

Shannon meets Blackwell and Le Cam: channels, codes, and statistical experiments

Maxim Raginsky

Abstract—The Blackwell–Le Cam decision theory provides an approximation framework for statistical experiments in terms of expected risks of optimal decision procedures. The Blackwell partial order formalizes an intuitive notion of which experiment of a given pair is “more informative” for the purposes of inference. The Le Cam deficiency is an approximation measure for any two statistical experiments (with the same parameter space), and it tells us how much we will lose if we base our decisions on one experiment rather than another. In this paper, we develop an extension of the Blackwell–Le Cam theory, starting from a partial ordering for channels introduced by Shannon. In particular, we define a new approximation measure for channels, which we call the Shannon deficiency, and use it to prove an approximation theorem for channel codes that extends an earlier result of Shannon. We also construct a broad class of deficiency-like measures for channels based on generalized divergences, relate them to several alternative notions of capacity, and prove new upper and lower bounds on the Le Cam deficiency.

I. INTRODUCTION

In two seminal papers written in the early 1950’s [1], [2], David Blackwell has introduced the concept of “comparison of statistical experiments” based on expected losses of statistical decision procedures. This work has sparked a great deal of interest and led to a number of alternative comparison criteria, including information-theoretic ones [3], [4]. Le Cam [5] has extended Blackwell’s theory to an *approximation* approach, which has since then been developed into a comprehensive theory of experiments (cf. [6], [7] for thorough expositions). This theory has led to deep results in mathematical statistics, such as the proof of asymptotic equivalence of nonparametric regression and function filtering in white Gaussian noise [8].

In information-theoretic terms, a statistical experiment is just a noisy communication channel with an uncoded input [9]; all the processing is done at the output. While this setup is the right one for statistics, in information theory we are also interested in such things as coding or modulation, which are interposed between the source and the channel. A coding/decoding comparison criterion for channels has been introduced in a short paper of Shannon [10]. Although Shannon’s comparison criterion bears certain similarities to Blackwell’s, it appears to have been developed independently (in fact, there is no reference to Blackwell in [10]) from a random coding perspective.

The present paper makes the following contributions:

- Following Le Cam [5], we augment the Shannon partial order with a norm-based *approximation* criterion. Thus,

This work was supported in part by DARPA under project KECOM.

The author is with the Department of Electrical and Computer Engineering, Duke University, Durham, NC, 27708, USA. E-mail: m.raginsky@duke.edu.

instead of asking which of two channels is “less noisy” or “more informative,” we ask how well one channel can be *approximated* by another. By allowing randomization at the input, as well as some shared randomness between the input and the output terminals, we define a new approximation measure for channels based on the Shannon ordering. In parallel to Le Cam’s terminology, we call it the *Shannon deficiency* of one channel w.r.t. another.

- Using this new notion of deficiency, we extend the result of Shannon [10], which says that for any code on a more noisy channel we can find a code on a less noisy channel that does at least as well, to an *arbitrary* pair of channels, where the notion of “doing at least as well” is replaced by “doing at least as well up to ε .”
- The Le Cam and Shannon deficiencies are based on the total variation distance. By replacing the total variation with any generalized divergence between probability distributions that obeys a data processing inequality (a *g-divergence* in the terminology of [11]), we obtain a new broad class of deficiency-like quantities, many of which can be related to various generalizations of the channel capacity and can be used to derive upper and lower bounds on the Le Cam and Shannon deficiencies.

II. NOTATION AND PRELIMINARIES

We take all alphabets to be *standard Borel* [12] (i.e., isomorphic to a Borel subspace of a Polish space). This covers virtually all settings of practical interest. Any such space X will always be endowed with its Borel σ -algebra $\mathcal{B}(X)$. The space of all probability measures on X will be denoted by $\mathcal{P}(X)$. The *total variation distance* between $P, Q \in \mathcal{P}(X)$ is

$$\|P - Q\| \triangleq \sup_{A \in \mathcal{B}(X)} |P(A) - Q(A)|. \quad (1)$$

A *Markov* (or *transition probability*) *kernel* between X and Y is a mapping $T : \mathcal{B}(Y) \times X \rightarrow [0, 1]$, such that $T(\cdot|x) \in \mathcal{P}(Y)$ for all $x \in X$ and $T(B|\cdot)$ is a measurable function on X for any $B \in \mathcal{B}(Y)$. We will denote the space of all such T by $\mathcal{M}(Y|X)$. If both X and Y are finite, then any $T \in \mathcal{M}(Y|X)$ is a stochastic matrix with elements $T(y|x)$, $(x, y) \in X \times Y$.

