

Properties of a Recent Upper Bound to the Mismatch Capacity

Ehsan Asadi Kangarshahi
University of Cambridge
ea460@cam.ac.uk

Albert Guillén i Fàbregas
ICREA & Universitat Pompeu Fabra
University of Cambridge
guillen@ieee.org

Abstract—We study several properties of the upper bound on the mismatch capacity problem we recently proposed. In particular, we show that the bound can be cast as a convex-concave saddlepoint problem enabling efficient computation. Moreover, as opposed to multiple achievability bound in the literature, we show that the multiletter version of this bound does not yield any gain. In addition, we show a necessary condition for the mismatch capacity to be strictly smaller than the channel capacity for binary-input channels.

I. INTRODUCTION AND PRELIMINARIES

We consider reliable communication over a discrete memoryless channel (DMC) W with a given decoding metric [1], [2]. This problem arises when the decoder uses a suboptimal decoding rule due to limited computational resources, or imperfect channel estimation. Moreover, it is shown in [2] that important problems in information theory, like zero-error capacity of a channel can be cast as instances of mismatched decoding problem. Multiple achievability results have been reported in the literature [1]–[4] (see also [5]). These results are derived by random-coding techniques, i.e. analyzing the average probability of error of mismatched decoder over an ensemble of codebooks. On the other hand, the only single-letter converse was given in [6], where it was claimed that for binary-input DMCs, the mismatch capacity was the achievable rate derived in [3], [4]. Reference [7] provided a counterexample to this converse invalidating its claim. Multiletter converse results were proposed in [8].

We assume input and output alphabets are $\mathcal{X} = \{1, 2, \dots, J\}$ and $\mathcal{Y} = \{1, 2, \dots, K\}$, respectively, with $J, K < \infty$. We denote the channel transition probability by $W(k|j), k \in \mathcal{Y}, j \in \mathcal{X}$. A codebook \mathcal{C}_n is defined as a set of M sequences $\mathcal{C}_n = \{\mathbf{x}(1), \mathbf{x}(2), \dots, \mathbf{x}(M)\}$, where $\mathbf{x}(m) = (x_1(m), x_2(m), \dots, x_n(m)) \in \mathcal{X}^n$, for $m \in \{1, 2, \dots, M\}$. A message $m \in \{1, 2, \dots, M\}$ is chosen equiprobably and $\mathbf{x}(m)$ is sent over the channel. The channel produces a noisy observation $\mathbf{y} = (y_1, y_2, \dots, y_n) \in \mathcal{Y}^n$ according to $W^n(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^n W(y_i|x_i)$. Upon observing $\mathbf{y} \in \mathcal{Y}^n$ the decoder produces an estimate of the transmitted message $\hat{m} \in \{1, 2, \dots, M\}$. The decoder that minimizes the

error probability is the maximum-likelihood (ML) decoder, that produces the message estimate \hat{m} according to

$$\hat{m} = \arg \max_{i \in \{1, 2, \dots, M\}} W^n(\mathbf{y}|\mathbf{x}(i)). \quad (1)$$

Rate $R > 0$ is achievable if for any $\epsilon > 0$ there exists a sequence of length- n codebooks $\{\mathcal{C}_n\}_{n=1}^{\infty}$ such that $|\mathcal{C}_n| \geq 2^{n(R-\epsilon)}$, and $\liminf_{n \rightarrow \infty} P_e(\mathcal{C}_n) = 0$. The capacity of W , denoted by $C(W)$, is defined as the largest achievable rate.

In multiple practical scenarios, it is not possible to use a decoder based on W^n and instead, the decoder produces the message estimate \hat{m} as

$$\hat{m} = \arg \max_{i \in \{1, 2, \dots, M\}} d(\mathbf{x}(i), \mathbf{y}), \quad (2)$$

where,

$$d(\mathbf{x}(i), \mathbf{y}) = \sum_{\ell=1}^n d(x_\ell(i), y_\ell) \quad (3)$$

The mismatch capacity $C_d(W)$ is defined as the largest achievable rate when the decoder is (2). Recently, we have shown that $C_d(W)$ is upper bounded by the following quantity,

$$\bar{R}_d(W) = \max_{P_X} \min_{\substack{P_{Y\hat{Y}|X} \in \mathcal{M}_{\max}(d) \\ P_{Y|X} = W}} I(P_X, P_{\hat{Y}|X}) \quad (4)$$

where $I(P_X, P_{\hat{Y}|X}) \triangleq I(X; \hat{Y})$ and the set $\mathcal{M}_{\max}(d)$ is given in the following definition.

Definition 1: Let $P_{Y\hat{Y}|X}$ be a joint conditional distribution and define the set $\mathcal{S}(k_1, k_2) \triangleq \{i \in \mathcal{X} | i = \arg \max_{i' \in \mathcal{X}} d(i', k_2) - d(i', k_1)\}$. We say that $P_{Y\hat{Y}|X}$ is a maximal joint conditional distribution if for all $(j, k_1, k_2) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y}$,

$$P_{Y\hat{Y}|X}(k_1, k_2|j) = 0 \text{ if } j \notin \mathcal{S}(k_1, k_2). \quad (5)$$

For a given decoding metric d , we define the set of maximal joint conditional distributions to be $\mathcal{M}_{\max}(d)$.

