# An Exploration of the Role of Principal Inertia Components in Information Theory

Flavio P. Calmon, Mayank Varia, Muriel Médard

### Abstract

The principal inertia components of the joint distribution of two random variables $X$ and $Y$ are inherently connected to how an observation of $Y$ is statistically related to a hidden variable $X$. In this paper, we explore this connection within an information theoretic framework. We show that, under certain symmetry conditions, the principal inertia components play an important role in estimating one-bit functions of $X$, namely $f(X)$, given an observation of $Y$. In particular, the principal inertia components bear an interpretation as filter coefficients in the linear transformation of $p_{f(X)|X}$ into $p_{f(X)|Y}$. This interpretation naturally leads to the conjecture that the mutual information between $f(X)$ and $Y$ is maximized when all the principal inertia components have equal value. We also study the role of the principal inertia components in the Markov chain $B \to X \to Y \to \widehat{B}$, where $B$ and $\widehat{B}$ are binary random variables. We illustrate our results for the setting where $X$ and $Y$ are binary strings and $Y$ is the result of sending $X$ through an additive noise binary channel.

## I. Introduction

Let $X$ and $Y$ be two discrete random variables with finite support $\mathcal{X}$ and $\mathcal{Y}$, respectively. $X$ and $Y$ are related through a conditional distribution (channel), denoted by $p_{Y|X}$. For each $x \in \mathcal{X}$, $p_{Y|X}(\cdot|x)$ will be a vector on the $|\mathcal{Y}|$-dimensional simplex, and the position of these vectors on the simplex will determine the nature of the relationship between $X$ and $Y$. If $p_{Y|X}$ is fixed, what can be learned about $X$ given an observation of $Y$, or the degree of accuracy of what can be inferred about $X$ *a posteriori*, will then depend on the marginal distribution $p_X$. The value $p_X(x)$, in turn, ponderates the corresponding vector $p_{Y|X}(\cdot|x)$ akin to a mass. As a simple example, if $|\mathcal{X}| = |\mathcal{Y}|$ and the vectors $p_{Y|X}(\cdot|x)$ are located on distinct corners of the simplex, then $X$ can be perfectly learned from $Y$. As another example, assume that the vectors $p_{Y|X}(\cdot|x)$ can be grouped into two clusters located near opposite corners of the simplex. If the sum of the masses induced by $p_X$ for each cluster is approximately $1/2$, then one may expect to reliably infer on the order of $1$ unbiased bit of $X$ from an observation of $Y$.

The above discussion naturally leads to considering the use of techniques borrowed from classical mechanics. For a given inertial frame of reference, the mechanical properties of a collection of distributed point masses can be characterized by the moments of inertia of the system. The moments of inertia measure how the weight of the point masses is distributed around the center of mass. An analogous metric exists for the distribution of the vectors $p_{Y|X}$ and masses $p_X$ in the simplex, and it is the subject of study of a branch of applied statistics called *correspondence analysis* ([1], [2]). In correspondence analysis, the joint distribution $p_{X,Y}$ is decomposed in terms of the *principal inertia components*, which, in some sense, are analogous to the moments of inertia of a collection of point masses. In mathematical probability, the study of principal inertia components dates back to Hirschfeld [3], Gebelein [4],

Sarmanov [5] and Rényi [6], and similar analysis have also recurrently appeared in the information theory and applied probability literature. We present the formal definition of principal inertia components and a short review of the relevant literature in the next section[1].

The distribution of the vectors $p_{Y|X}$ in the simplex or, equivalently, the principal inertia components of the joint distribution of $X$ and $Y$, is inherently connected to how an observation of $Y$ is statically related to $X$. In this paper, we explore this connection within an information theoretic framework. We show that, under certain assumptions, the principal inertia components play an important part in estimating a one-bit function of $X$, namely $f(X)$ where $f : \mathcal{X} \rightarrow \{0, 1\}$, given an observation of $Y$: they can be understood as the filter coefficients in the linear transformation of $p_{f(X)|X}$ into $p_{f(X)|Y}$. Alternatively, the principal inertia components can bear an interpretation as noise, in particular when $X$ and $Y$ are binary strings. We also show that maximizing the principal inertia components is equivalent to maximizing the first-order term of the Taylor series expansion of certain convex measures of information between $f(X)$ and $Y$. We conjecture that, for symmetric distributions of $X$ and $Y$ and a given upper bound on the value of the largest principal inertia component, $I(f(X); Y)$ is maximized when all the principal inertia components have the same value as the largest principal inertia component. This is equivalent to $Y$ being the result of passing $X$ through a $q$-ary symmetric channel. This conjecture, if proven, would imply that the conjecture made by Kumar and Courtade in [7].