Any $T \in \mathcal{M}(Y|X)$ induces a mapping $\mathcal{P}(X) \rightarrow \mathcal{P}(Y)$, which we also will denote with a slight abuse of notation by T ; it maps any $P \in \mathcal{P}(X)$ to $Q = TP \in \mathcal{P}(Y)$, where

$$Q(B) = TP(B) \triangleq \int_X T(B|x)P(dx), \quad \forall B \in \mathcal{B}(Y)$$

We will denote the composition of Markov kernels by juxtaposition. That is, for $T \in \mathcal{M}(Y|X)$ and $S \in \mathcal{M}(U|Y)$, their

composition $ST \in \mathcal{M}(\mathsf{U}|\mathsf{X})$ is defined by

$$ST(C|x) \triangleq \int_{\mathsf{Y}} S(C|y)T(dy|x), \quad \forall x \in \mathsf{X}, C \in \mathcal{B}(\mathsf{U}).$$

The total variation distance (1) induces a distance on $\mathcal{M}(\mathsf{Y}|\mathsf{X})$:

$$\|T - T'\|_{\infty} \triangleq \sup_{x \in \mathsf{X}} \|T(\cdot|x) - T'(\cdot|x)\|.$$

The following *contractivity properties* are easily proved:

$$\|S(T - T')\|_{\infty} \leq \|T - T'\|_{\infty}, \|S - S'\|_{\infty} \leq \|S - S'\|_{\infty}. \quad (2)$$

A *channel* with input alphabet X and output alphabet Y is simply a Markov kernel $W \in \mathcal{M}(\mathsf{Y}|\mathsf{X})$ [13]. Whenever we need to specify the input and the output alphabets explicitly, we will represent the channel as a triple $(\mathsf{X}, \mathsf{Y}, W)$. An M -code for $(\mathsf{X}, \mathsf{Y}, W)$ is a pair (E, D) , where $E \in \mathcal{M}(\mathsf{X}|\{1, \dots, M\})$ is a (possibly randomized) encoder and $D \in \mathcal{M}(\{1, \dots, M\}|\mathsf{Y})$ is a (possibly randomized) decoder. The *average probability of error* of (E, D) on W is

$$\bar{\varepsilon}_W(E, D) \triangleq \frac{1}{M} \sum_{j=1}^M (1 - DWE(j|j)).$$

We say that (E, D) is an (M, ε) -code if $\bar{\varepsilon}_W(E, D) \leq \varepsilon$.

III. COMPARISON OF CHANNELS

A. Blackwell ordering and Le Cam deficiency

Consider two channels $(\mathsf{X}, \mathsf{Y}, W)$ and $(\mathsf{X}, \mathsf{U}, W')$ with a common input alphabet X . When are we justified in saying that W is less noisy (or more informative) than W' ? One natural answer is if we can simulate a single use of W' by post-processing a single use of W . This is precisely the comparison criterion proposed by Blackwell [1], [2]. Following the terminology of [7], we will say that W is *Blackwell sufficient* for W' and write $W \succeq_{\mathsf{B}} W'$ or $W' \preceq_{\mathsf{B}} W$ if $W' = TW$ for some Markov kernel $T \in \mathcal{M}(\mathsf{U}|\mathsf{Y})$. Another way to state this is that W' can be realized as a *stochastic degrading* of W . If $W \succeq_{\mathsf{B}} W'$ and $W' \succeq_{\mathsf{B}} W$, we will say that they are *Blackwell equivalent* and write $W \sim_{\mathsf{B}} W'$.

Blackwell sufficiency induces only a *partial order* on the family of all channels with the same input alphabet X . In fact, most channels are incomparable. As an example, let $\mathsf{X} = \mathsf{Y} = \mathsf{U} = \mathbb{R}$, and consider the additive channels

$$W : Y = X + Z \quad \text{and} \quad W' : U = X + Z'$$

where Z and Z' are real-valued random variables independent of the input X . Then $W \succeq_{\mathsf{B}} W'$ if and only if the law of Z is a convolution factor of the law of Z' [14].