In this paper we study some properties of the upper bound. Specifically, in Section II we show that computing our upper bound is a convex-concave saddlepoint problem and we derive the optimality KKT conditions. In Section III, we show that the multiletter version of the upper bound coincides with the single-letter one. In Section IV we derive a sufficient condition for $C_d(W) < C(W)$ for binary-input channels.

This work was supported in part by the European Research Council under Grant 725411, and by the Spanish Ministry of Economy and Competitiveness under Grant TEC2016-78434-C3-1-R.

II. CONVEXITY ANALYSIS

In this section, we show that the optimization (4) is a convex-concave saddlepoint problem. First we argue that the constraints induce a convex set.

Lemma 1: For any channel W and metric d , the set of joint conditional distributions $P_{Y\hat{Y}|X}$ satisfying both $P_{Y\hat{Y}|X} \in \mathcal{M}_{\max}(d)$ and $P_{Y|X} = W$, is a convex set.

Proof: Let $P_{Y\hat{Y}|X}$ and $P'_{Y\hat{Y}|X}$ both satisfy the above constraints. Now for any $0 < \lambda < 1$ we have,

$$\lambda P_{Y\hat{Y}|X} + (1 - \lambda) P'_{Y\hat{Y}|X} = W. \quad (6)$$

In addition, if for some k_1, k_2 we have $j \notin \mathcal{S}(k_1, k_2)$, both $P_{Y\hat{Y}|X}(k_1, k_2|j)$ and $P'_{Y\hat{Y}|X}(k_1, k_2|j)$ are equal to zero, and so is any linear combination of them. Therefore,

$$\lambda P_{Y\hat{Y}|X} + (1 - \lambda) P'_{Y\hat{Y}|X} \in \mathcal{M}_{\max}(d). \quad (7)$$

Moreover, $I(P_X, P_{Y\hat{Y}|X})$ is convex in terms of $P_{Y\hat{Y}|X}$, and concave in terms of P_X . Since $P_{Y\hat{Y}|X}$ is a linear function of $P_{Y\hat{Y}|X}$, we get that $I(P_X, P_{Y\hat{Y}|X})$ is convex in terms of $P_{Y\hat{Y}|X}$. Therefore from the minimax theorem [9] we get,

$$\bar{R}_d(W) = \max_{P_X} \min_{\substack{P_{Y\hat{Y}|X} \in \mathcal{M}_{\max}(d) \\ P_{Y|X} = W}} I(P_X, P_{Y\hat{Y}|X}) \quad (8)$$

$$= \min_{\substack{P_{Y\hat{Y}|X} \in \mathcal{M}_{\max}(d) \\ P_{Y|X} = W}} \max_{P_X} I(P_X, P_{Y\hat{Y}|X}) \quad (9)$$

$$= \min_{\substack{P_{Y\hat{Y}|X} \in \mathcal{M}_{\max}(d) \\ P_{Y|X} = W}} C(P_{Y\hat{Y}|X}). \quad (10)$$

The rest of this section is devoted to deriving the KKT conditions for the optimization problem in (4). Given that $I(P_X, P_{Y\hat{Y}|X})$ is convex in $P_{Y\hat{Y}|X}$, and concave in P_X , then the KKT conditions are sufficient for global optimality. For convenience, we define $f(P_X, P_{Y\hat{Y}|X}) \triangleq I(P_X, P_{Y\hat{Y}|X})$ and rewrite the optimization problem in (4) as,

$$\bar{R}_d(W) = \max_{P_X} \min_{\substack{P_{Y\hat{Y}|X} \in \mathcal{M}_{\max}(d) \\ P_{Y|X} = W}} f(P_X, P_{Y\hat{Y}|X}). \quad (11)$$

Let $\hat{P}_X, \hat{P}_{Y\hat{Y}|X}$ be the optimal input and joint conditional distributions in (11) and $\hat{q}_{\hat{Y}}$ be the output distribution induced by \hat{P}_X and $\hat{P}_{Y\hat{Y}|X}$. Then for \hat{P}_X we have the following constraints:

$$\hat{P}_X(j) \geq 0, \quad \forall j \in \mathcal{X} \quad (12)$$

$$\sum_{j \in \mathcal{X}} \hat{P}_X(j) = 1. \quad (13)$$

Let $\mu_j, j = 1, 2, \dots, J$ be the Lagrange multipliers corresponding the inequalities in (12) and ρ be the Lagrange multiplier corresponding to (13). Therefore, from stationarity we have,

$$\left. \frac{\partial}{\partial P_X(j)} f(P_X, \hat{P}_{Y\hat{Y}|X}) \right|_{P_X = \hat{P}_X} = \rho + \mu_j \quad (14)$$

and from the complementary slackness [10] we have $\mu_j \hat{P}_X(j) = 0$ and from the dual feasibility we have $\mu_j \geq 0$ which leads to the separation of the equations of in two cases. If $\hat{P}_X(j) > 0$

$$\left. \frac{\partial}{\partial P_X(j)} f(P_X, \hat{P}_{Y\hat{Y}|X}) \right|_{P_X = \hat{P}_X} = \rho, \quad (15)$$

while when $\hat{P}_X(j) = 0$ we have

$$\left. \frac{\partial}{\partial P_X(j)} f(P_X, \hat{P}_{Y\hat{Y}|X}) \right|_{P_X = \hat{P}_X} \leq \rho. \quad (16)$$

