

Learning Joint Quantizers for Reconstruction and Prediction

Maxim Raginsky

Abstract—We consider the problem of empirical design of variable-rate quantizers for reconstruction and prediction. When a discriminative model (conditional distribution of the unobserved output given the observed input) is known or can be accurately estimated from a separate training set, we show that this problem reduces to designing a certain type of a generalized quantizer by means of empirical risk minimization on *unlabeled* input samples only. We derive a high-probability upper bound on the resulting expected performance of such a quantizer in terms of the training sample size and the complexity parameters of the reconstruction and the prediction problems. We also discuss two illustrative examples: binary classification with absolute loss and the information bottleneck.

I. INTRODUCTION

The problem of designing compressed representations of stochastic data sources subject to a reconstruction fidelity criterion is well-understood theoretically [1], [2]. There is also a wide variety of practical methods and algorithms which, though not necessarily optimal from a theoretical point of view, have been developed based on extensive domain knowledge and, show good empirical performance on many sources of interest (such as audio, images or video) [3]. On the other hand, the problem of designing compressed representations for inference or prediction is still lacking a theoretical framework comparable to that for lossy source coding. The main challenge is that there seems to be no universal way of quantifying the amount of “predictively relevant” information, although several proposals based either on decision-theoretic or on information-theoretic ideas have been put forward [4]–[8]. Moreover, one can easily come up with examples where the amount of information needed for reconstruction is much greater than what would be required for prediction, or vice versa.

It is often the case in practice that there is no readily available probabilistic model for the data source of interest. In such situations, one resorts to *empirical design* based on training data sampled from the source, and techniques from statistical learning theory can be used to analyze consistency and convergence to optimal performance (see [9] for applications to lossy source coding). In this paper, we apply learning-theoretic ideas to the problem of designing compressed representations that attempt to preserve enough information to enable both good reconstruction and good prediction. We show that this

problem can be reduced to designing a generalized quantizer with a modified fidelity criterion that incorporates both the reconstruction error term and the prediction error term, as well as the description length in bits. Our main result (Theorem 1) is a high-probability upper bound on the gap between the expected performance of an empirically learned generalized quantizer and the theoretical optimum. We also discuss two illustrative examples: binary classification with absolute loss and the information bottleneck.

II. PRELIMINARIES

We start by fixing some useful terminology and notation. We assume in the sequel that all spaces are standard Borel [10] (so that, in particular, all conditional distributions have regular versions), and that all maps under consideration are measurable with respect to appropriate σ -algebras.

A. Quantization

A variable-rate quantizer with input space X and reproduction space Z is a pair $q = (\varphi, \psi)$, where $\varphi : X \rightarrow S$ is the *encoder*, $\psi : S \rightarrow Z$ is the *decoder*, and $S \subseteq \{0, 1\}^*$ is a countable collection of binary strings satisfying Kraft’s inequality

$$\sum_{s \in S} 2^{-\text{len}(s)} \leq 1, \quad (1)$$

where $\text{len}(s)$ denotes the length of $s \in \{0, 1\}^*$ in bits. The quantizer q maps each point $x \in X$ to a reproduction $z = \psi \circ \varphi(x) \in Z$ (with a slight abuse of notation, we will also write $z = q(x)$). We are given a *distortion function* $d : X \times Z \rightarrow \mathbb{R}^+$, so that $d(x, z) = d(x, q(x))$ is the distortion incurred in representing x by its quantized version z . Also, $\text{len}(\varphi(x))$ is the *description length* of x . If q is fed with a random input X according to a probability law P , then the *expected distortion* and the *rate* of q are given by

$$D_P(q) \triangleq \mathbb{E}[d(X, q(X))] \text{ and } R_P(q) \triangleq \mathbb{E}[\text{len}(\varphi(X))]. \quad (2)$$

It is convenient [9], [11] to absorb the distortion and the rate into a single performance measure, the *Lagrangian*

$$L_P(q, \lambda) \triangleq D_P(q) + \lambda R_P(q) \quad (3)$$

$$= \mathbb{E}[d(X, q(X)) + \lambda \cdot \text{len}(\varphi(X))], \quad (4)$$

where $\lambda \geq 0$ is the Lagrange multiplier that controls the distortion-rate trade-off. Given a collection of quantizers \mathcal{Q} , the optimum Lagrangian performance on P is given by

$$L_P^*(\mathcal{Q}, \lambda) \triangleq \inf_{q \in \mathcal{Q}} L_P(q, \lambda). \quad (5)$$

The author is with the Department of Electrical and Computer Engineering and the Coordinated Science Laboratory, University of Illinois, Urbana, IL 61801, USA. E-mail: maxim@illinois.edu.