The scope of Blackwell's theory was expanded considerably by Le Cam [5], who introduced the notion of a *deficiency* of one statistical experiment w.r.t. another. We can do the same for channels, where in our information-theoretic context the deficiency of W w.r.t. W' will quantify how well any post-processing of W can be used to *approximate* W' . The (Le

Cam) *deficiency* of W w.r.t. W' is given by

$$\delta(W, W') \triangleq \inf_{T \in \mathcal{M}(\mathsf{U}|\mathsf{Y})} \|TW - W'\|_{\infty}. \quad (3)$$

According to the so-called *Le Cam randomization criterion* [6], [7], the infimum in (3) is actually achieved, i.e., there exists some $T^* \in \mathcal{M}(\mathsf{U}|\mathsf{Y})$, such that $\delta(W, W') = \|T^*W - W'\|_{\infty}$. Note that $\delta(W, W') = 0$ iff $W \succeq_{\mathsf{B}} W'$.

B. Shannon ordering and Shannon deficiency

The Shannon ordering criterion [10] also revolves around the simulation of one channel by another, but it allows for pre- and post-processing, as well as for shared randomness between the input and the output terminals. Consider two channels $(\mathsf{X}, \mathsf{Y}, W)$ and $(\mathsf{X}', \mathsf{Y}', W')$, where the input alphabets need not necessarily be the same. Adopting our earlier terminology and the definitions of [10], we say that W is *Shannon sufficient* for W' , and write $W \succeq_{\mathsf{S}} W'$, if there exist some $k \in \mathbb{N}$, a probability vector $\pi = (\pi_1, \dots, \pi_k)$, and k pairs $(T_i, S_i) \in \mathcal{M}(\mathsf{X}|\mathsf{X}') \times \mathcal{M}(\mathsf{Y}'|\mathsf{Y})$, $1 \leq i \leq k$, such that

$$W' = \sum_{i=1}^k \pi_i S_i W T_i.$$

In the special case $\mathsf{X} = \mathsf{X}'$ it is easy to see that $W \succeq_{\mathsf{B}} W'$ implies $W \succeq_{\mathsf{S}} W'$, but not necessarily *vice versa*. Thus, the Shannon ordering criterion is weaker than that of Blackwell. However, it is also only a partial order. To be able to compare any channel to any other, we adopt Le Cam's strategy and define the corresponding *Shannon deficiency* as follows:

$$\delta_{\mathsf{S}}(W, W') \triangleq \inf_{k \in \mathbb{N}} \inf_{\pi; (T_i, S_i)_{i=1}^k} \left\| \sum_{i=1}^k \pi_i S_i W T_i - W' \right\|_{\infty}. \quad (4)$$

The following bound is immediate:

$$\delta_{\mathsf{S}}(W, W') \leq \inf_{T \in \mathcal{M}(\mathsf{X}|\mathsf{X}')} \delta(TW, W').$$

When $\mathsf{X} = \mathsf{X}'$, it can be weakened to $\delta_{\mathsf{S}}(W, W') \leq \delta(W, W')$, which suggests that the Le Cam deficiency is a more stringent measure than the Shannon deficiency. Again, note that $\delta_{\mathsf{S}}(W, W') = 0$ iff $W \succeq_{\mathsf{S}} W'$.

IV. BASIC PROPERTIES OF CHANNEL DEFICIENCIES

A. Operational interpretation

In [10], Shannon gave an operational interpretation to the \succeq_{S} channel ordering: If $W \succeq_{\mathsf{S}} W'$, then for any (M, ε) -code for W' we can find an (M, ε) -code for W . The following theorem is an extension of Shannon's result to an arbitrary pair of channels in terms of the Shannon deficiency (4):

Theorem 1. *Consider two channels $(\mathsf{X}, \mathsf{Y}, W)$ and $(\mathsf{X}', \mathsf{Y}', W')$. If $\delta_{\mathsf{S}}(W, W') \leq \varepsilon'$, then for any (M, ε) -code for W' we can find an $(M, \varepsilon + \varepsilon')$ -code for W .*

