

# Time-Sharing Vs. Source-Splitting in the Slepian-Wolf Problem: Error Exponents Analysis

Todd P. Coleman, Muriel Médard; Michelle Effros

{colemant,medard}@mit.edu; effros@caltech.edu

Massachusetts Institute of Technology; California Institute of Technology  
Laboratory for Information and Decision Systems; Data Compression Laboratory  
77 Massachusetts Ave, 32-D626; 1200 East California Boulevard, MS 136-93  
Cambridge, MA 02139; Pasadena, CA 91125

## Abstract

*We discuss two approaches for decoding at arbitrary rates in the Slepian-Wolf problem - time sharing and source splitting - both of which rely on constituent vertex decoders. We consider the error exponents for both schemes and conclude that source-splitting is more robust at coding at arbitrary rates, as the error exponent for time-sharing degrades significantly at rates near vertices. As a by-product of our analysis, we exhibit an interesting connection between minimum mean-squared error estimation and error exponents.*

## 1 Introduction

In this setting we discuss alternative approaches to attain achievable rates on the dominant face of the Slepian-Wolf region. Here we will focus on the two-user setting, but this can naturally be generalized. Consider two sources  $(U^1, U^2)$  with joint probability distribution  $W(u^1, u^2)$ . The two-user Slepian-Wolf region  $\mathcal{R}[W]$  is given by

$$\mathcal{R}[W] = \left\{ \underline{R} \in \mathbb{R}_+^2 : \sum_{i \in S} R_i \geq H(U(S)|U(S^c)) \quad \forall S \subseteq \{1, 2\} \right\}$$

where  $U(S) = \{U^j, j \in S\}$ . Generally speaking, the decoder must find a pair of ‘jointly typical’ sequences [1, pp. 194-197] consistent with what is observed. This is in general a computationally difficult task. At *vertex* rate points, the joint search over both codebooks for a pair of ‘jointly typical’ sequences can be done successively. For instance, if users would like to communicate at the rate of  $(R_1, R_2) = (H(U^1), H(U^2|U^1))$ , then we note that communicating at a rate of  $H(U^1)$  can be done by simply entropy-encoding either a variable-rate lossless fashion or a near-lossless fixed-rate fashion. After successful decoding,  $U^1$  can be passed as side information to help decode  $U^2$  at a rate of  $H(U^2|U^1)$ . By exchanging the roles of  $U^1$  and  $U^2$ , it follows that the same approach applies to encoding at the vertex rate  $(R_1, R_2) = (H(U^1|U^2), H(U^2))$ . Recently, much attention has been paid to the construction of low-complexity decoders to achieve rates of  $R_2$  very close to  $H(U^2|U^1)$ .

A more interesting question concerns communicating at *any* rate in the achievable rate region - not necessarily vertices. The most efficient communication schemes minimize sum rate and thus attain rates lying on the *dominant face*,  $\mathcal{D}\{\mathcal{R}[W]\}$ , given by  $\{(R_1, R_2) \in \mathcal{R}[W] : R_1 + R_2 = H(U^1, U^2)\}$ . Two candidate approaches of using decoding strategies that rely upon vertex decoding are:

- time-sharing, where coding for a non-vertex point is done by coding a certain fraction  $\alpha \in [0, 1]$  of the time at one vertex, and the remaining fraction  $1 - \alpha$  of the time at the other vertex
- source-splitting [2, 3], where coding for a non-vertex point in a two-source problem is done by splitting one of the sources and coding at a vertex rate in the corresponding three-source problem

We would like to understand here the performance of the two candidate approaches at rates near the joint entropy boundary, in terms of error probability. We illustrate below that the source-splitting approach is more robust for decoding at arbitrary rates on the dominant face as compared to time-sharing, which can have significant error exponent penalty at rates close to vertices. As a by-product of our analysis, we show an interesting connection between information theory and estimation theory: the error exponent of vertex decoding in an arbitrary instance of the Slepian-Wolf problem depends on the inverse of the Fisher information of Gallager's  $\rho$ -parametrized tilted distribution.