Note that, because there is no other constraint on μ_j , all of the KKT conditions are summarized in (16) and (15). Moreover, computing the derivatives in (15) and (16) gives

$$\begin{aligned} \left. \frac{\partial}{\partial P_X(j)} f(P_X, \hat{P}_{Y\hat{Y}|X}) \right|_{P_X = \hat{P}_X} &= \sum_{k \in \mathcal{Y}} \hat{P}_{Y\hat{Y}|X}(k|j) \log \frac{\hat{P}_{Y\hat{Y}|X}(k|j)}{\hat{q}_{\hat{Y}}(k)} - 1. \end{aligned} \quad (17)$$

As for $\hat{P}_{Y\hat{Y}|X}$, we have the following constraints. For all $j, k_1, k_2 \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y}$,

$$\hat{P}_{Y\hat{Y}|X}(k_1, k_2|j) \geq 0, \quad (18)$$

$$\hat{P}_{Y\hat{Y}|X}(k_1, k_2|j) = 0, \text{ if } j \notin \mathcal{S}(k_1, k_2) \quad (19)$$

where (18) corresponds to $\hat{P}_{Y\hat{Y}|X}(k_1, k_2|j)$ being a distribution and (19) corresponds to $\hat{P}_{Y\hat{Y}|X}(k_1, k_2|j) \in \mathcal{M}_{\max}(d)$. Moreover from the constraint $P_{Y|X} = W$ we get for all $j, k_1 \in \mathcal{X} \times \mathcal{Y}$

$$\sum_{k_2} \hat{P}_{Y\hat{Y}|X}(k_1, k_2|j) = W(k_1|j). \quad (20)$$

For the ease of notation, we skip the step of explicitly considering a Lagrange multiplier for (18). However, after simplification, The following KKT conditions are equivalent to the full KKT conditions considering a Lagrange multiplier for (18). Details follow similarly to the above derivation. If we use a Lagrange multiplier λ_{j, k_1} for each of the conditions in (20), we have when $\hat{P}_{Y\hat{Y}|X}(k_1, k_2|j) > 0$

$$\left. \frac{\partial}{\partial P_{Y\hat{Y}|X}(k_1, k_2|j)} f(\hat{P}_X, P_{Y\hat{Y}|X}) \right|_{P_{Y\hat{Y}|X} = \hat{P}_{Y\hat{Y}|X}} = \lambda_{j, k_1} \quad (21)$$

and when $\hat{P}_{Y\hat{Y}|X}(k_1, k_2|j) = 0$ and $j \in \mathcal{S}(k_1, k_2)$ we have

$$\left. \frac{\partial}{\partial P_{Y\hat{Y}|X}(k_1, k_2|j)} f(\hat{P}_X, P_{Y\hat{Y}|X}) \right|_{P_{Y\hat{Y}|X} = \hat{P}_{Y\hat{Y}|X}} \geq \lambda_{j, k_1}. \quad (22)$$

Explicitly computing the derivative gives

$$\left. \frac{\partial}{\partial P_{Y\hat{Y}|X}(k_1, k_2|j)} f(\hat{P}_X, P_{Y\hat{Y}|X}) \right|_{P_{Y\hat{Y}|X} = \hat{P}_{Y\hat{Y}|X}} \quad (23)$$

$$= \hat{P}_X(j) \log \frac{\hat{P}_{Y\hat{Y}|X}(k_2|j)}{\hat{q}_{\hat{Y}}(k_2)}. \quad (24)$$

Summarizing, for the KKT optimality conditions of we get the following inequalities

1) For $\hat{P}_X(j) > 0$,

$$\sum_{k \in \mathcal{Y}} \hat{P}_{\hat{Y}|X}(k|j) \log \frac{\hat{P}_{\hat{Y}|X}(k|j)}{\hat{q}_{\hat{Y}}(k)} = 1 + \rho, \quad (25)$$

2) For $\hat{P}_X(j) = 0$,

$$\sum_{k \in \mathcal{Y}} \hat{P}_{\hat{Y}|X}(k|j) \log \frac{\hat{P}_{\hat{Y}|X}(k|j)}{\hat{q}_{\hat{Y}}(k)} \leq 1 + \rho, \quad (26)$$

3) For $\hat{P}_{Y\hat{Y}|X}(k_1, k_2|j) > 0$,

$$\hat{P}_X(j) \log \frac{\hat{P}_{\hat{Y}|X}(k_2|j)}{\hat{q}_{\hat{Y}}(k_2)} = \lambda_{j, k_1}, \quad (27)$$

4) For $\hat{P}_{Y\hat{Y}|X}(k_1, k_2|j) = 0$ and $j \in \mathcal{S}(k_1, k_2)$,

$$\hat{P}_X(j) \log \frac{\hat{P}_{\hat{Y}|X}(k_2|j)}{\hat{q}_{\hat{Y}}(k_2)} \geq \lambda_{j, k_1}. \quad (28)$$

In the next section, we employ the above KKT conditions to analyze the multiletter version of our bound.