Finally, we study the Markov chain $B \rightarrow X \rightarrow Y \rightarrow \widehat{B}$, where $B$ and $\widehat{B}$ are binary random variables, and the role of the principal inertia components in characterizing the relation between $B$ and $\widehat{B}$. We show that that this relation is linked to solving a non-linear maximization problem, which, in turn, can be solved when $\widehat{B}$ is an unbiased estimate of $B$, the joint distribution of $X$ and $Y$ is symmetric and $\Pr\{B = \widehat{B} = 0\} \geq \mathbb{E}[B]^2$. We illustrate this result for the setting where $X$ is a binary string and $Y$ is the result of sending $X$ through a memoryless binary symmetric channel. We note that this is a similar setting to the one considered by Anantharam *et al.* in [8].

The rest of the paper is organized as follows. Section II presents the notation and definitions used in this paper, and discusses some of the related literature. Section III introduces the notion of conforming distributions and ancillary results. Section IV presents results concerning the role of the principal inertia components in inferring one-bit functions of $X$ from an observation of $Y$, as well as the linear transformation of $p_X$ into $p_Y$ in certain symmetric settings. We argue that, in such settings, the principal inertia components can be viewed as filter coefficients in a linear transformation. In particular, results for binary channels with additive noise are derived using techniques inspired by Fourier analysis of Boolean functions. Furthermore, Section IV also introduces a conjecture that encompasses the one made by Kumar and Courtade in [7]. Finally, Section V provides further evidence for this conjecture by investigating the Markov chain $B \rightarrow X \rightarrow Y \rightarrow \widehat{B}$ where $B$ and $\widehat{B}$ are binary random variables.

## II. PRINCIPAL INERTIA COMPONENTS

### A. Notation

We denote matrices by bold capitalized letters (e.g. $\mathbf{A}$) and vectors by bold lower case letters (e.g. $\mathbf{x}$). The $i$-th component of a vector $\mathbf{x}$ is denoted by $\mathbf{x}_i$. Random variables are denoted by upper-case letters (e.g. $X$ and $Y$). We define $[n] \triangleq \{1, \ldots, n\}$.

Throughout the text we assume that $X$ and $Y$ are discrete random variables with finite support sets $\mathcal{X}$ and $\mathcal{Y}$. Unless otherwise specified, we let, without loss of generality, $\mathcal{X} = [m]$ and $\mathcal{Y} = [n]$. The joint distribution matrix of $\mathbf{P}$ is an $m \times n$ matrix with $(i, j)$-th entry equal to $p_{X,Y}(i, j)$. We denote by $\mathbf{p}_X$ (respectively, $\mathbf{p}_Y$) the vector

---

[1]We encourage the readers that are unfamiliar with the topic to skip ahead and read Section II and then return to this introduction.

with $i$-th entry equal to $p_X(i)$ (resp. $p_Y(i)$). $\mathbf{D}_X = \text{diag}(\mathbf{p}_X)$ and $\mathbf{D}_Y = \text{diag}(\mathbf{p}_Y)$ are matrices with diagonal entries equal to $\mathbf{p}_X$ and $\mathbf{p}_Y$, respectively, and all other entries equal to 0. The matrix $\mathbf{P}_{Y|X} \in \mathbb{R}^{m \times n}$ denotes the matrix with $(i,j)$-th entry equal to $p_{Y|X}(j|i)$. Note that $\mathbf{P} = \mathbf{D}_X \mathbf{P}_{Y|X}$.

For a given joint distribution matrix $\mathbf{P}$, the set of all vectors contained in the unit cube in $\mathbb{R}^n$ that satisfy $\|\mathbf{P}\mathbf{x}\|_1 = a$ is given by

$$\mathcal{C}^n(a, \mathbf{P}) \triangleq \{\mathbf{x} \in \mathbb{R}^n | 0 \leq \mathbf{x}_i \leq 1, \|\mathbf{P}\mathbf{x}\|_1 = a\}. \tag{1}$$

The set of all $m \times n$ probability distribution matrices is given by $\mathcal{P}_{m,n}$.