## 2 Error Exponents

Here we discuss the near-lossless fixed-to-fixed distributed data compression setting where  $n$  samples of the memoryless source  $\{(U_i^1, U_i^2)\}_{i=1}^n$  are separately encoded. For each source  $j$ , the  $\{U_i^j\}_{i=1}^n$  symbols will be mapped to  $2^{nR_j}$  output symbols. The error exponent for a particular coding scheme  $k$  will be denoted by

$$E^k(R_1, R_2) \triangleq \liminf_{n \rightarrow \infty} -\frac{1}{n} \log P_e^k(R_1, R_2).$$

As illustrated in the appendix (13), for random variables  $X, Y$  with joint distribution  $W$ , the error exponent  $E_{x|y}(R)$  for source coding  $X$  at rate  $R$  with side information  $Y$  has a flat slope at  $R = H(X|Y)$ :

$$\frac{d}{dR} \{E_{x|y}(R)\}_{R=H(X|Y)} = 0.$$

Thus to capture the behavior of the exponent at  $R = H(X|Y) + \delta$ , we must consider second order effects via a Taylor series expansion:

$$\begin{aligned} E_{x|y}(H(X|Y) + \delta) &= \frac{1}{2} \delta^2 \frac{d^2}{dR^2} \{E_{x|y}(R)\}_{R=H(X|Y)} + o(\delta^3) \\ &= \frac{1}{2} \delta^2 E''_{x|y}(H(X|Y)) + o(\delta^3). \end{aligned}$$

We will denote the error exponent for time-sharing as  $E^t(R_1, R_2)$  and that for source-splitting as  $E^s(R_1, R_2)$ . We are interested in the behavior of the error exponent at rates near the dominant face.

### 2.1 Time-Sharing

Time-sharing is one approach to attain any rate on the dominant face. For  $\alpha \in [0, 1]$ ,  $\alpha n$  of the samples are encoded near the vertex  $(R_1, R_2) = (H(U^1), H(U^2|U^1))$  and the remaining  $(1 - \alpha)n$  samples are encoded near the other vertex  $(R_1, R_2) = (H(U^1|U^2), H(U^2))$ . We will assume that decoding is done with the pipelined vertex decoding approach described above. Thus for the

decoding of the  $\alpha n$  symbol pairs at the rate  $(H(U^1) + \delta, H(U^2|U^1) + \delta)$ , we have

$$\begin{aligned} P_e^{t,\alpha} &\leq P\left(\left[\hat{U}^1\right]_1^{\alpha n} \neq \left[U^1\right]_1^{\alpha n}\right) + P\left(\left[\hat{U}^2\right]_1^{\alpha n} \neq \left[U^2\right]_1^{\alpha n} \mid \left[U^1\right]_1^{\alpha n}\right) \\ &= 2^{-n\alpha[E_{u^1}(H(U^1)+\delta)-o(n)]} + 2^{-n\alpha[E_{u^2|u^1}(H(U^2|U^1)+\delta)-o(n)]} \end{aligned}$$

For the decoding of the  $(1-\alpha)n$  symbol pairs at the rate  $(H(U^1|U^2) + \delta, H(U^2) + \delta)$ , we have

$$\begin{aligned} P_e^{t,1-\alpha} &\leq P\left(\left[\hat{U}^2\right]_{\alpha n+1}^n \neq \left[U^2\right]_{\alpha n+1}^n\right) + P\left(\left[\hat{U}^1\right]_{\alpha n+1}^n \neq \left[U^1\right]_{\alpha n+1}^n \mid \left[U^2\right]_{\alpha n+1}^n\right) \\ &= 2^{-n(1-\alpha)[E_{u^2}(H(U^2)+\delta)-o(n)]} + 2^{-n(1-\alpha)[E_{u^1|u^2}(H(U^1|U^2)+\delta)-o(n)]} \end{aligned}$$