III. MULTILETTER BOUND

In this section, we study the multiletter extension of the bound (4). In particular, we show that the multiletter version cannot improve on the single-letter bound. We define the ℓ -letter decoding metric $d^{(\ell)} : \mathcal{X}^\ell \times \mathcal{Y}^\ell \rightarrow \mathbb{R}$ as follows

$$d^{(\ell)}((x_1, x_2, \dots, x_\ell), (y_1, y_2, \dots, y_\ell)) = \sum_{i=1}^{\ell} d(x_i, y_i). \quad (29)$$

This decoding metric definition is consistent with the additive decoder we have defined in (3). We denote $\mathbf{j} \in \mathcal{X}^\ell$ and $\mathbf{k} \in \mathcal{Y}^\ell$ as the ℓ -letter inputs and outputs, respectively. Let $W^{(\ell)}$ denote a DMC over input alphabet \mathcal{X}^ℓ and output alphabet \mathcal{Y}^ℓ with the channel rule $W^{(\ell)}((y_1, y_2, \dots, y_\ell)|(x_1, x_2, \dots, x_\ell)) = \prod_{i=1}^{\ell} W(y_i|x_i)$. Additionally, we define $P_X^{(\ell)}$ and $P_{Y\hat{Y}|X}^{(\ell)}$ accordingly

$$P_X^{(\ell)}(x_1, \dots, x_\ell) = \prod_{i=1}^{\ell} P_X(x_i) \quad (30)$$

$$P_{Y\hat{Y}|X}^{(\ell)}((y_1, y_2, \dots, y_\ell), (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_\ell)|(x_1, x_2, \dots, x_\ell)) = \prod_{i=1}^{\ell} P_{Y\hat{Y}|X}(y_i, \hat{y}_i|x_i) \quad (31)$$

X^ℓ and Y^ℓ, \hat{Y}^ℓ denote random variables defined on alphabets $\mathcal{X}^\ell, \mathcal{Y}^\ell$ and \mathcal{Y}^ℓ , respectively. Moreover, $\mathcal{S}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2)$ is defined as

$$\mathcal{S}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2) \triangleq \{i \in \mathcal{X}^\ell \mid i = \arg \max_{i' \in \mathcal{X}^\ell} d^{(\ell)}(i', \mathbf{k}_2) - d^{(\ell)}(i', \mathbf{k}_1)\}. \quad (32)$$

In the following lemma we characterize the sets $\mathcal{S}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2)$ and relate them to $\mathcal{S}(k_{1,i}, k_{2,i}), i = 1, 2, \dots, \ell$.

Lemma 2:] For $\mathbf{j} \in \mathcal{X}^\ell, \mathbf{k}_1 \in \mathcal{Y}^\ell, \mathbf{k}_2 \in \mathcal{Y}^\ell$ we have that $\mathbf{j} \in \mathcal{S}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2)$ if and only if for all $1 \leq i \leq \ell$ we have

$$j_i \in \mathcal{S}(k_{1,i}, k_{2,i}). \quad (33)$$

Proof: We have

$$\arg \max_{\mathbf{j} \in \mathcal{X}^\ell} d^{(\ell)}(\mathbf{j}, \mathbf{k}_2) - d^{(\ell)}(\mathbf{j}, \mathbf{k}_1) \quad (34)$$

$$= \arg \max_{\mathbf{j} \in \mathcal{X}^\ell} \sum_{i=1}^{\ell} d(j_i, k_{2,i}) - d(j_i, k_{1,i}) \quad (35)$$

$$= \arg \max_{(j_1, j_2, \dots, j_\ell) \in \mathcal{X}^\ell} \sum_{i=1}^{\ell} d(j_i, k_{2,i}) - d(j_i, k_{1,i}) \quad (36)$$

From (36) we get that if $(j_1, j_2, \dots, j_\ell) \in \mathcal{S}(\mathbf{k}_1, \mathbf{k}_2)$ then for all $1 \leq i \leq \ell$ we should have $j_i \in \mathcal{S}(k_{1,i}, k_{2,i})$. Therefore,

$$\mathcal{S}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2) = \mathcal{S}(k_{1,1}, k_{2,1}) \times \mathcal{S}(k_{1,2}, k_{2,2}) \times \dots \times \mathcal{S}(k_{1,\ell}, k_{2,\ell}). \quad (37)$$

For the above ℓ -letter alphabets and distributions, the construction and analysis of the bound remains unchanged. Therefore, (4) remains valid for its ℓ -letter extension, which can be written as

$$\begin{aligned} \bar{R}_d^{(\ell)}(W) &\triangleq \frac{1}{\ell} \bar{R}_{d^{(\ell)}}(W^{(\ell)}) \\ &= \frac{1}{\ell} \max_{P_{X^\ell}} \min_{\substack{P_{Y^\ell \hat{Y}^\ell | X^\ell} \in \mathcal{M}_{\max}^{(d^{(\ell)})} \\ P_{Y^\ell \hat{Y}^\ell | X^\ell} = W^{(\ell)}}} I(p_{X^\ell}, P_{Y^\ell \hat{Y}^\ell | X^\ell}). \end{aligned} \quad (38)$$

We have the following result.