For $x^n \in \{-1, 1\}^n$ and $\mathcal{S} \subseteq [n]$, $\chi_{\mathcal{S}}(x^n) \triangleq \prod_{i \in \mathcal{S}} x_i$ (we consider $\chi_{\varnothing}(x) = 1$). For $y^n \in \{-1, 1\}^n$, $a^n = x^n \oplus y^n$ is the vector resulting from the entrywise product of $x^n$ and $y^n$, i.e. $a_i = x_i y_i$, $i \in [n]$.

Given two probability distributions $p_X$ and $q_X$ and $f(t)$ a smooth convex function defined for $t > 0$ with $f(1) = 0$, the $f$-divergence is defined as [9]

$$D_f(p_X \| q_X) \triangleq \sum_x q_X(x) f\left(\frac{p_X(x)}{q_X(x)}\right). \tag{2}$$

The $f$-information is given by

$$I_f(X; Y) \triangleq D_f(p_{X,Y} \| p_X p_Y). \tag{3}$$

When $f(x) = x \log(x)$, then $I_f(X; Y) = I(X; Y)$. A study of information metrics related to $f$-information was given in [10] in the context of channel coding converses.

## B. Principal Inertia Components and Decomposing the Joint Distribution Matrix

We briefly define in this section the *principal inertia decomposition* of the joint distribution matrix $\mathbf{P}$. The term "principal inertia" is borrowed from the correspondence analysis literature [1]. The study of the principal inertia components of the joint distribution of two random variables dates back to Hirshfield [3], Gebelein [4], Sarmanov [5] and Rényi [6], having appeared in the work of Witsenhausen [11], Ahlswede and Gács [12] and, more recently, Anantharam *et al.* [13], Polyanskiy [14] and Calmon *et al.* [15], among others. For an overview, we refer the reader to [13], [15].

**Definition 1.** We call the singular value decomposition $\mathbf{D}_X^{-1/2} \mathbf{P} \mathbf{D}_Y^{-1/2} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ the *principal inertia decomposition* of $X$ and $Y$, where $\mathbf{\Sigma}$ is a diagonal matrix with $\text{diag}(\mathbf{\Sigma}) = (1, \sigma_1, \ldots, \sigma_d)$ and $d = \min(m, n) - 1$. The values $\sigma_i^2$, $i = 1, \ldots, d$, are called the *principal inertia components* of $X$ and $Y$. In particular $\rho_m(X; Y) = \sigma_1$, where $\rho_m(X; Y)$ denotes the maximal correlation coefficient of $X$ and $Y$. The maximal correlation coefficient, in turn, is given by

$$\rho_m(X; Y) \triangleq \sup\left\{\mathbb{E}\left[f(X) g(Y)\right] | \mathbb{E}\left[f(X)\right] = \mathbb{E}\left[g(Y)\right] = 0, \mathbb{E}\left[f(X)^2\right] = \mathbb{E}\left[g(X)^2\right] = 1\right\}.$$

The values $\sigma_1, \ldots, \sigma_d$ in the previous definition are the spectrum of the conditional expectation operator $(Tf)(x) \triangleq \mathbb{E}\left[f(Y)|X = x\right]$, where $f : \mathcal{Y} \to \mathbb{R}$ [6]. Indeed, the spectrum of $T$ and the principal inertia components are entirely equivalent when $X$ and $Y$ have finite support sets. Nevertheless, the reader should note that the analysis based on the conditional expectation operator lends itself to more general settings, including random variables with continuous support. We do not pursue this matter further here, since our focus is on discrete random variables with finite support.

The principal inertia components satisfy the data processing inequality (see, for example, [14], [15], [16]): if $X \to Y \to Z$ and $\sigma_i$ are the principal inertia components of $X$ and $Y$ and $\widetilde{\sigma}_i$ are the principal inertia components

of $X$ and $Z$, then $\sum_{i=1}^{k} \widetilde{\sigma}_i^2 \leq \sum_{i=1}^{k} \sigma_i^2$ for all $k$. Furthermore, for a fixed marginal distribution $p_X$, $\sum_{i=1}^{k} \sigma_i^2$ is convex in $p_{Y|X}$. Note the joint distribution matrix $\mathbf{P}$ as can be written as

$$\mathbf{P} = \mathbf{D}_X^{1/2} \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{D}_Y^{1/2}. \tag{4}$$

## III. CONFORMING DISTRIBUTIONS

In this paper we shall recurrently use probability distribution matrices that are symmetric and positive-semidefinite. This motivates the following definition.