Research supported in part by DARPA under grant no. N66001-13-1-4004 and in part by NSF under CAREER award no. CCF-1254041.

It is easy to show that, for any quantizer $q = (\varphi, \psi)$ and a Lagrange multiplier λ , we can construct another quantizer $q' = (\varphi', \psi')$ with $\psi' = \psi$ and

$$\varphi'(x) \triangleq \arg \min_{s \in \mathcal{S}} \{d(x, \psi(s)) + \lambda \cdot \text{len}(s)\}, \quad (6)$$

so that $L_P(q') \leq L_P(q)$. Thus, we may always assume that our quantizers have (modified) nearest-neighbor encoders.

B. Prediction

In a statistical prediction problem, we have a random couple $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ with joint probability law \mathbb{P} , where X is the input (or feature) and Y is the output (or label). A *predictor* is a mapping $f : \mathcal{X} \rightarrow \mathcal{U}$, where \mathcal{U} is some *action space*. Thus, each input point $x \in \mathcal{X}$ is mapped to an action $u = f(x)$. We are given a *loss function* $\ell : \mathcal{X} \times \mathcal{Y} \times \mathcal{U} \rightarrow \mathbb{R}^+$, so that $\ell(x, y, u) = \ell(x, y, f(x))$ is the loss incurred when taking action u in response to the input-output pair (x, y) (note that the action u is allowed to depend only on the input x , so effectively we need to *predict* the corresponding input y without seeing it in order to settle on a good action, hence the term “prediction”). The *expected loss* of f is given by

$$\Lambda_{\mathbb{P}}(f) \triangleq \mathbb{E}[\ell(X, Y, f(X))]. \quad (7)$$

III. PROBLEM FORMULATION

We are interested in the following problem that combines the reconstruction and the prediction aspects: We have a random couple $(X, Y) \in \mathcal{X} \times \mathcal{Y}$, an action space \mathcal{U} , a collection \mathcal{Q} of quantizers with input and reproduction spaces both equal to \mathcal{X} , and a collection \mathcal{F} of predictors $f : \mathcal{X} \rightarrow \mathcal{U}$. Given a quantizer $q = (\varphi, \psi) \in \mathcal{Q}$ and a predictor $f \in \mathcal{F}$, we observe the input X , quantize it to $\hat{X} = q(X) = \psi \circ \varphi(X)$, and then take an action $U = f(\hat{X}) = f \circ q(X)$. Our goal is to jointly choose q and f to simultaneously guarantee small values of expected distortion, rate, and prediction loss. Once again, we adopt the Lagrangian viewpoint and quantify the performance of a pair $(q, f) \in \mathcal{Q} \times \mathcal{F}$ by the scalar quantity

$$\begin{aligned} \mathcal{L}_{\mathbb{P}}(q, f, \lambda, \mu) &\triangleq L_{\mathbb{P}_X}(q, \lambda) + \mu \cdot \Lambda_{\mathbb{P}}(f \circ q) \\ &\equiv \mathbb{E}[d(X, q(X)) + \lambda \cdot \text{len}(\varphi(X)) + \mu \cdot \ell(X, Y, f \circ q(X))], \end{aligned} \quad (8)$$

$$(9)$$

where $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+$ is the reconstruction distortion function, $\ell : \mathcal{X} \times \mathcal{Y} \times \mathcal{U} \rightarrow \mathbb{R}^+$ is the prediction loss function, and $\lambda, \mu \geq 0$ are the Lagrange multipliers. The optimal performance achievable by $\mathcal{Q} \times \mathcal{F}$ on \mathbb{P} is

$$\mathcal{L}_{\mathbb{P}}^*(\mathcal{Q}, \mathcal{F}, \lambda, \mu) \triangleq \inf_{q \in \mathcal{Q}} \inf_{f \in \mathcal{F}} \mathcal{L}_{\mathbb{P}}(q, f, \lambda, \mu). \quad (10)$$

A. Reduction to generalized quantization

We now reduce the above problem to a quantization problem with input space \mathcal{X} and reproduction space $\mathcal{Z} = \mathcal{X} \times \mathcal{U}$. Let $\mathbb{P}_{Y|X}$ denote the regular conditional distribution of Y given X , and define a new loss function $\tilde{\ell} : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}^+$ by

$$\tilde{\ell}(x, u) \triangleq \mathbb{E}[\ell(X, Y, u) | X = x] \quad (11)$$

$$\equiv \int_{\mathcal{Y}} \mathbb{P}_{Y|X}(dy|x) \ell(x, y, u) \quad (12)$$

(this is simply the conditional expectation of the loss $\ell(X, Y, u)$ w.r.t. the “discriminative model” $\mathbb{P}_{Y|X=x}$). Then for any q and f we can write