Proof: If $\delta_{\mathsf{S}}(W, W') \leq \varepsilon'$, we can find some $k \in \mathbb{N}$, a

probability vector $\pi = (\pi_1, \dots, \pi_k)$, and $(T_i, S_i)_{i=1}^k$, such that

$$\left\| \sum_{i=1}^k \pi_i S_i W T_i - W' \right\|_{\infty} \leq \varepsilon'. \quad (5)$$

Now let (E, D) be an (M, ε) -code for W' . For each $1 \leq i \leq k$, define an M -code (E_i, D_i) for W by letting $E_i = T_i E$ and $D_i = D S_i$. Let I be a random variable taking values in $\{1, \dots, k\}$ according to the distribution π . Then

$$\begin{aligned} & \mathbb{E}[\bar{\varepsilon}_W(E_I, D_I)] - \bar{\varepsilon}_{W'}(E, D) \\ & \leq \max_{j \in \{1, \dots, M\}} \sum_{\ell=1}^M \left| \mathbb{E}[D_I W E_I(\ell|j)] - D W' E(\ell|j) \right| \\ & = \|D(\mathbb{E}[S_I W T_I] - W')E\|_{\infty} \\ & \leq \|\mathbb{E}[S_I W T_I] - W'\|_{\infty} \\ & \leq \varepsilon', \end{aligned}$$

where the third step uses (2), and the last step is by (5). Since (E, D) is an (M, ε) -code for W' , this gives

$$\mathbb{E}[\bar{\varepsilon}_W(E_I, D_I)] \leq \varepsilon + \varepsilon'.$$

Therefore, there must exist at least one value $i^* \in \{1, \dots, k\}$, for which $\bar{\varepsilon}_W(E_{i^*}, D_{i^*}) \leq \varepsilon + \varepsilon'$. The theorem is proved. ■

B. Monotonicity of Le Cam deficiency under data processing

The following result is known (cf. [6]), but we give a self-contained proof, since we will need it later:

Theorem 2. Consider any three channels (X, Y, W) , (X, U, W') and (X, Z, V) . Then:

- 1) If $W \succeq_B W'$, then $\delta(W, V) \leq \delta(W', V)$.
- 2) If $W' \succeq_B V$, then $\delta(W, V) \leq \delta(W, W')$.

Equality holds in both cases if \succeq_B is replaced by \sim_B .

Proof: (1) Since $W \succeq_B W'$, there exists some $T \in \mathcal{M}(U|Y)$ such that $W' = TW$. Hence,

$$\begin{aligned} \delta(W', V) &= \inf_{S \in \mathcal{M}(Z|U)} \|S W' - V\|_{\infty} \\ &= \inf_{S \in \mathcal{M}(Z|U)} \|S T W - V\|_{\infty} \\ &\geq \inf_{S \in \mathcal{M}(Z|Y)} \|S W - V\|_{\infty} \\ &\equiv \delta(W, V). \end{aligned}$$

(2) Again, since $W' \succeq_B V$, there exists some $T \in \mathcal{M}(Z|U)$ such that $V = T W'$. Then

$$\begin{aligned} \delta(W, V) &= \inf_{S \in \mathcal{M}(Z|Y)} \|S W - T W'\|_{\infty} \\ &\leq \inf_{S \in \mathcal{M}(U|Y)} \|T S W - T W'\|_{\infty} \\ &\leq \inf_{S \in \mathcal{M}(U|Y)} \|S W - W'\|_{\infty} \\ &\equiv \delta(W, W'), \end{aligned}$$

where the third step is by (2).

C. Comparison to extremal channels

A channel with a *finite* input alphabet and an arbitrary output alphabet is called *semicontinuous* (SC) [15, Ch. 5]. For a fixed finite X , we have two *extremal* channels, W_0 and W_1 , such that $W_0 \preceq_B W \preceq_B W_1$ for any (X, Y, W) . These are unique up to Blackwell equivalence, and can be represented by $W_0(\cdot|x) = e_{x_0}$ for some fixed $x_0 \in X$ and by $W_1(\cdot|x) = e_x$, where e_x is the probability distribution that puts all mass on $x \in X$.