**Definition 2.** A joint distribution $p_{X,Y}$ is said to be *conforming* if the corresponding matrix $\mathbf{P}$ satisfies $\mathbf{P} = \mathbf{P}^T$ and $\mathbf{P}$ is positive-semidefinite.

**Remark 1.** If $X$ and $Y$ have a conforming joint distribution, then they have the same marginal distribution. Consequently, $\mathbf{D} \triangleq \mathbf{D}_X = \mathbf{D}_Y$, and $\mathbf{P} = \mathbf{D}^{1/2} \mathbf{U} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{D}^{1/2}$.

Symmetric channels[2] are closely related to conforming probability distributions. We shall illustrate this relation in the next lemma and in Section IV.

**Lemma 1.** *If $\mathbf{P}$ is conforming, then the corresponding conditional distribution matrix $\mathbf{P}_{Y|X}$ is positive semidefinite. Furthermore, for any symmetric channel $\mathbf{P}_{Y|X} = \mathbf{P}_{Y|X}^T$, there is an input distribution $\mathbf{p}_X$ (namely, the uniform distribution) such that the principal inertia components of $\mathbf{P} = \mathbf{D}_X \mathbf{P}_{Y|X}$ correspond to the square of the eigenvalues of $\mathbf{P}_{Y|X}$. In this case, if $\mathbf{P}_{Y|X}$ is also positive-semidefinite, then $\mathbf{P}$ is conforming.*

*Proof:* Let $\mathbf{P}$ be conforming and $\mathcal{X} = \mathcal{Y} = [m]$. Then $\mathbf{P}_{Y|X} = \mathbf{D}^{-1/2} \mathbf{U} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{D}^{1/2} = \mathbf{Q} \boldsymbol{\Sigma} \mathbf{Q}^{-1}$, where $\mathbf{Q} = \mathbf{D}^{-1/2} \mathbf{U}$. It follows that $\mathrm{diag}(\boldsymbol{\Sigma})$ are the eigenvalues of $\mathbf{P}_{Y|X}$, and, consequently, $\mathbf{P}_{Y|X}$ is positive semidefinite.

Now let $\mathbf{P}_{Y|X} = \mathbf{P}_{Y|X}^T = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T$. The entries of $\boldsymbol{\Lambda}$ here are the eigenvalues of $\mathbf{P}_{Y|X}$ and not necessarily positive. Since $\mathbf{P}_{Y|X}$ is symmetric, it is also doubly stochastic, and for $X$ uniformly distributed $Y$ is also uniformly distributed. Therefore, $\mathbf{P}$ is symmetric, and $\mathbf{P} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T / m$. It follows directly that the principal inertia components of $\mathbf{P}$ are exactly the diagonal entries of $\boldsymbol{\Lambda}^2$, and if $P_{Y|X}$ is positive-semidefinite then $\mathbf{P}$ is conforming. ∎

The $q$-ary symmetric channel, defined below, is of particular interest to some of the results derived in the following sections.

**Definition 3.** The $q$-ary symmetric channel with crossover probability $\epsilon \leq 1 - q^{-1}$, also denoted as $(\epsilon, q)$-SC, is defined as the channel with input $X$ and output $Y$ where $\mathcal{X} = \mathcal{Y} = [q]$ and

$$p_{Y|X}(y|x) = \begin{cases} 1 - \epsilon & \text{if } x = y \\ \frac{\epsilon}{q-1} & \text{if } x \neq y. \end{cases}$$

Let $X$ and $Y$ have a conforming joint distribution matrix with $\mathcal{X} = \mathcal{Y} = [q]$ and principal inertia components $\sigma_1^2, \ldots, \sigma_d^2$. The following lemma shows that conforming $\mathbf{P}$ can be transformed into the joint distribution of a $q$-ary symmetric channel with input distribution $p_X$ by setting $\sigma_1^2 = \sigma_2^2 = \cdots = \sigma_d^2$, i.e. making all principal inertia components equal to the largest one.