$$\begin{aligned} \mathcal{L}_{\mathbb{P}}(q, f, \lambda, \mu) &= \mathbb{E}[d(X, q(X)) + \lambda \cdot \text{len}(\varphi(X))] \\ &\quad + \mu \cdot \mathbb{E}\left[\mathbb{E}[\ell(X, Y, f \circ q(X)) | X]\right] \\ &= \mathbb{E}\left[d(X, q(X)) + \mu \cdot \tilde{\ell}(X, f \circ q(X))\right] + \lambda \cdot \mathbb{E}[\text{len}(\varphi(X))]. \end{aligned} \quad (13)$$

$$(14)$$

For any $q = (\varphi, \psi) \in \mathcal{Q}$ and $f \in \mathcal{F}$, we can define a new quantizer $\tilde{q} = (\tilde{\varphi}, \tilde{\psi})$ with

- input space \mathcal{X} ,
- reproduction space $\mathcal{Z} = \mathcal{X} \times \mathcal{U}$,
- encoder $\tilde{\varphi} = \varphi$, and
- decoder $\tilde{\psi} : \mathcal{S} \rightarrow \mathcal{Z}$ with $\tilde{\psi}(s) = (\psi(s), f \circ \psi(s))$.

If we also define a distortion function $d_{\mu} : \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R}^+$ by

$$d_{\mu}(x, z) = d_{\mu}(x, (\hat{x}, u)) \triangleq d(x, \hat{x}) + \mu \cdot \tilde{\ell}(x, u), \quad (15)$$

we can recognize the expression in (14) as the Lagrangian

$$L_{\mathbb{P}_X}(\tilde{q}, \lambda; \mu) \triangleq \mathbb{E}[d_{\mu}(X, \tilde{q}(X))] + \lambda \cdot \mathbb{E}[\text{len}(\tilde{\varphi}(X))] \quad (16)$$

of the *generalized quantizer* $\tilde{q}(X) = (q(X), f \circ q(X))$.

In a nutshell, if we let $\tilde{\mathcal{Q}}$ denote the collection of all generalized quantizers \tilde{q} of the above form, then we have

$$\mathcal{L}_{\mathbb{P}}^*(\mathcal{Q}, \mathcal{F}, \lambda, \mu) = L_{\mathbb{P}_X}^*(\tilde{\mathcal{Q}}, \lambda; \mu). \quad (17)$$

Note that we have effectively removed the label Y from the picture: it enters indirectly, through the “posterior expected loss” $\tilde{\ell}(x, u) = \mathbb{E}[\ell(X, Y, u) | X = x]$. This reduction is conceptually similar to indirect rate-distortion problems (see [12] and references therein).

IV. EMPIRICAL QUANTIZER DESIGN

In this paper, we are interested in the situation where the joint distribution \mathbb{P} is unknown, and we must instead *learn* a suitable quantizer-predictor pair on the basis of a training sequence $\{(X_i, Y_i)\}_{i=1}^m$ of i.i.d. samples from \mathbb{P} . Of course, when we do not know \mathbb{P} , the above reduction to generalized quantization does not hold – indeed, we cannot compute $\tilde{\ell}(\cdot, \cdot)$ exactly since we do not know $\mathbb{P}_{Y|X}$. However, we may first learn a good discriminative model from a large labeled sample, and then use this learned model in a plug-in fashion on another sample when designing our quantizer and predictor. Moreover, once we have obtained such a discriminative model, we no longer need labeled data – our problem reduces to empirical quantizer design (albeit with an augmented reproduction space), which can be done only with unlabeled input samples.

This may seem like a roundabout way of learning a quantizer-predictor pair: first we learn a very good randomized predictor, and then we effectively forget (or marginalize away) what we have learned in order to design a quantizer and an inferior predictor. However, this separation-based approach makes sense as soon as we consider the case when the

feature X is first encoded into a binary representation at some remote observation station, and then this binary representation is transmitted to a decision-maker who will use it to both reconstruct the feature and to take a decision based on the reconstruction. In this setting, the encoder has strictly better information than the decision-maker, which is actually helpful, since we can guarantee that the quantized representation of the feature retains as much of this information as possible for a given choice of the Lagrange multipliers.