Proposition 1:

$$\bar{R}_d^{(\ell)}(W) = \bar{R}_d(W). \quad (40)$$

Proof: Given that $I(P_X, P_{Y\hat{Y}|X})$ is convex in $P_{Y\hat{Y}|X}$, and concave in P_X , the KKT conditions are also sufficient for global optimality. Similarly, $f(p_{X^\ell}, P_{Y^\ell \hat{Y}^\ell | X^\ell})$ is convex in terms of p_{X^ℓ} and concave in terms of $P_{Y^\ell \hat{Y}^\ell | X^\ell}$. Here we use the optimality conditions derived in the past section to show that if $P_X^*, P_{Y\hat{Y}|X}^*$ are the optimal distributions for the single-letter bound then $\hat{P}_X^{(\ell)}, \hat{P}_{Y\hat{Y}|X}^{(\ell)}$ defined in (30) and (31) are optimal distributions for the multiletter version. As a result, if we find a feasible pair $P_{Y^\ell \hat{Y}^\ell | X^\ell}, P_{X^\ell}$ such that when fixing $P_{Y^\ell \hat{Y}^\ell | X^\ell}$, the input distribution P_{X^ℓ} is a maximizer of $f(\cdot, P_{Y^\ell \hat{Y}^\ell | X^\ell})$, and when fixing P_{X^ℓ} , the joint conditional distribution $P_{Y^\ell \hat{Y}^\ell | X^\ell}$ is a minimizer of $f(p_{X^\ell}, \cdot)$, then the pair $(P_{Y^\ell \hat{Y}^\ell | X^\ell}, P_{X^\ell})$ is a saddlepoint.

We need to show that if $\hat{P}_X, \hat{P}_{Y\hat{Y}|X}$ is a saddlepoint for the single-letter case, then, $\hat{P}_X^{(\ell)}, \hat{P}_{Y\hat{Y}|X}^{(\ell)}$ is a saddlepoint for the multiletter bound. Based on the aforementioned argument, it is sufficient to show that $\hat{P}_{Y\hat{Y}|X}^{(\ell)}$ is a minimizer of (39) by fixing

$\widehat{P}_X^{(\ell)}$. This is because it is known that $\frac{1}{\ell}C(\widehat{P}_{\widehat{Y}|X}^{(\ell)}) = C(P_{\widehat{Y}|X})$, i.e., the product distribution $P_X^{*(\ell)}$ achieves $C(\widehat{P}_{\widehat{Y}|X}^{(\ell)})$.

In the following lemma we prove that by fixing $\widehat{P}_X^{(\ell)}$, then $\widehat{P}_{\widehat{Y}|X}^{(\ell)}$ satisfies the KKT conditions and hence, it is a minimizer of (39). Before stating the result we recall that the multiletter counterparts of the single-letter KKT conditions given in (27) and (28) hold. Moreover, as in the single-letter case, the multiletter KKT conditions are sufficient for global optimality, because the function $f(\widehat{P}_X^{(\ell)}, \cdot)$ is concave. Using Lemma 3 below completes the proof. ■

Lemma 3: Let $\widehat{P}_X, \widehat{P}_{\widehat{Y}|X}$ be a saddlepoint for optimization problem (4). Set $P_{X^\ell} = \widehat{P}_X^{(\ell)}$. Then, the joint conditional distribution $\widehat{P}_{\widehat{Y}|X}^{(\ell)}$ is a minimizer of

$$\min_{\substack{P_{Y^\ell \widehat{Y}^\ell | X^\ell} \in \mathcal{M}_{\max}(d^{(\ell)}) \\ P_{Y^\ell | X^\ell} = W^{(\ell)}}} f(\widehat{P}_X^{(\ell)}, P_{Y^\ell \widehat{Y}^\ell | X^\ell}). \quad (41)$$

Proof: We should show that by setting $P_{X^\ell} = \widehat{P}_X^{(\ell)}$, the multiletter versions of the KKT conditions (27) and (28) hold for $\widehat{P}_{\widehat{Y}|X}^{(\ell)}$. Generalizing the conditions of (27) and (28) to the multiletter case, and setting $P_{Y^\ell \widehat{Y}^\ell | X^\ell} = \widehat{P}_{\widehat{Y}|X}^{(\ell)}$, we should show that for all $\mathbf{j}, \mathbf{k}_1 \in \mathcal{X}^\ell \times \mathcal{Y}^\ell$ there exist $\lambda_{\mathbf{j}, \mathbf{k}_1}$ such that the conditions below are fulfilled. If we show this, then the Lemma is proved because these are precisely the conditions for the minimizer of (41).