**Lemma 2.** *Let $\mathbf{P}$ be a conforming joint distribution matrix of $X$ and $Y$, with $X$ and $Y$ uniformly distributed, $\mathcal{X} = \mathcal{Y} = [q]$, $\mathbf{P} = q^{-1} \mathbf{U} \boldsymbol{\Sigma} \mathbf{U}^T$ and $\boldsymbol{\Sigma} = \mathrm{diag}(1, \sigma_1, \ldots, \sigma_d)$. For $\widetilde{\boldsymbol{\Sigma}} = \mathrm{diag}(1, \sigma_1, \ldots, \sigma_1)$, let $X$ and $\widetilde{Y}$ have*

---

[2] We say that a channel is symmetric if $\mathcal{X} = \mathcal{Y} = [m]$ and $p_{Y|X}(i|j) = p_{Y|X}(j|i) \; \forall i, j \in [m]$.

*joint distribution* $\widetilde{\mathbf{P}} = q^{-1}\mathbf{U}\widetilde{\boldsymbol{\Sigma}}\mathbf{U}^T$. *Then,* $\widetilde{Y}$ *is the result of passing* $X$ *through a* $(\epsilon, q)$-*SC, with*

$$\epsilon = \frac{(q-1)(1-\rho_m(X;Y))}{q}. \tag{5}$$

*Proof:* The first column of $\mathbf{U}$ is $\mathbf{p}_X^{1/2}$ and, since $X$ is uniformly distributed, $\mathbf{p}_X^{1/2} = q^{-1/2}\mathbf{1}$. Therefore

$$\begin{aligned}\widetilde{\mathbf{P}} &= q^{-1}\mathbf{U}\widetilde{\boldsymbol{\Sigma}}\mathbf{U}^T \\ &= q^{-1}\sigma_1\mathbf{I} + q^{-2}(1-\sigma_1)\mathbf{11}^T.\end{aligned} \tag{6}$$

Consequently, $\widetilde{\mathbf{P}}$ has diagonal entries equal to $(1 + (q-1)\sigma_1)/q^2$ and all other entries equal to $(1-\sigma_1)/q^2$. The result follows by noting that $\sigma_1 = \rho_m(X;Y)$. ∎

**Remark 2.** For $X$, $Y$ and $\widetilde{Y}$ given in the previous lemma, a natural question that arises is whether $Y$ is a degraded version of $\widetilde{Y}$, i.e. $X \to \widetilde{Y} \to Y$. Unfortunately, this is **not true** in general, since the matrix $\mathbf{U}\widetilde{\boldsymbol{\Sigma}}^{-1}\boldsymbol{\Sigma}\mathbf{U}^T$ does not necessarily contain only positive entries, although it is doubly-stochastic. However, since the principal inertia components of $X$ and $\widetilde{Y}$ upper bound the principal inertia components of $X$ and $Y$, it is natural to expect that, at least in some sense, $\widetilde{Y}$ is more informative about $X$ than $Y$. This intuition is indeed correct for certain estimation problems where a one-bit function of $X$ is to be inferred from a single observation $Y$ or $\widetilde{Y}$, and will be investigated in the next section.

## IV. ONE-BIT FUNCTIONS AND CHANNEL TRANSFORMATIONS

Let $B \to X \to Y$, where $B$ is a binary random variable. When $X$ and $Y$ have a conforming probability distribution, the principal inertia components of $X$ and $Y$ have a particularly interesting interpretation: they can be understood as the filter coefficients in the linear transformation of $p_{B|X}$ into $p_{B|Y}$. In order to see why this is the case, consider the joint distribution of $B$ and $Y$, denoted here by $\mathbf{Q}$, given by

$$\mathbf{Q} = [\mathbf{f} \quad 1-\mathbf{f}]^T\mathbf{P} = [\mathbf{f} \quad 1-\mathbf{f}]^T\mathbf{P}_{X|Y}\mathbf{D}_Y = [\mathbf{g} \quad 1-\mathbf{g}]^T\mathbf{D}_Y, \tag{7}$$

where $\mathbf{f} \in \mathbb{R}^m$ and $\mathbf{g} \in \mathbb{R}^n$ are column-vectors with $\mathbf{f}_i = p_{B|X}(0|i)$ and $\mathbf{g}_j = p_{B|Y}(0|j)$. In particular, if $B$ is a deterministic function of $X$, $\mathbf{f} \in \{0,1\}^m$.