Thus it follows that for  $(R_1, R_2) \in \mathcal{D}$ ,

$$\begin{aligned} E^t(R_1 + \delta, R_2 + \delta) &= \min \left[ \alpha E_{u^1}(H(U^1) + \delta), \alpha E_{u^2|u^1}(H(U^2|U^1) + \delta), \right. \\ &\quad \left. (1-\alpha)E_{u^2}(H(U^2) + \delta), (1-\alpha)E_{u^1|u^2}(H(U^1|U^2) + \delta) \right] \\ &= \frac{1}{2}\delta^2 \min \left[ \alpha E''_{u^1}(H(U^1)), \alpha E''_{u^2|u^1}(H(U^2|U^1)), \right. \\ &\quad \left. (1-\alpha)E''_{u^2}(H(U^2)), (1-\alpha)E''_{u^1|u^2}(H(U^1|U^2)) \right] + o(\delta^3) \end{aligned}$$

where  $\alpha$  satisfies

$$R_1 = \alpha H(U^1) + (1-\alpha)H(U^1|U^2). \quad (1)$$

## 2.2 Source-Splitting

*Source-splitting* transforms a point on the dominant face of the two-source problem to a vertex point in a three-source problem. This is done by

$$U_i^1 \mapsto \left( \begin{array}{l} U_i^{1a} = f_a(U_i^1) \\ U_i^{1b} = f_b(U_i^1) \end{array} \right) \mapsto U_i = f(U_i^{1a}, U_i^{1b}) \quad (2)$$

where the functions  $f_a: \mathcal{U}_1 \rightarrow \mathcal{U}_1$ ,  $f_b: \mathcal{U}_1 \rightarrow \mathcal{U}_1$  and  $f: \mathcal{U}_1 \rightarrow \mathcal{U}_1$  satisfy

$$f(f_a(u), f_b(u)) = u \quad \forall u \in \mathcal{U}_1.$$

As an example [3], this can be done as follows:

$$U_i \mapsto \left( \begin{array}{l} U_i^a = \min(\pi(U_i), T) \\ U_i^b = \max(\pi(U_i), T) - T \end{array} \right) \mapsto U_i = \pi^{-1}(U_i^a + U_i^b), \quad (3)$$

where  $T \in \mathcal{U}$  operates as a threshold and  $\pi \in \Pi(\mathcal{U})$  is a permutation operator.

If we have two discrete memoryless sources  $(U^1, U^2)$  drawn according to  $P(u^1, u^2)$ , then we can split  $U^1$  to form  $(U^{1a}, U^{1b})$  as shown in (3). At this point, we have three sources, each of which can be encoded separately at rates  $R_{1a}, R_{1b}, R_2$ . We note that because  $U \leftrightarrow (U^{1a}, U^{1b})$ ,  $H(U^1, U^2) = H(U^{1a}, U^{1b}, U^2)$ . Through the chain rule for entropy, we consider the rates

$$R_{1a} = H(U^{1a}) \quad (4a)$$

$$R_2 = H(U^2|U^{1a}) \quad (4b)$$

$$R_{1b} = H(U^{1b}|U^2, U^{1a}) \quad (4c)$$

$$R_1 = R_{1a} + R_{1b}. \quad (4d)$$

For any nontrivial split,  $(R_1, R_2)$  is not a vertex in  $\mathcal{R}[P(u^1, u^2)]$ , but  $(R_{1a}, R_2, R_{1b})$  is a vertex in  $\mathcal{R}[P(u^{1a}, u^2, u^{1b})]$ . This directly implies a parallelizable encoding strategy and pipelined single-user decoding strategy that operates with the complexity of a smaller-alphabet decoder.

This corresponds to a vertex in the  $U^{1a}, U^{1b}, U^2$  problem by using the chain rule for entropy and encoding at rates given by (4). Also in this scheme, coding to attain a non-vertex is mapped to coding at a vertex. Here we also assume that decoding is done with the pipelined vertex decoding approach. Thus for the decoding, we have

$$\begin{aligned} P_e^s &\leq P(\hat{U}^{1a} \neq \underline{U}^{1a}) + P(\hat{U}^2 \neq \underline{U}^2|\underline{U}^{1a}) + P(\hat{U}^{1b} \neq \underline{U}^{1b}|\underline{U}^{1a}, \underline{U}^2) \\ &= 2^{-n[E_{u^{1a}}(H(U^{1a}) + \frac{1}{2}\delta) - o(n)]} + 2^{-n[E_{u^2|u^{1a}}(H(U^2|U^{1a}) + \delta) - o(n)]} \\ &\quad + 2^{-n[E_{u^{1b}|u^2, u^{1a}}(H(U^{1b}|U^2, U^{1a}) + \frac{1}{2}\delta) - o(n)]} \end{aligned}$$

