\lambda}^{*} = \begin{cases} \sigma_{Y}^{2} - C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y} + \frac{\lambda}{2} \log \frac{2eC_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y}}{\lambda \log e} & \text{if } \frac{\lambda \log e}{2} < C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y} \\ \sigma_{Y}^{2} & \text{if } \frac{\lambda \log e}{2} \ge C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y}, \end{cases}$$
(10)

where the optimal P_{XY}^* is Gaussian with its mean and covariance matrix specified in 9, and the optimal estimate

$$\hat{Y} = \begin{cases} aC_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{X} + b + Z & \text{if } \frac{\lambda \log e}{2} < C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y} \\ \mu_Y & \text{if } \frac{\lambda \log e}{2} \ge C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y}, \end{cases}$$

where

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

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

Note that this setting does not satisfy the conditions in Proposition \blacksquare We directly analyze \blacksquare to obtain the optimal $P_{\mathbf{X}Y}^*$. Given the optimal $P_{\mathbf{X}Y}^*$, Theorem \blacksquare and Lemma \blacksquare can be used to construct the scheme. Operationally, $e_n(x,w)$ is a random quantizer of $aC_{\mathbf{X}Y}^T\Sigma_{\mathbf{X}}^{-1}\mathbf{X} + b$ such that the quantization noise follows $N(0, \sigma_Z^2)$. 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, $\mu_X=\mu_Y=0$, $\sigma_X^2=\sigma_Y^2=1$, $C_{XY}=0.95$ for the scheme constructed using the Poisson functional representation in [4], with the lower bound given by the minimax risk-information cost \tilde{L}_{λ}^* , and the upper bound given in Theorem 2.

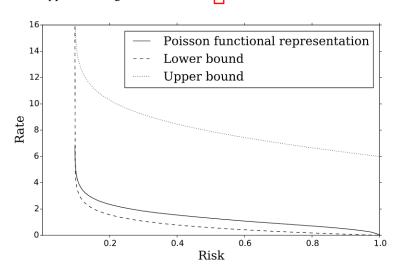


Fig. 3. Tradeoff between the rate and the risk in rate-constrained minimax linear regression.

Proof of Proposition 2 Without loss of generality, assume $\mu_{\mathbf{X}} = \mathbf{0}$ and $\mu_Y = 0$. We first prove " \leq " in (10). For this, fix $P_{\hat{Y}|\mathbf{X}}$ 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

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

$$I_{P}(\mathbf{X}; \hat{Y}) = h(\hat{Y}) - h(\hat{Y}|\mathbf{X})$$

$$\leq \frac{1}{2} \log \left(\frac{2C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y}}{\lambda \log e} \right).$$

Therefore.

$$\inf_{P_{\hat{Y}|X}} \sup_{P \in \Gamma} \left(\mathsf{E}_P \left[\ell(\hat{Y}, Y) \right] + \lambda I_P(X; \hat{Y}) \right) \le \mathsf{R.H.S.} \text{ of } \boxed{10}.$$

It can also be checked that the above relation holds when $\frac{\lambda \log e}{2} \ge C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y}$, and thus we have proved " \le " in (10). To prove " \geq " in (10), fix a Gaussian P_{XY} with its mean and covariance matrix specified in (9) and consider an arbitrary $P_{\hat{Y}|\mathbf{X}}$. We have

$$\begin{split} & \mathsf{E}_{P} \left[\ell(\hat{Y}, Y) \right] = \mathsf{E}_{P} \left[(Y - \hat{Y})^{2} \right] \\ & = \sigma_{Y}^{2} - C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y} + \mathsf{E}_{P} \left[(\hat{Y} - C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} \mathbf{X})^{2} \right], \text{ and } \\ & I_{P}(X; \hat{Y}) = I_{P} \left(C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} \mathbf{X}; \hat{Y} \right) \\ & \geq h \left(C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} \mathbf{X} \right) - h \left(C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} \mathbf{X} - \hat{Y} \right) \\ & \geq \frac{1}{2} \log C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y} - \frac{1}{2} \log \mathsf{E}_{P} \left[(\hat{Y} - C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} \mathbf{X})^{2} \right]. \end{split}$$

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

$$\begin{split} & \mathsf{E}_{P}\left[\ell(\hat{Y},Y)\right] + \lambda I_{P}(X;\hat{Y}) \\ & \geq \sigma_{Y}^{2} - C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y} + \frac{\lambda}{2} \log C_{\mathbf{X}Y}^{T} \Sigma_{\mathbf{X}}^{-1} C_{\mathbf{X}Y} + \gamma - \frac{\lambda \log \gamma}{2} \\ & \geq \mathsf{R.H.S. of (10)}, \end{split}$$

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 2.