If $\mathbf{P}$ is conforming and $\mathcal{X} = \mathcal{Y} = [m]$, then $\mathbf{P} = \mathbf{D}^{1/2}\mathbf{U}\boldsymbol{\Sigma}\mathbf{U}^T\mathbf{D}^{1/2}$, where $\mathbf{D} = \mathbf{D}_X = \mathbf{D}_Y$. Assuming $\mathbf{D}$ fixed, the joint distribution $\mathbf{Q}$ is entirely specified by the linear transformation of $\mathbf{f}$ into $\mathbf{g}$. Denoting $\mathbf{T} \triangleq \mathbf{U}^T\mathbf{D}^{1/2}$, this transformation is done in three steps:

1) (Linear transform) $\widehat{\mathbf{f}} \triangleq \mathbf{Tf}$,
2) (Filter) $\widehat{\mathbf{g}} \triangleq \boldsymbol{\Sigma}\widehat{\mathbf{f}}$, where the diagonal of $\boldsymbol{\Sigma}^2$ are the principal inertia components of $X$ and $Y$,
3) (Inverse transform) $\mathbf{g} = \mathbf{T}^{-1}\widehat{\mathbf{g}}$.

Note that $\widehat{\mathbf{f}}_1 = \widehat{\mathbf{g}}_1 = 1 - \mathbb{E}[B]$ and $\widehat{\mathbf{g}} = \mathbf{Tg}$. Consequently, the principal inertia coefficients of $X$ and $Y$ bear an interpretation as the filter coefficients in the linear transformation of $p_{B|X}(0|\cdot)$ into $p_{B|Y}(0|\cdot)$.

A similar interpretation can be made for symmetric channels, where $\mathbf{P}_{Y|X} = \mathbf{P}_{Y|X}^T = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T$ and $\mathbf{P}_{Y|X}$ acts as the matrix of the linear transformation of $\mathbf{p}_X$ into $\mathbf{p}_Y$. Note that $\mathbf{p}_Y = \mathbf{P}_{Y|X}\mathbf{p}_X$, and, consequently, $\mathbf{p}_X$ is transformed into $\mathbf{p}_Y$ in the same three steps as before:

1) (Linear transform) $\widehat{\mathbf{p}_X} = \mathbf{U}^T\mathbf{p}_X$,
2) (Filter) $\widehat{\mathbf{p}_Y} = \boldsymbol{\Lambda}\widehat{\mathbf{p}_X}$, where the diagonal of $\boldsymbol{\Lambda}^2$ are the principal inertia components of $X$ and $Y$ in the particular case when $X$ is uniformly distributed (Lemma 1),
3) (Inverse transform) $\mathbf{p}_Y = \mathbf{U}\widehat{\mathbf{p}_Y}$.

From this perspective, the vector $\mathbf{z} = \mathbf{U\Lambda 1}m^{-1/2}$ can be understood as a proxy for the "noise effect" of the channel. Note that $\sum_i \mathbf{z}_i = 1$. However, the entries of $\mathbf{z}$ are not necessarily positive, and $\mathbf{z}$ might not be a *de facto* probability distribution.

We now illustrate these ideas by investigating binary channels with additive noise in the next section, where $\mathbf{T}$ will correspond to the well-known Walsh-Hadamard transform matrix.

### A. Example: Binary Additive Noise Channels

In this example, let $\mathcal{X}^n, \mathcal{Y}^n \subseteq \{-1, 1\}^n$ be the support sets of $X^n$ and $Y^n$, respectively. We define two sets of channels that transform $X^n$ into $Y^n$. In each set definition, we assume the conditions for $p_{Y^n|X^n}$ to be a valid probability distribution (i.e. non-negativity and unit sum).

**Definition 4.** The set of *parity-changing channels* of block-length $n$, denoted by $\mathcal{A}_n$, is defined as:

$$\mathcal{A}_n \triangleq \left\{ p_{Y^n|X^n} \mid \forall \mathcal{S} \subseteq [n], \ \exists c_{\mathcal{S}} \in [-1, 1] \text{ s.t. } \mathbb{E}\left[\chi_{\mathcal{S}}(Y^n)|X^n\right] = c_{\mathcal{S}}\chi_{\mathcal{S}}(X^n) \right\}. \tag{8}$$

The set of all *binary additive noise channels* is given by

$$\mathcal{B}_n \triangleq \left\{ p_{Y^n|X^n} \mid \exists Z^n \text{ s.t. } Y^n = X^n \oplus Z^n, \ \operatorname{supp}(Z^n) \subseteq \{-1, 1\}^n, Z^n \perp\!\!\!\perp X^n \right\}.$$

The definition of parity-changing channels is inspired by results from the literature on Fourier analysis of Boolean functions. For an overview of the topic, we refer the reader to the survey [17]. The set of binary additive noise channels, in turn, is widely used in the information theory literature. The following theorem shows that both characterizations are equivalent.