Thus it follows that for  $(R_1, R_2) \in \mathcal{D}$ ,

$$\begin{aligned} E^s(R_1 + \delta, R_2 + \delta) &= \min \left[ E_{u^{1a}} \left( H(U^{1a}) + \frac{1}{2}\delta \right), E_{u^2|u^{1a}} \left( H(U^2|U^{1a}) + \delta \right), \right. \\ &\quad \left. E_{u^{1b}|u^2, u^{1a}} \left( H(U^{1b}|U^2, U^{1a}) + \frac{1}{2}\delta \right) \right] \\ &= \frac{1}{2}\delta^2 \min \left[ \frac{1}{4}E''_{u^{1a}} \left( H(U^{1a}) \right), E''_{u^2|u^{1a}} \left( H(U^2|U^{1a}) \right), \right. \\ &\quad \left. \frac{1}{4}E''_{u^{1b}|u^2, u^{1a}} \left( H(U^{1b}|U^2, U^{1a}) \right) \right] + o(\delta^3). \end{aligned}$$

Note that to attain a rate of  $(R_1 + \delta, R_2 + \delta)$  we have to allocate  $\frac{1}{2}\delta$  extra rate to  $U^{1a}$  and  $\frac{1}{2}\delta$  to  $U^{1b}$ , as compared to the usual  $\delta$  to  $U^1$  in the time-sharing case.

### 2.3 Comparison

It is the purpose of this discussion to observe how the error exponents behave for the two approaches when coding at rates  $(R_1 + \delta, R_2 + \delta)$  where  $(R_1, R_2) \in \mathcal{D}$ . From the onset it is not clear which approach has better exponents - both cases exhibit error exponent degradation. In the case of time-sharing, error exponent degradation is caused by a reduction in the effective block length by

factors of  $\alpha$  and  $1 - \alpha$ . In the source-splitting scenario, error exponent degradation arises because of error propagation in decoding three sources rather than two, along with the reduction by a factor of  $\frac{1}{4}$  due to the splitting operation. Furthermore, the comparison is not straightforward because the source-splitting operation creates a new joint distribution on the three sources, as compared to the original joint distribution on the two sources. Although these distributions are in some sense equivalent because  $(U_i^1, U_i^2)$  and  $(U_i^{1a}, U_i^2, U_i^{1b})$  form a bijection, the behavior of the error exponent's second derivative involves more complicated functions of the distribution than just entropy:

$$\textbf{Lemma 2.1} \quad E''_{x|y}(H(X|Y)) = \frac{1}{-H(X|Y)^2 + \sum_{x,y} Q(x,y) \log^2[Q(x|y)]}.$$

Proof details are in the appendix. What is interesting about this denominator is that it can be characterized in terms of the Fisher information of Gallager's  $\rho$ -tilted distribution [4]. In particular, if we define  $Q_\rho(x, y)$  as the product of  $Q_\rho(y)$ , given by (7), and  $Q_\rho(x|y)$ , given by (8), then we can calculate the Fisher information of this parametrized probability distribution:

$$F(\rho) = \sum_{x,y} \left( \frac{d \log Q_\rho(x,y)}{d\rho} \right)^2 Q_\rho(x,y) \quad (5)$$