Given an arbitrary SC channel (X, Y, W) , let us define

$$\bar{\delta}(W) \triangleq \delta(W, W_1) \quad \text{and} \quad \underline{\delta}(W) \triangleq \delta(W_0, W).$$

Theorem 3 (Torgersen [16]).

$$\bar{\delta}(W) = \inf_{T \in \mathcal{M}(X|Y)} \max_{x \in X} (1 - T W(x|x)) \quad (6)$$

$$\underline{\delta}(W) = \inf_{Q \in \mathcal{P}(Y)} \max_{x \in X} \|Q - W(\cdot|x)\| \quad (7)$$

In operational terms, $\bar{\delta}(W)$ is the best maximal probability of error achievable on W with an arbitrary decoder and an identity encoder, $E(x|x) = 1$ for all $x \in X$. On the other hand, $\underline{\delta}(W)$ has a geometric interpretation as the “radius” of the set $\{W(\cdot|x)\}_{x \in X}$ w.r.t. the total variation distance. The unique $Q^* \in \mathcal{P}(Y)$ that attains the infimum in (7) can be thought as the “center” of this set. We will come back to this interpretation later and extend it to more general notions of “information radius” [9], [17].

We can also use the Shannon deficiency to define $\bar{\delta}_S(W)$ and $\underline{\delta}_S(W)$. The former can be bounded as $\bar{\delta}_S(W) \leq \bar{\delta}(W)$. For the latter, it is easy to show from definitions that it is equal to the Le Cam deficiency $\underline{\delta}(W)$.

V. GENERALIZED DEFICIENCIES

By using other statistical or information-theoretic measures of (dis)similarity between probability distributions instead of the total variation distance, it is possible to construct a wide variety of deficiency-like quantities. The resulting alternative (or generalized) deficiencies are of interest in their own right, and they can also be used to bound the Le Cam deficiencies from above and from below. We start by defining a broad class of deficiencies based on a very minimalistic requirement: monotonicity under data processing. Following recent work of Polyanskiy and Verdú [11], we call any mapping $\mathcal{D} : \mathcal{P}(X) \times \mathcal{P}(X) \rightarrow \mathbb{R}$ a *g-divergence* if

$$\mathcal{D}(TP \| TQ) \leq \mathcal{D}(P \| Q) \quad (8)$$

for any $P, Q \in \mathcal{P}(X)$ and any Markov kernel $T \in \mathcal{M}(Y|X)$. Given the channels (X, Y, W) and (X, U, W') , let us define the corresponding *g-deficiency* of W w.r.t. W' as

$$\delta_{\mathcal{D}}(W, W') \triangleq \inf_{T \in \mathcal{M}(U|Y)} \sup_{x \in X} \mathcal{D}(W'(\cdot|x) \| T W(\cdot|x)).$$

The data-processing inequality (8) is enough to ensure the same monotonicity properties as the Le Cam deficiency:

Theorem 4. The analog of Theorem 2 holds if we replace δ with $\delta_{\mathcal{D}}$ throughout. ■

As we shall see next, particular choices of a g -divergence with additional properties (such as convexity) lead to many useful alternative notions of deficiency.

A. f -deficiencies

A wide class of g -divergences is formed by the so-called f -divergences of Csiszár (cf. [18] and references therein). Let $f : (0, \infty) \rightarrow \mathbb{R}$ be a convex function with $f(1) = 0$. Let P and Q be two probability measures on a space X . Then the f -divergence between P and Q is defined as

$$D_f(P\|Q) \triangleq \int_X \frac{dQ}{d\mu} f\left(\frac{dP/d\mu}{dQ/d\mu}\right) d\mu, \quad (9)$$

where μ is any σ -finite measure on X that dominates both P and Q ¹, and we use the conventions

$$0 \cdot f(0/0) = 0, \quad 0 \cdot f(t/0) = t \lim_{s \searrow 0} s f(1/s), \forall t > 0. \quad (10)$$

With proper choice of f , we recover many of the frequently used measures of divergence, for example:

- $f(t) = t \log t$: the relative entropy $D(P\|Q)$.
- $f(t) = |t - 1|$: the total variation distance $\|P - Q\|$.