**Theorem 1.** $\mathcal{A}_n = \mathcal{B}_n$.

*Proof:* Let $Y^n = X^n \oplus Z^n$ for some $Z^n$ distributed over $\{-1, 1\}^n$ and independent of $X^n$. Thus

$$\begin{aligned}
\mathbb{E}\left[\chi_{\mathcal{S}}(Y^n)|X^n\right] &= \mathbb{E}\left[\chi_{\mathcal{S}}(Z^n \oplus X^n) \mid X^n\right] \\
&= \mathbb{E}\left[\chi_{\mathcal{S}}(X^n)\chi_{\mathcal{S}}(Z^n) \mid X^n\right] \\
&= \chi_{\mathcal{S}}(X^n)\mathbb{E}\left[\chi_{\mathcal{S}}(Z^n)\right],
\end{aligned}$$

where the last equality follows from the assumption that $X^n \perp\!\!\!\perp Z^n$. By letting $c_{\mathcal{S}} = \mathbb{E}\left[\chi_{\mathcal{S}}(Z^n)\right]$, it follows that $p_{Y^n|X^n} \in \mathcal{A}_n$ and, consequently, $\mathcal{B}_n \subseteq \mathcal{A}_n$.

Now let $y_n$ be fixed and $\delta_{y^n} : \{-1, 1\}^n \to \{0, 1\}$ be given by

$$\delta_{y^n}(x^n) = \begin{cases} 1, & x^n = y^n, \\ 0, & \text{otherwise.} \end{cases}$$

Since the function $\delta_{y^n}$ has Boolean inputs, it can be expressed in terms of its Fourier expansion [17, Prop. 1.1] as

$$\delta_{y^n}(x^n) = \sum_{\mathcal{S} \subseteq [n]} \widehat{d}_{\mathcal{S}}\chi_{\mathcal{S}}(x^n).$$

Now let $p_{Y^n|X^n} \in \mathcal{A}_n$. Observe that $p_{Y^n|X^n}(y^n|x^n) = \mathbb{E}\left[\delta_{y^n}(Y^n) \mid X^n = x^n\right]$ and, for $z^n \in \{-1,1\}^n$,

$$
\begin{aligned}
p_{Y^n|X^n}(y^n \oplus z^n | x^n \oplus z^n) &= \mathbb{E}\left[\delta_{y^n \oplus z^n}(Y^n) \mid X^n = x^n \oplus z^n\right] \\
&= \mathbb{E}\left[\delta_{y^n}(Y^n \oplus z^n) \mid X^n = x^n \oplus z^n\right] \\
&= \mathbb{E}\left[\sum_{\mathcal{S} \subseteq [n]} \widehat{d}_{\mathcal{S}} \chi_{\mathcal{S}}(Y^n \oplus z^n) \mid X^n = x^n \oplus z^n\right] \\
&= \mathbb{E}\left[\sum_{\mathcal{S} \subseteq [n]} \widehat{d}_{\mathcal{S}} \chi_{\mathcal{S}}(Y^n)\chi_{\mathcal{S}}(z^n) \mid X^n = x^n \oplus z^n\right] \\
&\overset{(a)}{=} \sum_{\mathcal{S} \subseteq [n]} c_{\mathcal{S}} \widehat{d}_{\mathcal{S}} \chi_{\mathcal{S}}(x^n \oplus z^n)\chi_{\mathcal{S}}(z^n) \\
&= \sum_{\mathcal{S} \subseteq [n]} c_{\mathcal{S}} \widehat{d}_{\mathcal{S}} \chi_{\mathcal{S}}(x^n) \\
&\overset{(b)}{=} \mathbb{E}\left[\sum_{\mathcal{S} \subseteq [n]} \widehat{d}_{\mathcal{S}} \chi_{\mathcal{S}}(Y^n) | X^n = x^n\right] \\
&= \mathbb{E}\left[\delta_{y^n}(Y^n) \mid X^n = x^n\right] \\
&= p_{Y^n|X^n}(y^n|x^n).
\end{aligned}
$$

Equalities $(a)$ and $(b)$ follow from the definition of $\mathcal{A}_n$. By defining the distribution of $Z^n$ as $p_{Z^n}(z^n) \triangleq p_{Y^n|X^n}(z^n|\mathbf{1}^n)$, where $\mathbf{1}^n$ is the vector with all entries equal to 1, it follows that $Z^n = X^n \oplus Y^n$, $Z^n \perp\!\!\!\perp X^n$ and $p_{Y^n|X^n} \subseteq \mathcal{B}_n$.