From now on we will assume that the conditional distribution $\mathbb{P}_{Y|X}$ is known (at least to the extent needed to compute $\tilde{\ell}$), and focus on generalized quantizer design with *unlabeled* input samples. Given n i.i.d. samples X_1, \dots, X_n from the marginal distribution \mathbb{P}_X of the input variable (or feature) X , we seek a solution to the following Empirical Risk Minimization (ERM) problem:

$$\tilde{q}_n \equiv (\tilde{q}_n, \tilde{f}_n) = \arg \min_{\tilde{q} \in \tilde{\mathcal{Q}}} \widehat{L}_n(\tilde{q}, \lambda; \mu), \quad (18)$$

where we have defined the empirical Lagrangian

$$\widehat{L}_n(\tilde{q}, \lambda; \mu) = \frac{1}{n} \sum_{i=1}^n \{d_\mu(X_i, \tilde{q}(X_i)) + \lambda \cdot \ln(\tilde{\varphi}(X_i))\} \quad (19)$$

for any quantizer $\tilde{q} = (q, f \circ q)$.

V. THE MAIN RESULT AND SOME EXAMPLES

Our main result gives a high-probability bound on the generalization performance of the ERM solution (18):

Theorem 1. *Suppose that both the reconstruction distortion d and the prediction loss ℓ are uniformly bounded:*

$$d_{\max} \triangleq \sup_{x \in \mathcal{X}} \sup_{\hat{x} \in \mathcal{X}} d(x, \hat{x}) < +\infty \quad (20)$$

$$\ell_{\max} \triangleq \sup_{x \in \mathcal{X}} \sup_{y \in \mathcal{Y}} \sup_{u \in \mathcal{U}} \ell(x, y, u) < +\infty. \quad (21)$$

Let \mathcal{A} denote the collection of all sets of the form

$$\left\{ x \in \mathcal{X} : d(x, \hat{x}) + \mu \tilde{\ell}(x, u) \leq r \right\} \quad (22)$$

for all $\hat{x} \in \mathcal{X}$, $u \in \mathcal{U}$, $r > 0$, and suppose that \mathcal{A} is a Vapnik–Chervonenkis (VC) class with VC dimension V [13]. Let \mathcal{Q} be the collection of all quantizers on \mathcal{X} , and let \mathcal{F} be the collection of all mappings $f : \mathcal{X} \rightarrow \mathcal{U}$. Then there exists some absolute constant $C > 0$, such that, for any $\lambda > 0$ and $\mu \geq 0$,

$$\begin{aligned} \mathcal{L}_{\mathbb{P}}(\tilde{q}_n, \tilde{f}_n, \lambda, \mu) &\leq \mathcal{L}_{\mathbb{P}}^*(\mathcal{Q}, \mathcal{F}, \lambda, \mu) \\ &\quad + T \left(C \sqrt{\frac{VN \log N}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}} \right) \end{aligned} \quad (23)$$

with probability at least $1 - \delta$, where

$$T = 3(d_{\max} + \mu \ell_{\max}) \quad \text{and} \quad N = 2^{2T/3\lambda}. \quad (24)$$

Remark 1. When $\ell \equiv 0$, we recover existing finite-sample bounds for empirically optimal variable-rate quantizers [9].

Remark 2. As we had mentioned earlier, the discriminative model $\mathbb{P}_{Y|X}$ must be learned separately. For instance, suppose we use an independent labeled sample of size m to construct an estimate $\tilde{\ell}_m$ of $\tilde{\ell}$, such that $\|\tilde{\ell} - \tilde{\ell}_m\|_{\infty} \leq \varepsilon_m$ with probability at least $1 - \delta$. Then we can plug this estimate into the ERM procedure (18) as a proxy for $\tilde{\ell}$. As a consequence, the error bound of Theorem 1 will pick up an extra term $\mu \varepsilon_m$, and the new bound will hold with probability at least $1 - 2\delta$.

The proof of the theorem is given in the next section. Here, we discuss two specific examples.

A. Binary classification with absolute loss

Let \mathcal{X} be a compact subset of \mathbb{R}^d , $\mathcal{Y} = \{0, 1\}$, and $\mathcal{U} = [0, 1]$. The reconstruction fidelity is measured by the squared Euclidean distance, $d(x, \hat{x}) = \|x - \hat{x}\|^2$, whereas the prediction error is measured by the absolute loss $\ell(x, y, u) = |y - u|$ (notice that the loss does not directly depend on the feature x). The boundedness conditions of Theorem 1 are automatically satisfied. Defining the usual regression function $\eta(x) \triangleq \mathbb{E}[Y|X = x] \equiv \mathbb{P}(Y = 1|X = x)$ [13], we can write

$$\tilde{\ell}(x, u) = u + (1 - 2u)\eta(x). \quad (25)$$

Thus, we are interested in bounding the VC dimension of the collection \mathcal{A} of all sets of the form

$$\left\{ x \in \mathcal{X} : \|x - \hat{x}\|^2 + \mu(u + (1 - 2u)\eta(x)) \leq r \right\} \quad (26)$$

for all $\hat{x} \in \mathcal{X}$, $u \in [0, 1]$, $r > 0$. The following lemma follows from a well-known result that an m -dimensional space of real-valued functions is a VC-subgraph class with VC dimension $\leq m$ [13, Theorem 13.9]:

Lemma 1. *Suppose that η can be represented as a linear combination $\sum_{j=1}^k \alpha_j \xi_j$ of k linearly independent functions ξ_1, \dots, ξ_k on \mathcal{X} . Then $V(\mathcal{A}) \leq d + 2k + 3$.*

B. Information bottleneck

Again, let \mathcal{X} be a compact subset of \mathbb{R}^d , let \mathcal{Y} be arbitrary, and suppose that all conditional distributions $\mathbb{P}_{Y|X=x}$, $x \in \mathcal{X}$, are dominated by a given σ -finite measure ν on \mathcal{Y} . We again take the squared Euclidean distance $d(x, \hat{x}) = \|x - \hat{x}\|^2$ as our reconstruction criterion. As our action space \mathcal{U} we take a suitable subset of the space of all ν -densities $u : \mathcal{Y} \rightarrow \mathbb{R}^+$ with $\int u d\nu = 1$, and consider the prediction loss of the form

$$\ell(x, y, u) = \phi \left(\frac{u(y)}{p_x(y)} \right), \quad (27)$$

where ϕ is a convex function on \mathbb{R}^+ and where $p_x \triangleq d\mathbb{P}_{Y|X=x}/d\nu$. Then

$$\tilde{\ell}(x, u) = \int_{\mathcal{Y}} \nu(dy) p_x(y) \phi \left(\frac{u(y)}{p_x(y)} \right) \equiv D_\phi(u||p_x), \quad (28)$$

which we recognize as the ϕ -divergence between the densities u and p_x [14], [15]. For example, if we let $\phi(t) = t \log t$, then we get the Kullback–Leibler divergence

$$\tilde{\ell}(x, u) = D(p_x||u); \quad (29)$$

with $\phi(t) = (1 - \sqrt{t})^2$, we get the squared Hellinger distance

$$\tilde{\ell}(x, u) = \int_{\mathcal{Y}} \nu(dy) \left(\sqrt{p_x(y)} - \sqrt{u(y)} \right)^2 = H^2(p_x, u). \quad (30)$$

Loss functions of this type underlie the so-called information bottleneck method [4], [6] for constructing compressed representations of stochastic data sources that retain as much “predictively relevant” information as possible (see also [8] for a different application of ϕ -divergences in the context of surrogate loss functions for binary classification). In fact, if we ignore reconstruction entirely by letting $d(\cdot, \cdot) = 0$ and set $\lambda \equiv 1$, then we end up with an optimization problem of the form

$$\inf_{q, f} \mathbb{E} \left\{ \text{len}(\varphi(X)) + \mu \cdot D_\phi(f_{q(X)} \| p_X) \right\}, \quad (31)$$

where, for notational consistency, we wrote $f_{q(x)}$ for the image of $q(x)$ under $f : \mathcal{X} \rightarrow \mathcal{U}$. The Shannon-theoretic counterpart of this operational problem,

$$\inf_U \{ I(X; U) + \mu \mathbb{E} [D(U \| p_X)] \}, \quad (32)$$

where I is the mutual information [16] and the infimum is over all \mathcal{U} -valued random variables (random densities) jointly distributed with X , yields the original information bottleneck formulation. By including a data reconstruction fidelity (distortion) term, we obtain a generalization of information bottleneck that can be used to design optimal compressed representations for joint reconstruction and prediction. (See also [7] for related results on quantizers with a fixed finite number of levels.)

The boundedness condition on d is satisfied by compactness of \mathcal{X} , whereas the boundedness of ℓ can be ensured via appropriate restrictions on the likelihood ratios u/p_x and the function ϕ . We seek an upper bound on the VC dimension of the collection \mathcal{A} of all sets of the form

$$\{x \in \mathcal{X} : \|x - \hat{x}\|^2 + \mu D_\phi(u \| p_x) \leq r\} \quad (33)$$

for all $\hat{x} \in \mathcal{X}$, $u \in \mathcal{U}$, $r > 0$. The following lemma can be proved using similar methods as Lemma 1:

Lemma 2. *Let \mathcal{U} be chosen in such a way that one can find k linearly independent functions ξ_1, \dots, ξ_k on \mathcal{X} and a sequence Φ_1, Φ_2, \dots of functions in $L^1(\nu)$ with the following property: any function of the form $h(x, y) = p_x(y)\phi(u(y)/p_x(y))$ for some $u \in \mathcal{U}$ admits a bilinear representation*

$$h(x, y) = \sum_{j=1}^k \sum_{m=1}^{\infty} c_{jm} \theta_j(x) \Phi_m(y),$$

such that the sequence $\{c_{jm} \|\Phi_m\|_{L^1(\nu)}\}_{m=1}^{\infty}$ is summable for every j . Then $V(\mathcal{A}) \leq d + k + 2$.