The optimal scheme in the above example corresponds to compressing and communicating the minimax optimal rateunconstrained predictor $\bar{Y} = C_{\mathbf{X}Y}^T \Sigma_{\mathbf{X}}^{-1} (\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}}) + \mu_Y$, since the optimal \tilde{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 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

We assume $\mathcal{Y} = \hat{\mathcal{Y}} = \{1, \dots, k\}$ and \mathcal{X} are finite, $\ell(\hat{y}, y) = \mathbf{1}\{\hat{y} \neq y\}$, and Γ is closed and convex. The following proposition gives the minimax risk-information cost and the optimal estimator.

Proposition 3. Consider the setting described above. The minimax risk-information cost is given by

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

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

such that $\pi(y_1) = y_2$ and $P_{X,Y} \in \Gamma \Leftrightarrow P_{\tau(X),\pi(Y)} \in \Gamma$),

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

We can see that when $\lambda \to 0$, $P_{\hat{Y}|X}^*$ tends to the maximum a posteriori estimator (under \bar{P}^* , the worst-case distribution when $\lambda = 0$).

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

$$\begin{split} &\inf_{P_{\hat{Y}|X}^*} \tilde{L}_{\lambda}(P_{\hat{Y}|X}, P) \\ &= \inf_{P_{\hat{Y}|X}^*} \left(\mathsf{E}_P \left[\ell(\hat{Y}, Y) \right] + \lambda I_P(X; \hat{Y}) \right) \\ &= \inf_{P_{\hat{Y}|X}^*} \left(P\{\hat{Y} \neq Y\} + \lambda \inf_{\tilde{P}_{\hat{Y}}^*} \int D(P_{\hat{Y}|X=x} \parallel \tilde{P}_{\hat{Y}}) dP(x) \right) \\ &= \inf_{\tilde{P}_{\hat{Y}}, P_{\hat{Y}|X}^*} \left(P\{\hat{Y} \neq Y\} + \lambda \int D(P_{\hat{Y}|X=x} \parallel \tilde{P}_{\hat{Y}}) dP(x) \right) \\ &= 1 + \lambda \inf_{\tilde{P}_{\hat{Y}}, P_{\hat{Y}|X}^*} \mathsf{E}_P \left(\sum_y P_{\hat{Y}|X}(y|X) \left(\log \frac{P_{\hat{Y}|X}(y|X)}{\tilde{P}_{\hat{Y}}(y)} - \lambda^{-1} P_{Y|X}(y|X) \right) \right) \\ &= 1 + \lambda \inf_{\tilde{P}_{\hat{Y}}^*} \mathsf{E}_P \left(\sum_y P_{\hat{Y}|X}(y|X) \left(\log \frac{P_{\hat{Y}|X}(y|X)}{\tilde{P}_{\hat{Y}}(y)} - \lambda^{-1} P_{Y|X}(y|X) \tilde{P}_{\hat{Y}}(y|X) \right) \right) \\ &= 1 + \lambda \inf_{\tilde{P}_{\hat{Y}}^*} \mathsf{E}_P \left(\sum_y P_{\hat{Y}|X}(y|X) \left(\log \frac{P_{\hat{Y}|X}(y|X)}{2^{\lambda^{-1}P_{Y|X}(y|X)} \tilde{P}_{\hat{Y}}(y) / \sum_{y'} 2^{\lambda^{-1}P_{Y|X}(y'|X)} \tilde{P}_{\hat{Y}}(y') \right) - \log \sum_y 2^{\lambda^{-1}P_{Y|X}(y|X)} \tilde{P}_{\hat{Y}}(y) \right) \\ &\stackrel{(a)}{=} 1 + \lambda \inf_{\tilde{P}_{\hat{Y}}^*} \mathsf{E}_P \left(-\log \sum_y 2^{\lambda^{-1}P_{Y|X}(y|X)} \tilde{P}_{\hat{Y}}(y) \right), \end{split}$$

where (a) is due to that relative entropy is nonnegative, and equality is attained when $P_{\hat{Y}|X}(y|x) \propto 2^{\lambda^{-1}P_{Y|X}(y|X)}\tilde{P}_{\hat{Y}}(y)$. Next we consider the case in which Γ is symmetric. Consider the minimax rate-information cost

$$\tilde{L}_{\lambda}^{*} = \inf_{P_{\hat{Y}|X}} \sup_{P \in \Gamma} \tilde{L}_{\lambda}(P_{\hat{Y}|X}, P) = \inf_{P_{\hat{Y}|X}} \sup_{P \in \Gamma} \bigg(\mathsf{E}_{P} \left[\ell(\hat{Y}, Y) \right] + \lambda I_{P}(X; \hat{Y}) \bigg).$$