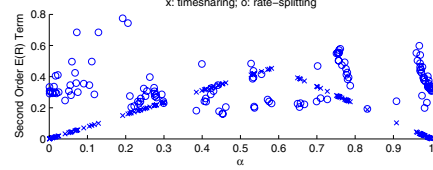
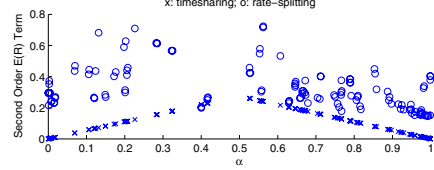
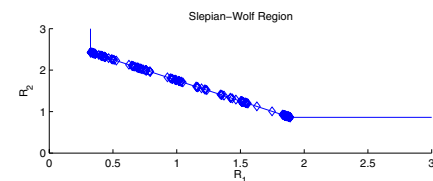
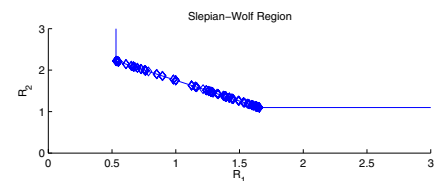
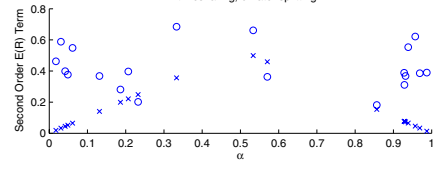
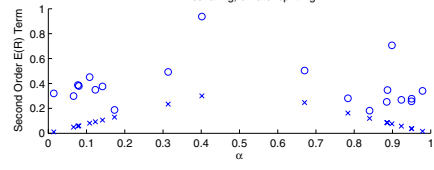
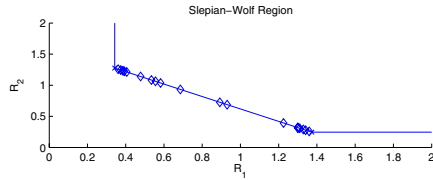
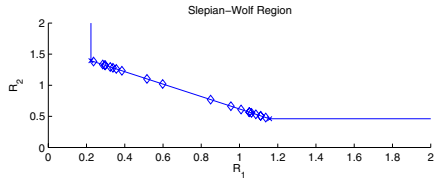
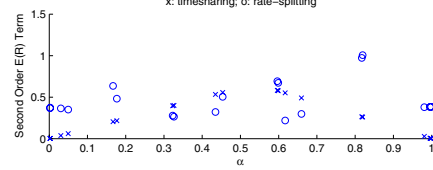
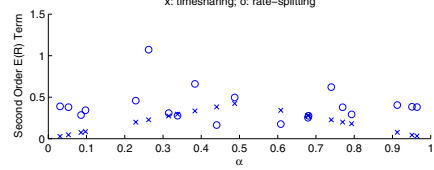
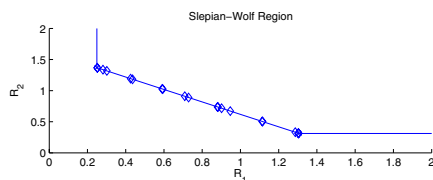
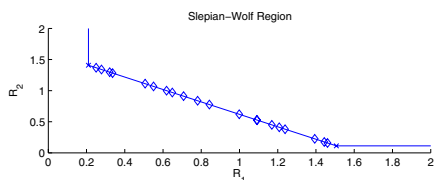
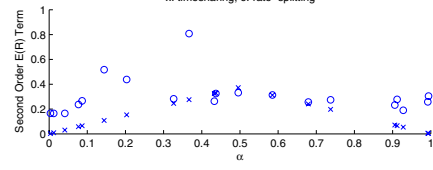
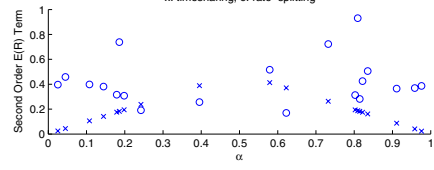
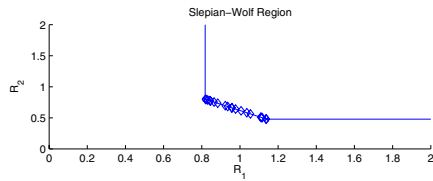
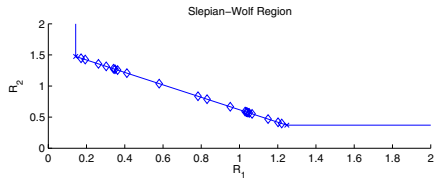
Then we can characterize the error exponent's second derivative in terms of the inverse of the above Fisher information quantity:

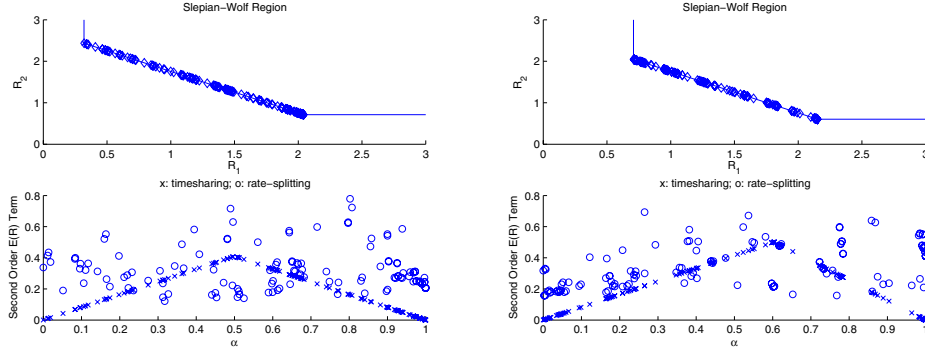
$$\textbf{Lemma 2.2} \quad E''_{x|y}(H(X|Y)) = \frac{1}{F(\rho)} \Big|_{\rho=0}.$$

Proof details are in the appendix. This leads to another interesting connection between information-theoretic quantities (error exponents) and estimation theoretic ones (MMSE and Fisher Information). In particular, the connection relates the error exponent's second derivative to the inverse of a Fisher information - which is a bound on minimum mean-squared error. The common thread appears to lie in the information geometry [5] interpretation of the Kullback-Leibler distance. However in our opinion, an in-depth understanding of this relation remains to be found.

### 3 Examples on Randomly Constructed Joint Probability Distributions

Here we randomly construct joint probability distributions  $W$  on  $(U^1, U^2)$  and compare  $E^s(R_1 + \delta, R_2 + \delta)$  with  $E^t(R_1 + \delta, R_2 + \delta)$ . In the figure pairs below, the top figure in each pair shows the Slepian-Wolf achievable rate region and the target rate points on the dominant face. The bottom figure in each pair shows  $E^s(R_1 + \delta, R_2 + \delta)$  and  $E^t(R_1 + \delta, R_2 + \delta)$  as a function of  $\alpha$ , where  $\alpha$  satisfies (1). For the splitting case, splitting is done according to (3). The takeaway theme from all these examples is that the minimum  $E^s(R_1 + \delta, R_2 + \delta)$  for points  $(R_1, R_2) \in \mathcal{D}$  is bounded away from 0 whereas  $E^t(R_1 + \delta, R_2 + \delta)$  decays to 0 linearly as  $\alpha$  approaches 0 or 1. Consequently at rates close to vertices, the second order source-splitting exponent significantly dominates that of time-sharing. At rates halfway between vertices, in some cases source-splitting wins, and in other cases time-sharing does. We were not able to find many cases where the second-order time-sharing exponent significantly dominates that of source-splitting. Thus in terms of error exponents, source-splitting appears to be more robust across various rates than time-sharing.





## Acknowledgment

The authors would like to thank Ralf Koetter for his initial suggestion to consider this comparison.

## A Proof of Lemma 2.1

Here we consider the ML decoding error exponent for source decoding  $x$  when side information  $y$  is known at the decoder. Denote  $P_e(y)$  to be the error probability conditioned upon receiving  $y$ . Then from [4] we have that

$$\begin{aligned} \frac{-1}{n} \log P_e(y) &\geq E_{x|y}(R, y) \triangleq \max_{0 \leq \rho \leq 1} \rho R - E_{0,x|y}(\rho, y) \\ E_{0,x|y}(\rho, y) &\triangleq (1 + \rho) \log \left[ \sum_x Q(x|y)^{\frac{1}{1+\rho}} \right] \end{aligned} \quad (6)$$

For future reference, let us define the tilted distributions

$$Q_\rho(x|y) \triangleq \frac{Q(x|y)^{\frac{1}{1+\rho}}}{\sum_x Q(x|y)^{\frac{1}{1+\rho}}} \quad (7)$$

$$Q_\rho(y) \triangleq \frac{P(y) \left[ \sum_x Q(x|y)^{\frac{1}{1+\rho}} \right]^{1+\rho}}{\sum_y P(y) \left[ \sum_x Q(x|y)^{\frac{1}{1+\rho}} \right]^{1+\rho}}. \quad (8)$$