i) When $\widehat{P}_{\widehat{Y}|X}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2 | \mathbf{j}) > 0$ we must have,

$$\left. \frac{\partial}{\partial P_{Y^\ell \widehat{Y}^\ell | X^\ell}(\mathbf{k}_1, \mathbf{k}_2 | \mathbf{j})} f(\widehat{P}_X^{(\ell)}, P_{Y^\ell \widehat{Y}^\ell | X^\ell}) \right|_{P_{Y^\ell \widehat{Y}^\ell | X^\ell} = \widehat{P}_{\widehat{Y}|X}^{(\ell)}} = \lambda_{\mathbf{j}, \mathbf{k}_1}. \quad (42)$$

ii) When $\widehat{P}_{\widehat{Y}|X}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2 | \mathbf{j}) = 0$ and $\mathbf{j} \in \mathcal{S}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2)$ we must have that,

$$\left. \frac{\partial}{\partial P_{Y^\ell \widehat{Y}^\ell | X^\ell}(\mathbf{k}_1, \mathbf{k}_2 | \mathbf{j})} f(\widehat{P}_X^{(\ell)}, P_{Y^\ell \widehat{Y}^\ell | X^\ell}) \right|_{P_{Y^\ell \widehat{Y}^\ell | X^\ell} = \widehat{P}_{\widehat{Y}|X}^{(\ell)}} \geq \lambda_{\mathbf{j}, \mathbf{k}_1}. \quad (43)$$

Similarly to (23), the derivative in (42) and (43) is,

$$\left. \frac{\partial}{\partial P_{Y^\ell \widehat{Y}^\ell | X^\ell}(\mathbf{k}_1, \mathbf{k}_2 | \mathbf{j})} f(\widehat{P}_X^{(\ell)}, P_{Y^\ell \widehat{Y}^\ell | X^\ell}) \right|_{P_{Y^\ell \widehat{Y}^\ell | X^\ell} = \widehat{P}_{\widehat{Y}|X}^{(\ell)}} = \widehat{P}_X^{(\ell)}(\mathbf{j}) \log \frac{\widehat{P}_{\widehat{Y}|X}^{(\ell)}(\mathbf{k}_1 | \mathbf{j})}{\widehat{q}_{\widehat{Y}}^{(\ell)}(\mathbf{k}_1)} \quad (44)$$

which, by using that $P_{Y^\ell \widehat{Y}^\ell | X^\ell} = \widehat{P}_{\widehat{Y}|X}^{(\ell)}$, $\widehat{P}_X^{(\ell)}$ and $\widehat{q}_{\widehat{Y}}^{(\ell)}$ are product distributions, gives,

$$\begin{aligned} & \widehat{P}_X^{(\ell)}(\mathbf{j}) \log \frac{\widehat{P}_{\widehat{Y}|X}^{(\ell)}(\mathbf{k}_1 | \mathbf{j})}{\widehat{q}_{\widehat{Y}}^{(\ell)}(\mathbf{k}_1)} \\ &= \widehat{P}_X(j_1) \widehat{P}_X(j_2) \cdots \widehat{P}_X(j_\ell) \left(\sum_{i=1}^{\ell} \log \frac{\widehat{P}_{\widehat{Y}|X}(k_{2,i} | j_i)}{\widehat{q}_{\widehat{Y}}(k_{2,i})} \right) \end{aligned} \quad (45)$$

In order to show that there exist some coefficients $\lambda_{\mathbf{j}, \mathbf{k}_1}$ satisfying both (42) and (43), we make a particular choice and show that the choice satisfies both (42) and (43). To this end, define

$$\lambda_{\mathbf{j}, \mathbf{k}_1} = \begin{cases} 0 & \widehat{P}_X(\mathbf{j}) = 0 \\ \prod_{i=1}^{\ell} \widehat{P}_X(j_i) \left(\sum_{i=1}^{\ell} \frac{\lambda_{j_i, k_{1,i}}}{\widehat{P}_X(j_i)} \right) & \widehat{P}_X(\mathbf{j}) \neq 0 \end{cases} \quad (46)$$

where $\lambda_{j_i, k_{1,i}}$ is the single-letter Lagrange multiplier corresponding to j_i and $k_{1,i}$.

Now, excluding the cases where $\widehat{P}_X(j_1) \widehat{P}_X(j_2) \cdots \widehat{P}_X(j_\ell) = 0$ where from (45), (42) and (43) the KKT conditions clearly hold, we have two cases i) When $\widehat{P}_{\widehat{Y}|X}^{(\ell)}(\mathbf{j}, \mathbf{k}_1, \mathbf{k}_2) > 0$, then for all $1 \leq i \leq \ell$ we must have $\widehat{P}_{\widehat{Y}|X}(k_{1,i}, k_{2,i} | j_i) > 0$ and therefore, (27) is valid. We have to verify that this implies that (42) is also valid. Thus,

$$\begin{aligned} & \left. \frac{\partial}{\partial P_{Y^\ell \widehat{Y}^\ell | X^\ell}(\mathbf{k}_1, \mathbf{k}_2 | \mathbf{j})} f(\widehat{P}_X, P_{Y^\ell \widehat{Y}^\ell | X^\ell}) \right|_{P_{Y^\ell \widehat{Y}^\ell | X^\ell} = \widehat{P}_{\widehat{Y}|X}^{(\ell)}} \\ &= \widehat{P}_X(j_1) \widehat{P}_X(j_2) \cdots \widehat{P}_X(j_\ell) \left(\sum_{i=1}^{\ell} \log \frac{\widehat{P}_{\widehat{Y}|X}(k_{2,i} | j_i)}{\widehat{q}_{\widehat{Y}}(k_{2,i})} \right) \end{aligned} \quad (47)$$

$$\begin{aligned} &= \widehat{P}_X(j_1) \widehat{P}_X(j_2) \cdots \widehat{P}_X(j_\ell) \left(\sum_{i=1}^{\ell} \frac{\lambda_{j_i, k_{1,i}}}{\widehat{P}_X(j_i)} \right) \\ &= \lambda_{\mathbf{j}, \mathbf{k}_1} \end{aligned} \quad (49)$$

where (48) holds from the single-letter optimality in (27).