Thus, given two channels (X, Y, W) and (X, U, W') , we can define the f -deficiency of W w.r.t. W' as

$$\delta_f(W, W') \triangleq \inf_{T \in \mathcal{M}(U|Y)} \sup_{x \in X} D_f(W'(\cdot|x)\|TW(\cdot|x))$$

For an SC channel (X, Y, W) , we have an analog of Theorem 3:

Theorem 5. Suppose that $\lim_{t \searrow 0} f(t) = 0$. Then

$$\bar{\delta}_f(W) = \inf_{T \in \mathcal{M}(X|Y)} \max_{x \in X} f^*(TW(x|x)) \quad (11)$$

$$\underline{\delta}_f(W) = \inf_{Q \in \mathcal{P}(Y)} \max_{x \in X} D_f(W(\cdot|x)\|Q), \quad (12)$$

where $f^*(t) \triangleq t f(1/t)$ is the Csiszár conjugate of f [18].

Remark 1. The f -deficiency $\underline{\delta}_f(W)$ is the *absolute f -informativity* (or the *f -radius*) of W [9].

Proof: In (9), let μ be the counting measure on X . Then

$$\begin{aligned} \bar{\delta}_f(W) &= \inf_T \max_{x \in X} \sum_{x' \in X} TW(x'|x) \cdot f\left(\frac{e_x(x')}{TW(x'|x)}\right) \\ &= \inf_T \max_{x \in X} TW(x|x) \cdot f\left(\frac{1}{TW(x|x)}\right) \end{aligned} \quad (13)$$

$$= \inf_T \max_{x \in X} f^*(TW(x|x)), \quad (14)$$

where (13) uses (10) and the fact that $f(t) \rightarrow 0$ as $t \searrow 0$, while (14) uses the definition of f^* . This proves (11); the simple proof of (12) is omitted. ■

For the special case $f(t) = t \log t$ (natural logarithms), let us use the term *I-deficiency* and write δ_I .

¹It is not hard to show that the value of $D_f(P\|Q)$ does not depend on the choice of the dominating measure μ .

Theorem 6. For any two channels (X, Y, W) and (X, U, W') ,

$$\delta(W, W') \leq \sqrt{\delta_I(W, W')/2}. \quad (15)$$

For any SC channel (X, Y, W) ,

$$\bar{\delta}_I(W) = -\log(1 - \bar{\delta}(W)) \quad \text{and} \quad \underline{\delta}_I(W) = C(W),$$

where $C(W) = \max_{P \in \mathcal{P}(X)} I(P, W)$ is the Shannon capacity of W . Moreover, if (X, Y, W) and (X, U, W') are SC, then $\delta_I(W, W') \leq \log |X|$.

Proof: The bound (15) follows from Pinsker's inequality. The expression for $\bar{\delta}_I(W)$ follows from Theorems 5 and 3. Next, using (12) and the minimax theorem, we have

$$\begin{aligned} \underline{\delta}_I(W) &= \max_{P \in \mathcal{P}(X)} \min_{Q \in \mathcal{P}(Y)} \sum_{x \in X} P(x) D(W(\cdot|x)\|Q) \\ &= \max_{P \in \mathcal{P}(X)} I(P, W) \equiv C(W), \end{aligned}$$

Finally, by data processing we have $\delta_I(W, W') \leq \delta_I(W_0, W_1) = C(W_1) = \log |X|$. ■

B. Rényi deficiencies

The *Rényi divergence* of order $\lambda \in (0, \infty) \setminus \{1\}$ between P and Q is [19]

$$D_\lambda(P\|Q) \triangleq \frac{1}{\lambda - 1} \log \int \frac{dQ}{d\mu} \left(\frac{dP/d\mu}{dQ/d\mu}\right)^\lambda d\mu.$$

Although it is not an f -divergence (but rather a monotone transformation of one, the so-called Hellinger divergence [18]), it is a g -divergence [11]. Hence we can define the *Rényi deficiency* of order λ , which we will denote by $\delta_\lambda(\cdot, \cdot)$.

Theorem 7. Let (X, Y, W) and (X, U, W') be SC. Then

$$\begin{aligned} \underline{\delta}_\lambda(W) &= C_\lambda(W) \\ &= \max_{P \in \mathcal{P}(X)} \min_{Q \in \mathcal{P}(Y)} \sum_{x \in X} P(x) D_\lambda(W(\cdot|x)\|Q), \end{aligned}$$

the *order- λ capacity* of W [19], and $\delta_\lambda(W, W') \leq \log |X|$.