Differentiating with respect to  $\rho$  to find a stationary point, we can relate  $E_{x|y}(R, y)$  and  $R$  parametrically in terms of  $\rho$ :

$$R = \frac{\partial E_{0,x|y}(\rho, y)}{\partial \rho} = H(X_\rho|y)$$

where the second equality above can be verified with calculation, as mentioned in [4]. Now consider averaging over a memoryless  $y$ . Again, from [4], we have:

$$\begin{aligned} \frac{-\log P_e}{n} &\geq E_{x|y}(R) \triangleq \max_{0 \leq \rho \leq 1} \rho R - E_0(\rho), \\ E_0(\rho) &\triangleq \log \left( \sum_y P(y) \left[ \sum_x Q(x|y)^{\frac{1}{1+\rho}} \right]^{1+\rho} \right) \end{aligned} \quad (9)$$

We now relate  $E_{x|y}(R)$  and  $R$  in terms of  $\rho$  through differentiation:

$$R = \frac{\partial E_0(\rho)}{\partial \rho} = H(X_\rho|Y_\rho) \quad (10)$$

$$\begin{aligned} \Rightarrow E_{x|y}(R) &= \rho H(X_\rho|Y_\rho) - E_0(\rho) \\ \Rightarrow \frac{\partial}{\partial \rho} E_{x|y}(R) &= \rho \frac{\partial}{\partial \rho} H(X_\rho|Y_\rho) \\ \Rightarrow E'_{x|y}(R) &\triangleq \frac{dE_{x|y}(R)}{dR} = \frac{\partial E_{x|y}(R)}{\partial \rho} = \rho \end{aligned} \quad (11)$$

$$\Rightarrow E''_{x|y}(R) = \frac{\frac{\partial E'_{x|y}(R)}{\partial \rho}}{\frac{\partial R}{\partial \rho}} = \frac{1}{\frac{\partial}{\partial \rho} H(X_\rho|Y_\rho)} \quad (12)$$

where the second equality in (10) can be verified with tedious calculations, as mentioned in [4]. Note from (10),(7) (8), and (11) that

$$E'_{x|y}(H(X|Y)) = 0. \quad (13)$$

$$\begin{aligned} \Rightarrow \frac{\partial}{\partial \rho} \{H(X_\rho|Y_\rho)\} &= \sum_y \frac{\partial}{\partial \rho} \{Q_\rho(y)H(X_\rho|y)\} \\ &= \sum_y H(X_\rho|y) \frac{\partial Q_\rho(y)}{\partial \rho} + \sum_y Q_\rho(y) \frac{\partial}{\partial \rho} \{H(X_\rho|y)\}. \end{aligned} \quad (14)$$

To address  $\frac{\partial Q_\rho(y)}{\partial \rho}$ , note that for any differentiable function  $g$ ,  $g'(\rho) = g(\rho) \frac{d}{d\rho} \{\log g(\rho)\}$ .

$$\begin{aligned} \Rightarrow \log Q_\rho(y) &= \log P(y) + (1 + \rho) \log \left[ \sum_x Q(x|y)^{\frac{1}{1+\rho}} \right] \\ &\quad - \log \left[ \sum_y P(y) \left[ \sum_x Q(x|y)^{\frac{1}{1+\rho}} \right]^{1+\rho} \right] \\ &= \log P(y) + E_{0,x|y}(\rho, y) - E_0(\rho) \end{aligned} \quad (15)$$

$$\Rightarrow \frac{\partial \log Q_\rho(y)}{\partial \rho} = H(X_\rho|y) - H(X_\rho|Y_\rho) \quad (16)$$

$$\Rightarrow \frac{\partial Q_\rho(y)}{\partial \rho} \Big|_{\rho=0} = Q(y) [H(X|y) - H(X|Y)] \quad (17)$$

where (15) is due to (6), (9). As for  $\frac{\partial H(X_\rho|y)}{\partial \rho}$ , note that

$$H(X_\rho|y) = -D(Q_\rho(x|y) \parallel \mathbb{U}) + \log |\mathcal{X}| \quad (18)$$