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

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

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

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

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

$$\begin{split} &\tilde{L}_{\lambda}^{*} = \inf_{P_{\hat{Y}|X} \text{ symm.}} \sup_{P \in \Gamma \text{ symm.}} \tilde{L}_{\lambda}(P_{\hat{Y}|X}, P) \\ &= \inf_{P_{\hat{Y}|X} \text{ symm.}} \sup_{P \in \Gamma \text{ symm.}} \left(\mathbb{E}_{P} \left[\ell(\hat{Y}, Y) \right] + \lambda I_{P}(X; \hat{Y}) \right) \\ &= \inf_{P_{\hat{Y}|X} \text{ symm.}} \sup_{P \in \Gamma \text{ symm.}} \left(P\{\hat{Y} \neq Y\} + \lambda (\log k - H_{P}(\hat{Y}|X)) \right) \\ &= 1 + \lambda \log k + \lambda \inf_{P_{\hat{Y}|X} \text{ symm.}} \sup_{P \in \Gamma \text{ symm.}} \mathbb{E}_{P} \left(\sum_{y} P_{\hat{Y}|X}(y|X) \left(\log P_{\hat{Y}|X}(y|X) - \lambda^{-1} P_{Y|X}(y|X) \right) \right) \\ &= 1 + \lambda \log k + \lambda \inf_{P_{\hat{Y}|X} \text{ symm.}} \sup_{P \in \Gamma \text{ symm.}} \mathbb{E}_{P} \left(\sum_{y} P_{\hat{Y}|X}(y|X) \log \frac{P_{\hat{Y}|X}(y|X)}{2^{\lambda^{-1}P_{Y|X}(y|X)} / \sum_{y'} 2^{\lambda^{-1}P_{Y|X}(y'|X)}} - \log \sum_{y} 2^{\lambda^{-1}P_{Y|X}(y|X)} \right) \\ &\geq 1 + \lambda \log k + \lambda \inf_{P_{\hat{Y}|X} \text{ symm.}} \sup_{P \in \Gamma \text{ symm.}} \mathbb{E}_{P} \left(-\log \sum_{y} 2^{\lambda^{-1}P_{Y|X}(y|X)} \right) \\ &= \sup_{P \in \Gamma \text{ symm.}} \left(1 + \lambda \log k - \lambda \mathbb{E}_{P} \log \sum_{y} 2^{\lambda^{-1}P_{Y|X}(y|X)} \right), \end{split}$$

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

$$1 + \lambda \log k - \lambda \, \mathsf{E}_P \log \sum_{y} 2^{\lambda^{-1} P_{Y|X}(y|X)} = \inf_{P_{\hat{Y}|X}} \left(P\{\hat{Y} \neq Y\} + \lambda (\log k - H_P(\hat{Y}|X)) \right)$$

is an infimum of affine functions of P, hence it is concave in P. Also it is symmetric about π and τ , hence

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

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

To show that the estimate-compress approach is not always optimal, let $\ell(\hat{y}, y) = \mathbf{1}\{\hat{y} \neq y\}$, $\mathcal{Y} = \mathcal{Y}_1 \cup \mathcal{Y}_2$, where $\mathcal{Y}_1 \cap \mathcal{Y}_2 = \emptyset$ and $|\mathcal{Y}_i| = k_i$ is finite. 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. By Proposition 3, the optimal risk-information cost is

$$1 - \lambda \log \max \left\{ \frac{1}{k_1} (2^{\lambda^{-1} q_1} - 1) + 1, \frac{1}{k_2} (2^{\lambda^{-1} q_2} - 1) + 1 \right\}, \tag{11}$$

and the optimal estimator is

$$P_{\hat{Y}|X_{1},X_{2}}^{*}(\hat{y}|x_{1},x_{2}) = \begin{cases} \frac{2^{\lambda^{-1}q_{1}}}{2^{\lambda^{-1}q_{1}+k_{1}-1}} & \text{if } \hat{y} = x_{1} \\ \frac{1}{2^{\lambda^{-1}q_{1}+k_{1}-1}} & \text{if } \hat{y} \in \mathcal{Y}_{1} \setminus \{x_{1}\} \\ 0 & \text{if } \hat{y} \in \mathcal{Y}_{2} \end{cases}$$
(12)

if $\frac{1}{k_1}(2^{\lambda^{-1}q_1}-1)+1\geq \frac{1}{k_2}(2^{\lambda^{-1}q_2}-1)+1$, and similar for the other case. Assume $q_1>q_2$, then the optimal MAP estimate is $Y=X_1$. An estimate-compress approach would either communicate a compressed version of $\bar{Y}=X_1$ as in (12), or output any element in \mathcal{Y}_2 (giving a risk $1-q_2k_2^{-1}$). The risk-information cost achieved by this approach is

$$\min\left\{1 - \lambda \log\left(\frac{1}{k_1}(2^{\lambda^{-1}q_1} - 1) + 1\right), 1 - q_2 k_2^{-1}\right\} = 1 - \lambda \log \max\left\{\frac{1}{k_1}(2^{\lambda^{-1}q_1} - 1) + 1, 2^{\lambda^{-1}q_2 k_2^{-1}}\right\}. \tag{13}$$

Now, if $k_1 \gg k_2$, the optimal rate constrained descriptor communicates a lossy version of X_2 instead, and the risk of estimate-compress in (13) is larger than (11).