VI. THE PROOF OF THEOREM 1

We will need the following lemma, similar in spirit to Lemma A.3 in [17] (see also Lemma 10 in [9]):

Lemma 3. *Given two positive integers N, L , let $\tilde{\mathcal{Q}}(N, L) \subseteq \tilde{\mathcal{Q}}$ denote the collection of all quantizers $\tilde{q} = (\tilde{\varphi}, \tilde{\psi})$ with modified nearest-neighbor encoders*

$$\tilde{\varphi}(x) = \arg \min_{s \in \mathcal{S}} \left\{ d_\mu(x, \tilde{\psi}(s)) + \lambda \cdot \text{len}(s) \right\}, \quad (34)$$

such that $|\mathcal{S}| \leq N$ and $\text{len}(s) \leq L$ for all $s \in \mathcal{S}$. Then for any probability distribution P on \mathcal{X} we have

$$L_P^*(\tilde{\mathcal{Q}}, \lambda; \mu) = L_P^*(\tilde{\mathcal{Q}}(N, L), \lambda; \mu) \quad (35)$$

with

$$\log_2 N = L = \frac{2(d_{\max} + \mu\lambda_{\max})}{\lambda}. \quad (36)$$

Proof: Given P , let \tilde{q}^* with encoder $\tilde{\varphi}^* : \mathcal{X} \rightarrow \mathcal{S}$ and decoder $\tilde{\psi}^* : \mathcal{S} \rightarrow \mathcal{X} \times \mathcal{U}$ achieve (or come arbitrarily close to) the optimum performance on the left-hand side of (35). Let s_0 denote the shortest binary string in \mathcal{S} , i.e., $\text{len}(s_0) = \min_{s \in \mathcal{S}} \text{len}(s)$. Since $\tilde{\varphi}^*$ is a modified nearest-neighbor encoder, for any $s \in \mathcal{S}$ and any $x \in \mathcal{X}$ such that $s = \tilde{\varphi}^*(x)$ we have

$$d_\mu(x, \tilde{\psi}^*(s)) + \lambda \cdot \text{len}(s) \leq d_\mu(x, \tilde{\psi}^*(s_0)) + \lambda \cdot \text{len}(s_0). \quad (37)$$

Thus, for any $s \in \mathcal{S}$,

$$\text{len}(s) \leq \frac{d_{\max} + \mu\ell_{\max}}{\lambda} + \text{len}(s_0). \quad (38)$$

On the other hand, consider a (very) suboptimal quantizer \tilde{q}_0 with encoder $\tilde{\varphi}_0(x) = \mathbf{e}$ (the empty binary string of length 0) and decoder $\psi_0(\mathbf{e}) = \hat{x}_0 \times u_0$ for some arbitrary choice of $\hat{x}_0 \in \mathcal{X}$ and $u_0 \in \mathcal{U}$. Then $\tilde{q}_0 \in \tilde{\mathcal{Q}}$, and

$$\lambda \cdot \text{len}(s_0) \leq L_P^*(\tilde{\mathcal{Q}}, \lambda; \mu) \leq L_P(\tilde{q}_0, \lambda; \mu) \leq d_{\max} + \mu\ell_{\max}. \quad (39)$$

Consequently, for any $s \in \mathcal{S}$ we have

$$\text{len}(s) \leq \frac{2(d_{\max} + \mu\ell_{\max})}{\lambda}. \quad (40)$$

Since the binary strings in \mathcal{S} must satisfy Kraft's inequality, we have

$$1 \geq \sum_{s \in \mathcal{S}} 2^{-\text{len}(s)} \geq |\mathcal{S}| 2^{-2(d_{\max} + \mu\ell_{\max})/\lambda}. \quad (41)$$

Thus, the optimal quantizer \tilde{q}^* has no more than $N = 2^{2(d_{\max} + \mu\ell_{\max})/\lambda}$ levels, and the description length of any $x \in \mathcal{X}$ does not exceed $L = 2(d_{\max} + \mu\ell_{\max})/\lambda$. ■