$$\Rightarrow \frac{\partial}{\partial \rho} \{H(X_\rho|y)\} = -\frac{\partial}{\partial \rho} \{D(Q_\rho(x|y) \parallel \mathbb{U})\} \quad (19)$$

where  $\mathbb{U}(x) \triangleq \frac{1}{|\mathcal{X}|}$ . By [6], for any two distributions  $W_0$  and  $W_1$ , the distribution  $W_t$

$$W_t(x) \triangleq \frac{W_0(x)^{1-t} W_1(x)^t}{\sum_a W_0(a)^{1-t} W_1(a)^t} \quad (20)$$



relates the Kullback-Leibler divergence to the Fisher information  $F(t)$  according to:

$$\frac{d}{dt} \{D(W_t \| W_0)\} = tF(t), \quad (21)$$

$$F(t) \triangleq \sum_x \left( \frac{d \log W_t(x)}{dt} \right)^2 W_t(x) \quad (22)$$

We would like to characterize  $Q_\rho(x|y)$  in terms of a  $W_t$  of the form (20):

$0 \leq \rho \leq \infty$ $\rho = 0 : Q_\rho(x y) = Q(x y)$ $\rho = \infty : Q_\rho(x y) = \frac{1}{ \mathcal{X} }$	$\longleftrightarrow$ $t = \frac{1}{1+\rho}$	$0 \leq t \leq 1$ $t = 1 : W_t(x) = Q(x y)$ $t = 0 : W_t(x) = \frac{1}{ \mathcal{X} }$
---	---	--

$$\begin{aligned} \Rightarrow \frac{\partial H(X_\rho|y)}{\partial \rho} &= -\frac{\partial}{\partial \rho} \{D(Q_\rho(x|y) \| \mathbb{U})\} \\ &= -\frac{dt}{d\rho} \frac{d}{dt} \{D(W_t \| W_0)\} \Big|_{t=\frac{1}{1+\rho}} \\ &= \frac{1}{(1+\rho)^3} F\left(\frac{1}{1+\rho}\right) \\ \Rightarrow \frac{\partial H(X_\rho|y)}{\partial \rho} \Big|_{\rho=0} &= F(1) = -H(X|y)^2 + \sum_x Q(x|y) \log^2[Q(x|y)] \end{aligned} \quad (23)$$

Thus it follows from (12),(14), (17), and (23) that

$E''_{x y}(H(X Y)) = \frac{1}{-H(X Y)^2 + \sum_{x,y} Q(x,y) \log^2[Q(x y)]}$
--

## B Proof of Lemma 2.2

$$\begin{aligned} \log Q_\rho(x,y) &= \log Q_\rho(x|y) + \log Q_\rho(y) \\ &= \frac{1}{1+\rho} [\log Q(x|y) - E_{0,x|y}(\rho,y)] + \log P(y) + E_{0,x|y}(\rho,y) - E_0(\rho) \\ \Rightarrow \frac{\partial \log Q_\rho(x,y)}{\partial \rho} &= \frac{-1}{1+\rho} \frac{\partial E_{0,x|y}(\rho,y)}{\partial \rho} + \frac{-1}{(1+\rho)^2} [\log Q(x|y) - E_{0,x|y}(\rho,y)] \\ &\quad + \frac{\partial E_{0,x|y}(\rho,y)}{\partial \rho} - \frac{\partial E_0(\rho)}{\partial \rho} \\ &= \frac{-1}{1+\rho} [H(X_\rho|y) + \log Q_\rho(x|y)] + H(X_\rho|y) - H(X_\rho|Y_\rho) \end{aligned} \quad (24)$$

Thus

$$\begin{aligned}\left(\frac{\partial \log Q_{\rho}(x,y)}{\partial \rho}\right)^2 \Big|_{\rho=0} &= (\log Q(x|y) + H(X|Y))^2 \\ &= H(X|Y)^2 + 2 \log Q(x|y) H(X|Y) + \log^2 [Q(x|y)] \\ \Rightarrow F(\rho) \Big|_{\rho=0} &= -H(X|Y)^2 + \sum_{x,y} Q(x,y) \log^2 [Q(x|y)]\end{aligned}$$