ii) When $\widehat{P}_{\widehat{Y}|X}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2 | \mathbf{j}) = 0$ and $\mathbf{j} \in \mathcal{S}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2)$, as a result of Lemma 2, we have that $\mathcal{S}^{(\ell)}(\mathbf{k}_1, \mathbf{k}_2)$ is a product set, i.e., for all $1 \leq i \leq \ell$,

$$j_i \in \mathcal{S}(k_{1,i}, k_{2,i}). \quad (50)$$

Moreover, either $\widehat{P}_{\widehat{Y}|X}(k_{1,i}, k_{2,i} | j_i) > 0$ where (27) is satisfied or $\widehat{P}_{\widehat{Y}|X}(k_{1,i}, k_{2,i} | j_i) = 0$ where (28) is satisfied. Now, with these assumptions, we should verify that (43) is valid. We have,

$$\begin{aligned} & \left. \frac{\partial}{\partial P_{Y^\ell \widehat{Y}^\ell | X^\ell}(\mathbf{k}_1, \mathbf{k}_2 | \mathbf{j})} f(\widehat{P}_X, P_{Y^\ell \widehat{Y}^\ell | X^\ell}) \right|_{P_{Y^\ell \widehat{Y}^\ell | X^\ell} = \widehat{P}_{\widehat{Y}|X}^{(\ell)}} \\ &= \widehat{P}_X(j_1) \widehat{P}_X(j_2) \cdots \widehat{P}_X(j_\ell) \left(\sum_{i=1}^{\ell} \log \frac{\widehat{P}_{\widehat{Y}|X}(k_{2,i} | j_i)}{\widehat{q}_{\widehat{Y}}(k_{2,i})} \right) \end{aligned} \quad (51)$$

$$\geq \widehat{P}_X(j_1) \widehat{P}_X(j_2) \cdots \widehat{P}_X(j_\ell) \left(\sum_{i=1}^{\ell} \frac{\lambda_{j_i, k_{1,i}}}{\widehat{P}_X(j_i)} \right) \quad (52)$$

$$= \lambda_{\mathbf{j}, \mathbf{k}_1} \quad (53)$$

where (52) is true because of the single-letter optimality in (27) and (28). ■

IV. BINARY-INPUT CHANNELS

In [2], the authors state that for any DMC and decoding metric $d(x, y)$, the mismatch capacity $C_d(W)$ remains unaltered for a decoder with metric $\tilde{d}(x, y) = d(x, y) + a(x) + b(y)$, where $a(x), b(y)$ are functions of the input and output, respectively. This property suggests that for binary-input channels, the mismatch capacity $C_d(W)$ is only a function of the metric difference $d(1, y) - d(2, y)$. In this section, we show a necessary condition for $C_d(W) < C(W)$ for binary-input channels based on the above observation.

Definition 2: We say that two sequences $\{\alpha_i\}_{i=1}^K$ and $\{\beta_i\}_{i=1}^K$ have the same order if for all $1 \leq i_1, i_2 \leq K$

$$\alpha_{i_1} \geq \alpha_{i_2} \Rightarrow \beta_{i_1} \geq \beta_{i_2}. \quad (54)$$

We have the following result.

Theorem 1: Assume that $W(k|j) > 0$, for all $j = 1, 2, k = 1, \dots, K$. If the sequences $\{\log W(k|1) - \log W(k|2)\}_{k=1}^K$ and $\{d(1, k) - d(2, k)\}_{k=1}^K$ do not have the same order, then $\bar{R}_d(W) < C(W)$.

Proof: Without loss of generality, we assume that the sequence $\{d(1, k) - d(2, k)\}_{k=1}^K$ is non-decreasing, i.e., for $k_1 \leq k_2$,

$$d(1, k_1) - d(2, k_1) \leq d(1, k_2) - d(2, k_2). \quad (55)$$

This assumption simplifies the evaluation of the sets $\mathcal{S}(\cdot, \cdot)$. For $k_1 = k_2$ we have $\mathcal{S}(k_1, k_2) = \{1, 2\}$. Moreover, when $k_1 < k_2$ from (55) and Definition 1, we have that $1 \in \mathcal{S}(k_1, k_2)$ and $2 \in \mathcal{S}(k_2, k_1)$.

We prove a slightly stronger result. In particular, we prove that the condition $C_d(W) = C(W)$ implies that sequences

$$\left\{ \hat{P}_X(1) \log \frac{W(k|1)}{\hat{q}_Y(k)} \right\}_{k=1}^K, \left\{ -\hat{P}_X(2) \log \frac{W(k|2)}{\hat{q}_Y(k)} \right\}_{k=1}^K \quad (56)$$

both should have the same order as the decoding metric difference sequence $\{d(1, k) - d(2, k)\}_{k=1}^K$, where recall that the notation \hat{P}_X refers to the capacity-achieving distribution of W .