In view of Lemma 3, both the empirically optimal quantizer \tilde{q}_n and the theoretically optimal quantizer \tilde{q}^* (i.e., the one that achieves $L_{\mathbb{P}_X}^*(\lambda; \mu)$) belong to the class $\tilde{\mathcal{Q}}(N, L)$ with N and L given in (36). That is,

$$\hat{L}_n(\tilde{q}_n, \lambda; \mu) = \inf_{\tilde{q} \in \tilde{\mathcal{Q}}(N, L)} L_{\hat{\mathbb{P}}_n}(\tilde{q}, \lambda; \mu), \quad (42)$$

$$L_{\mathbb{P}_X}(\tilde{q}^*, \lambda; \mu) = \inf_{\tilde{q} \in \tilde{\mathcal{Q}}(N, L)} L_{\mathbb{P}_X}(\tilde{q}, \lambda; \mu), \quad (43)$$

where $\hat{P}_n = (1/n) \sum_{i=1}^n \delta_{X_i}$ is the empirical distribution induced by the sample X^n . By a standard argument we obtain

$$\begin{aligned} & L_{\mathbb{P}_X}(\tilde{q}_n, \lambda; \mu) - L_{\mathbb{P}_X}(\tilde{q}^*, \lambda; \mu) \\ & \leq 2 \sup_{\tilde{q} \in \tilde{\mathcal{Q}}(N, L)} \left| L_{\hat{P}_n}(\tilde{q}, \lambda; \mu) - L_{\mathbb{P}_X}(\tilde{q}, \lambda; \mu) \right|. \end{aligned} \quad (44)$$

Since all quantizers under consideration have modified nearest-neighbor encoders, for any distribution P on X and any $\tilde{q} = (\tilde{\varphi}, \tilde{\psi})$ we have $L_P(\tilde{q}, \lambda; \mu) = \mathbb{E}_P[g_{\tilde{q}}(X)]$, where

$$g_{\tilde{q}}(x) \triangleq \min_{s \in \mathsf{S}} \left\{ d_\mu(x, \tilde{\psi}(s)) + \lambda \cdot \text{len}(s) \right\}. \quad (45)$$

Let \mathcal{G} denote the class of all functions $g_{\tilde{q}}$ as \tilde{q} ranges over $\tilde{\mathcal{Q}}(N, L)$. Denoting the supremum on the right-hand side of (44) by Δ_n , we can write

$$\Delta_n = \sup_{g \in \mathcal{G}} \left| \mathbb{E}_{\hat{P}_n} g - \mathbb{E}_{\mathbb{P}_X} g \right|. \quad (46)$$

Any g in \mathcal{G} is bounded between 0 and $T \triangleq 3(d_{\max} + \mu \ell_{\max})$, so the uniform deviation Δ_n has bounded differences with $c_1 = \dots = c_n = T/n$. Thus, by McDiarmid's inequality [13],

$$\mathcal{L}_{\mathbb{P}_X}(\tilde{q}_n, \lambda, \mu) \leq \mathcal{L}_{\mathbb{P}_X}^*(\tilde{q}^*, \lambda, \mu) + 2\mathbb{E}\Delta_n + T\sqrt{\frac{2\log(1/\delta)}{n}}, \quad (47)$$

with probability at least $1 - \delta$. It remains to upper-bound the expectation $\mathbb{E}\Delta_n$.

Using the identity $\mathbb{E}W = \int_0^T \mathbb{P}(W > t) dt$ for any random variable W taking values in $[0, T]$ a.s., we can write

$$\mathbb{E}\Delta_n = \mathbb{E} \left[\sup_{g \in \mathcal{G}} \left| \int_0^T \left(\hat{P}_n(g > t) - \mathbb{P}_X(g > t) \right) dt \right| \right] \quad (48)$$

$$\leq T \mathbb{E} \left[\sup_{t > 0} \sup_{g \in \mathcal{G}} \left| \hat{P}_n(g > t) - \mathbb{P}_X(g > t) \right| \right] \quad (49)$$

$$= T \mathbb{E} \left[\sup_{A \in \tilde{\mathcal{A}}_N} \left| \hat{P}_n(A) - \mathbb{P}_X(A) \right| \right], \quad (50)$$

where $\tilde{\mathcal{A}}_N$ denotes the collection of all indicator functions of the form $\mathbf{1}_{\{g > t\}}$ for all $g \in \mathcal{G}$ and all $t > 0$. A standard estimate from empirical process theory [18] then gives

$$\mathbb{E}\Delta_n \leq CT\sqrt{\frac{V(\tilde{\mathcal{A}}_N)}{n}} \quad (51)$$

for some universal constant $C > 0$, where $V(\cdot)$ denotes the VC dimension. Thus, we need to bound $V(\tilde{\mathcal{A}}_N)$.