Proof: An argument similar to the proof of (12) leads to

$$\underline{\delta}_\lambda(W) = \min_{Q \in \mathcal{P}(Y)} \max_{x \in X} D_\lambda(W(\cdot|x)\|Q).$$

From Proposition 1 in [19], the right-hand side is equal to $C_\lambda(W)$. By data processing, $\delta_\lambda(W, W') \leq \delta_\lambda(W_0, W_1) = C_\lambda(W_1)$. From [19, p. 28], $C_\lambda(W_1) = \max_{P \in \mathcal{P}(X)} H(P) = \log |X|$, where $H(P)$ is the Shannon entropy of P . ■

If both X and Y are finite and $\lambda = 1/(1 + \rho)$ for some $\rho > 0$, $\rho \neq 1$, then we can express $\underline{\delta}_\lambda(W)$ in terms of Gallager's function $E_0(\rho, P, W)$ [20] as

$$\underline{\delta}_\lambda(W) = \frac{1}{\rho} \max_{P \in \mathcal{P}(X)} E_0(\rho, P, W).$$

C. Neyman–Pearson deficiencies

An interesting result from [11] is that there exist g -divergences that are not representable as monotone functions of f -divergences. One example given there uses the *Neyman–Pearson functions* [6]. Given two probability measures P and

Q on a space X , let us define for any $0 < \alpha < 1$ the number

$$\beta_\alpha(P, Q) \triangleq \inf_{T \in \mathcal{M}(\{0,1\}|X)} \{TQ(1) : TP(1) \geq \alpha\}.$$

Then $-\beta_\alpha(TP, TQ) \leq -\beta_\alpha(P, Q)$ for any $T \in \mathcal{M}(Y|X)$ [11]. Using this fact, let us define the *Neyman–Pearson α -divergence* between P and Q by

$$\mathcal{D}_\alpha(P||Q) \triangleq \alpha - \beta_\alpha(P, Q);$$

we will denote the corresponding g -deficiency by $\delta_\alpha(\cdot, \cdot)$.

Theorem 8. *The Le Cam deficiency $\delta(W, W')$ between any two channels (X, Y, W) and (X, U, W') satisfies the bound*

$$\delta(W, W') \geq \delta_\alpha(W, W')$$

for any $0 < \alpha < 1$. Moreover, if (X, Y, W) and (X, U, W') are SC, then $\delta_\alpha(W, W') \leq \alpha(1 - 1/|X|)$.

Proof: We start by noting that $\|P - Q\| = \sup_{0 < \alpha < 1} \mathcal{D}_\alpha(P||Q)$ [6, p. 36]. Then

$$\begin{aligned} \delta(W, W') &= \inf_T \sup_{x \in X} \|TW(\cdot|x) - W'(\cdot|x)\| \\ &= \inf_T \sup_{0 < \alpha < 1} \sup_{x \in X} \mathcal{D}_\alpha(W'(\cdot|x)||TW(\cdot|x)) \\ &\geq \sup_{0 < \alpha < 1} \mathcal{D}_\alpha(W, W'). \end{aligned}$$

If (X, Y, W) and (X, U, W') are SC, then by data processing $\delta_\alpha(W, W') \leq \delta_\alpha(W_0, W_1)$, which in turn is equal to

$$\alpha - \sup_{Q \in \mathcal{P}(X)} \min_{x \in X} \beta_\alpha(e_x, Q) = \alpha(1 - 1/|X|)$$

(see Lemma 1 in the Appendix). ■

VI. CONCLUSION

We have extended the notion of comparison between statistical experiments [1], [2], as well as an approximation theory for experiments [5], to coding/decoding scenarios common in information theory. Starting from a comparison criterion for channels due to Shannon [10], we have defined an approximation measure for channels, which we have termed the *Shannon deficiency*, mimicking the terminology of [5]. An interesting, though challenging, future direction would be to develop an approximation-based perspective on the coding theorems of information theory, similar to Le Cam's programme in asymptotic statistics.