Now assume that $C_d(W) = C(W)$. Therefore, $\hat{P}_X, P_{Y\hat{Y}|X} = P_{Y|X}$ must be a saddle point of (9). As a result, the KKT conditions in (27) (28) must hold. Observe that

$$P_{Y\hat{Y}|X}(k_1, k_2|j) = \begin{cases} W(k_1|j) & k_1 = k_2 \\ 0 & k_1 \neq k_2. \end{cases} \quad (57)$$

Therefore, combining the KKT conditions in (27) (28) we have,

1) If $k_1 = k_2$, for both $j = 1, 2$ we have

$$\hat{P}_X(j) \log \frac{W(k_1|j)}{\hat{q}_Y(k_1)} = \lambda_{j, k_1} \quad (58)$$

2) If $k_1 < k_2$ we know $1 \in \mathcal{S}(k_1, k_2)$ and $2 \in \mathcal{S}(k_2, k_1)$

$$\hat{P}_X(1) \log \frac{W(k_2|1)}{\hat{q}_Y(k_2)} \geq \lambda_{1, k_1} \quad (59)$$

$$\hat{P}_X(2) \log \frac{W(k_1|2)}{\hat{q}_Y(k_1)} \geq \lambda_{2, k_2} \quad (60)$$

Therefore we get if $k_1 < k_2$

$$\hat{P}_X(1) \log \frac{W(k_2|1)}{\hat{q}_Y(k_2)} \geq \lambda_{1, k_1} = \hat{P}_X(1) \log \frac{W(k_1|1)}{\hat{q}_Y(k_1)} \quad (61)$$

$$\hat{P}_X(2) \log \frac{W(k_1|2)}{\hat{q}_Y(k_1)} \geq \lambda_{2, k_2} = \hat{P}_X(2) \log \frac{W(k_2|2)}{\hat{q}_Y(k_2)}. \quad (62)$$

Therefore we get that $\left\{ \hat{P}_X(1) \log \frac{W(k|1)}{\hat{q}_Y(k)} \right\}_{k=1}^K$ and $-\left\{ \hat{P}_X(2) \log \frac{W(k|2)}{\hat{q}_Y(k)} \right\}_{k=1}^K$ are both non-decreasing sequences and so is any linear combination of them with positive coefficients. Therefore, since

$$\begin{aligned} \log W(k|1) - \log W(k|2) &= \frac{1}{\hat{P}_X(1)} \left(\hat{P}_X(1) \log \frac{W(k|1)}{\hat{q}_Y(k)} \right) \\ &\quad - \frac{1}{\hat{P}_X(2)} \left(\hat{P}_X(2) \log \frac{W(k|2)}{\hat{q}_Y(k)} \right) \end{aligned} \quad (63)$$

we conclude that the sequence $\{\log W(k|1) - \log W(k|2)\}_{k=1}^K$ is a non-decreasing sequence. \blacksquare

REFERENCES

- [1] N. Merhav, G. Kaplan, A. Lapidoth, and S. S. Shitz, "On information rates for mismatched decoders," *IEEE Trans. Inf. Theory*, vol. 40, pp. 1953–1967, Nov. 1994.
- [2] I. Csiszár and P. Narayan, "Channel capacity for a given decoding metric," *IEEE Trans. Inf. Theory*, vol. 41, pp. 35–43, Jan. 1995.
- [3] J. Y. N. Hui, "Fundamental issues of multiple accessing," Ph.D. dissertation, Massachusetts Institute of Technology, 1983.
- [4] I. Csiszár and J. Körner, "Graph decomposition: A new key to coding theorems," *IEEE Trans. Inf. Theory*, vol. 27, pp. 5–12, Jan. 1981.
- [5] J. Scarlett, "Reliable communication under mismatched decoding," Ph.D. dissertation, Ph. D. dissertation, University of Cambridge, 2014, [Online: <http://itc.upf.edu/biblio/1061>, 2014.
- [6] V. B. Balakirsky, "A converse coding theorem for mismatched decoding at the output of binary-input memoryless channels," *IEEE Trans. Inf. Theory*, vol. 41, no. 6, pp. 1889–1902, 1995.
- [7] J. Scarlett, A. Somekh-Baruch, A. Martínez, and A. Guillén i Fàbregas, "A counter-example to the mismatched decoding converse for binary-input discrete memoryless channels," *IEEE Trans. Inf. Theory*, vol. 61, pp. 5387–5395, Oct. 2015.
- [8] A. Somekh-Baruch, "Converse theorems for the DMC with mismatched decoding," *IEEE Trans. Inf. Theory*, vol. 64, pp. 6196–6207, Sept. 2018.
- [9] J. von Neumann, "Zur Theorie der Gesellschaftsspiele," *Math. Ann.*, vol. 100, pp. 295–320, 1928.
- [10] S. Boyd and L. Vandenberghe, "Convex Optimization," *Cambridge University Press*, 2004.