To that end, consider an arbitrary $g = g_{\tilde{q}} \in \mathcal{G}$ for some $\tilde{q} \in \tilde{\mathcal{Q}}(N, L)$. For any $x \in \mathsf{X}$, $g(x) > t$ if and only if

$$d_\mu(x, \tilde{\psi}(s)) > t - \lambda \text{len}(s), \quad \forall s \in \mathsf{S}. \quad (52)$$

Since $|\mathsf{S}| \leq N$, there are two possibilities to consider for the event $E = \{g_{\tilde{q}} > t\}$:

- 1) If $t - \lambda \text{len}(s) \geq 0$ for at least one $s \in \mathsf{S}$, then E is an intersection of the complements of at most N sets from the collection \mathcal{A} of all sets of the form (22).

- 2) If $t < \lambda \text{len}(s)$ for all $s \in \mathsf{S}$, then $E = \mathsf{X}$.

This means that $\tilde{\mathcal{A}}_N \subseteq \bar{\mathcal{A}}_N \cup \{\mathsf{X}\}$, where $\bar{\mathcal{A}}_N$ denotes the collection of all intersections of N complements of sets from \mathcal{A} . Then, using Theorem 1.1 in [19], we have

$$V(\tilde{\mathcal{A}}_N) \leq 3VN \log(4N), \quad (53)$$

where V is the VC dimension of the class \mathcal{A} . Plugging the estimate (53) into (51) and using the resulting bound in (47), we conclude that (23) holds with probability at least $1 - \delta$, and the theorem is proved.

REFERENCES

- [1] T. Berger, *Rate Distortion Theory*. Englewood Cliffs, NJ: Prentice-Hall, 1971.
- [2] R. M. Gray, *Source Coding Theory*. Boston: Kluwer, 1990.
- [3] A. Gersho and R. M. Gray, *Vector Quantization and Signal Compression*. Boston: Kluwer, 1992.
- [4] N. Tishby, F. Pereira, and W. Bialek, "The information bottleneck method," in *Proc. 37th Annu. Allerton Conf. Communication, Control and Computing*, 1999, pp. 368–377.
- [5] A. Banerjee, S. Merugu, I. S. Dhillon, and J. Ghosh, "Clustering with Bregman divergences," *J. Machine Learning Res.*, vol. 6, pp. 1705–1749, 2005.
- [6] P. Harremoës and N. Tishby, "The Information Bottleneck revisited or how to choose a good distortion measure," in *Proc. IEEE Int. Symp. Inform. Theory*, 2007, pp. 566–570.
- [7] S. Lazebnik and M. Raginsky, "Supervised learning of quantizer codebooks by information loss minimization," *IEEE Trans. Pattern Anal. Machine Intelligence*, vol. 31, no. 7, pp. 1294–1309, July 2009.
- [8] X. Nguyen, M. J. Wainwright, and M. I. Jordan, "On surrogate loss functions and f -divergences," *Ann. Statist.*, vol. 37, no. 2, pp. 876–904, 2009.
- [9] T. Linder, "Learning-theoretic methods in vector quantization," in *Principles of Nonparametric Learning*, L. Györfi, Ed. New York: Springer-Verlag, 2001.
- [10] K. R. Parthasarathy, *Probability Measures on Metric Spaces*. New York: Academic Press, 1967.
- [11] P. A. Chou, T. Lookabaugh, and R. M. Gray, "Entropy-constrained vector quantization," *IEEE Trans. Acoust. Speech Signal Processing*, vol. 37, no. 1, pp. 31–42, January 1989.
- [12] H. S. Witsenhausen, "Indirect rate distortion problems," *IEEE Trans. Inform. Theory*, vol. 26, no. 5, pp. 518–521, 1980.
- [13] L. Devroye, L. Györfi, and G. Lugosi, *A Probabilistic Theory of Pattern Recognition*. New York: Springer-Verlag, 1996.
- [14] I. Csiszár, "Information-type measures of difference of probability distributions and indirect observations," *Stud. Sci. Math. Hung.*, vol. 2, pp. 299–318, 1967.
- [15] F. Liese and I. Vajda, "On divergences and informations in statistics and information theory," *IEEE Trans. Inform. Theory*, vol. 52, no. 10, pp. 4394–4412, October 2006.
- [16] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. New York: Wiley, 2006.
- [17] M. Raginsky, "Joint universal lossy coding and identification of stationary mixing sources with general alphabets," *IEEE Trans. Inform. Theory*, vol. 55, no. 5, pp. 1945–1960, May 2009.
- [18] R. M. Dudley, "Central limit theorems for empirical measures," *Ann. Probab.*, vol. 6, pp. 898–929, 1978.
- [19] A. van der Vaart and J. A. Wellner, "A note on bounds for VC dimensions," in *High Dimensional Probability V: The Luminy Volume*, C. Houdré, V. Koltchinskii, D. M. Mason, and M. Peligrad, Eds. Beachwood, OH: Institute of Mathematical Statistics, 2009, pp. 103–107.