REFERENCES

- [1] D. Blackwell, "Comparison of experiments," in *Proc. 2nd Berkeley Symp. Math. Statist. Prob.*, 1951, pp. 93–102.
- [2] —, "Equivalent comparisons of experiments," *Ann. Math. Statist.*, vol. 24, pp. 265–272, 1953.
- [3] D. V. Lindley, "On a measure of the information provided by an experiment," *Ann. Math. Statist.*, vol. 27, pp. 986–1005, 1956.
- [4] P. K. Goel and M. H. DeGroot, "Comparison of experiments and information measures," *Ann. Statist.*, vol. 7, no. 5, pp. 1066–1077, 1979.
- [5] L. Le Cam, "Sufficiency and approximate sufficiency," *Ann. Math. Statist.*, vol. 35, pp. 1419–1455, 1964.
- [6] E. Torgersen, *Comparison of Statistical Experiments*. Cambridge Univ. Press, 1991.
- [7] A. N. Shiryaev and V. G. Spokoiny, *Statistical Experiments and Decisions: Asymptotic Theory*. Singapore: World Scientific, 2000.

- [8] L. D. Brown and M. G. Low, "Asymptotic equivalence of nonparametric regression and white noise," *Ann. Statist.*, vol. 24, no. 6, pp. 2384–2398, 1996.
- [9] I. Csiszár, "A class of measures of informativity of observation channels," *Periodica Math. Hungar.*, vol. 2, no. 1–4, pp. 191–213, 1972.
- [10] C. E. Shannon, "A note on a partial ordering for communication channels," *Inform. Control*, vol. 1, pp. 390–397, 1958.
- [11] Y. Polyanskiy and S. Verdú, "Arimoto channel coding converse and Rényi divergence," in *Proc. 48th Annu. Allerton Conf. on Commun., Control, and Comput.*, 2010, pp. 1327–1333.
- [12] K. R. Parthasarathy, *Probability Measures on Metric Spaces*. New York: Academic Press, 1967.
- [13] R. L. Dobrushin, "A general formulation of the basic Shannon theorem in information theory," *Uspekhi Math. Nauk*, vol. 14, no. 6, pp. 3–103, 1959.
- [14] M. Stone, "Non-equivalent comparisons of experiments and their use for experiments involving location parameters," *Ann. Statist.*, vol. 32, pp. 326–332, 1961.
- [15] A. Feinstein, *Information Theory*. New York: McGraw-Hill, 1958.
- [16] E. N. Torgersen, "Measures of information based on comparison with total information and with total ignorance," *Ann. Statist.*, vol. 9, no. 3, pp. 638–657, 1981.
- [17] R. Sibson, "Information radius," *Z. Wahrsch. verw. Geb.*, vol. 14, pp. 149–160, 1969.
- [18] 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.
- [19] I. Csiszár, "Generalized cutoff rates and Rényi's information measures," *IEEE Trans. Inform. Theory*, vol. 41, no. 1, pp. 26–34, January 1995.
- [20] R. G. Gallager, *Information Theory and Reliable Communication*. New York: Wiley, 1968.

APPENDIX

Lemma 1. *For a finite alphabet X ,*

$$\sup_{Q \in \mathcal{P}(X)} \min_{x \in X} \beta_\alpha(e_x, Q) = \alpha/|X|.$$

Proof: Fix some $Q \in \mathcal{P}(X)$ and $x \in X$. Then

$$\beta_\alpha(e_x, Q) = \min_{T \in \mathcal{M}(\{0,1\}|X)} \{TQ(1) : Te_x(1) \geq \alpha\}. \quad (16)$$

Since $Te_x(1) = \sum_{x' \in X} T(1|x')e_x(x') = T(1|x)$, we can rewrite (16) as

$$\beta_\alpha(e_x, Q) = \min_{T \in \mathcal{M}(\{0,1\}|X)} \{TQ(1) : T(1|x) \geq \alpha\}.$$

For any feasible T , we must have

$$TQ(1) = \sum_{x' \in X} T(1|x')Q(x') \geq T(1|x)Q(x) \geq \alpha Q(x),$$

which means that $\beta_\alpha(e_x, Q) \geq \alpha Q(x)$. On the other hand, choosing T^* with $T^*(1|x) = \alpha$ and $T^*(1|x') = 0$ for all $x' \neq x$, we have $T^*Q(1) = \alpha Q(x)$, so $\beta_\alpha(e_x, Q) \leq \alpha Q(x)$. Thus, $\beta_\alpha(e_x, Q) = \alpha Q(x)$. This and the fact that

$$\max_{Q \in \mathcal{P}(X)} \min_{x \in X} Q(x) = \min_{x \in X} \max_{Q \in \mathcal{P}(X)} Q(x) = 1/|X|$$

completes the proof. ■