Moreover, the gap between the rates needed by the two approaches for a fixed risk can be unbounded. Take $q_1 = 1 - q_2 = 2/3$, $k_2 = 2$, $k_1 \ge 15$. 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 compressing X_1 (by passing it through a symmetric channel with $P\{\hat{Y} = X_1\} = 1/2$) to achieve a risk 2/3, which 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 the optimal scheme, the lower bound obtained from the optimal risk-information tradeoff (11), the upper bound of the optimal rate by Theorem 1 and the risk-information tradeoff for the estimate-compress approach (13) for $q_1 = 1 - q_2 = 2/3$, $k_1 = 2^{32}$, $k_2 = 2$. Note that the optimal scheme is to perform time sharing (using common randomness) 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 using 0 bit with 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).

V. ACKNOWLEDGEMENTS

This work was partially supported by a gift from Huawei Technologies and by the Center for Science of Information (CSoI), an NSF Science and Technology Center, under grant agreement CCF-0939370.

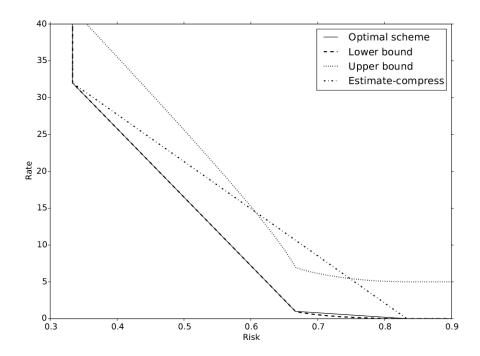


Fig. 4. Tradeoff between the rate and risk in rate-constrained minimax linear classification for the optimal scheme, lower bound [11], upper bounnd by Theorem [1] and estimate-compress approach [13].

REFERENCES

- [1] F. Farnia and D. Tse, "A minimax approach to supervised learning," in Advances in Neural Information Processing Systems, 2016, pp. 4240-4248.
- H. Namkoong and J. C. Duchi, "Variance-based regularization with convex objectives," in Advances in Neural Information Processing Systems, 2017, pp. 2975–2984.
- [3] J. Lee and M. Raginsky, "Minimax statistical learning and domain adaptation with Wasserstein distances," arXiv preprint arXiv:1705.07815, 2017.
- [4] C. T. Li and A. El Gamal, "Strong functional representation lemma and applications to coding theorems," in *Proc. IEEE Int. Symp. Inf. Theory*, June 2017, pp. 589–593.
- [5] 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.
- [6] M. Braverman and A. Garg, "Public vs private coin in bounded-round information," in *International Colloquium on Automata, Languages, and Programming*. Springer, 2014, pp. 502–513.
- [7] N. Tishby, F. C. Pereira, and W. Bialek, "The information bottleneck method," arXiv preprint physics/0004057, 2000.
- [8] P. Harremoës and N. Tishby, "The information bottleneck revisited or how to choose a good distortion measure," in *Information Theory*, 2007. ISIT 2007. IEEE International Symposium on. IEEE, 2007, pp. 566–570.
- [9] T. A. Courtade and T. Weissman, "Multiterminal source coding under logarithmic loss," *IEEE Transactions on Information Theory*, vol. 60, no. 1, pp. 740–761, 2014.
- [10] A. Dembo and T. Weissman, "The minimax distortion redundancy in noisy source coding," *IEEE Transactions on Information Theory*, vol. 49, no. 11, pp. 3020–3030, 2003.
- [11] D. Haussler, "A general minimax result for relative entropy," IEEE Transactions on Information Theory, vol. 43, no. 4, pp. 1276–1280, Jul 1997.
- [12] R. G. Gallager, "Source coding with side information and universal coding," Technical Report LIDS-P-937, MIT Laboratory for Information and Decision Systems, 1979.
- [13] E. Posner, "Random coding strategies for minimum entropy," IEEE Transactions on Information Theory, vol. 21, no. 4, pp. 388-391, Jul 1975.
- [14] M. Sion, "On general minimax theorems," Pacific Journal of mathematics, vol. 8, no. 1, pp. 171-176, 1958.
- [15] P. Elias, "Universal codeword sets and representations of the integers," IEEE Trans. Inf. Theory, vol. 21, no. 2, pp. 194-203, Mar 1975.