■

The previous theorem suggests that there is a correspondence between the coefficients $c_{\mathcal{S}}$ in (8) and the distribution of the additive noise $Z^n$ in the definition of $\mathcal{B}_n$. The next result shows that this is indeed the case and, when $X^n$ is uniformly distributed, the coefficients $c_{\mathcal{S}}^2$ correspond to the principal inertia components between $X^n$ and $Y^n$.

**Theorem 2.** Let $p_{Y^n|X^n} \in \mathcal{B}_n$, and $X^n \sim p_{X^n}$. Then $\mathbf{P}_{X^n,Y^n} = \mathbf{D}_{X^n}\mathbf{H}_{2^n}\mathbf{\Lambda}\mathbf{H}_{2^n}$, where $\mathbf{H}_l$ is the $l \times l$ normalized Hadamard matrix (i.e. $\mathbf{H}_l^2 = \mathbf{I}$). Furthermore, for $Z^n \sim p_{Z^n}$, $\mathrm{diag}\,(\mathbf{\Lambda}) = 2^{n/2}\mathbf{H}_{2^n}\mathbf{p}_{Z^n}$, and the diagonal entries of $\mathbf{\Lambda}$ are equal to $c_{\mathcal{S}}$ in (8). Finally, if $X$ is uniformly distributed, then $c_{\mathcal{S}}^2$ are the principal inertia components of $X^n$ and $Y^n$.

*Proof:* Let $p_{Y^n|X^n} \in \mathcal{A}_n$ be given. From Theorem 1 and the definition of $\mathcal{A}_n$, it follows that $\chi_{\mathcal{S}}(Y^n)$ is a right eigenvector of $p_{Y^n|X^n}$ with corresponding eigenvalue $c_{\mathcal{S}}$. Since $\chi_{\mathcal{S}}(Y^n)2^{-n/2}$ corresponds to a row of $\mathbf{H}_{2^n}$ for each $\mathcal{S}$ (due to the Kronecker product construction of the Hadamard matrix) and $\mathbf{H}_{2^n}^2 = \mathbf{I}$, then $\mathbf{P}_{X^n,Y^n} = \mathbf{D}_{X^n}\mathbf{H}_{2^n}\mathbf{\Lambda}\mathbf{H}_{2^n}$. Finally, note that $\mathbf{p}_Z^T = 2^{-n/2}\mathbf{1}^T\mathbf{\Lambda}\mathbf{H}_{2^n}$. From Lemma 1, it follows that $c_{\mathcal{S}}^2$ are the principal inertia components of $X^n$ and $Y^n$ if $X^n$ is uniformly distributed. ■

**Remark 3.** Theorem 2 indicates that one possible method for estimating the distribution of the additive binary noise $Z^n$ is to estimate its effect on the parity bits of $X^n$ and $Y^n$. In this case, we are estimating the coefficients $c_{\mathcal{S}}$ of the Walsh-Hadamard transform of $p_{Z^n}$. This approach was studied by Raginsky *et al.* in [18].

Theorem 2 illustrates the filtering role of the principal inertia components, discussed in the beginning of this section. If $X^n$ is uniform, and using the same notation as in (7), then the vector of conditional probabilities $\mathbf{f}$ is transformed into the vector of *a posteriori* probabilities $\mathbf{g}$ by: (i) taking the Hadamard transform of $\mathbf{f}$, (ii) filtering

the transformed vector according to the coefficients $c_{\mathcal{S}}$, where $\mathcal{S} \in [n]$, and (iii) taking the inverse Hadamard transform. The same rationale applies to the transformation of $\mathbf{p}_X$ into $\mathbf{p}_Y$ in binary additive channels.

### B. Quantifying the Information of a Boolean Function of the Input of a Noisy Channel

We now investigate the connection between the principal inertia components and $f$-information in the context of one-bit functions of $X$. Recall from the discussion in the beginning of this section and, in particular, equation (7), that for a binary $B$ and $B \to X \to Y$, the distribution of $B$ and $Y$ is entirely specified by the transformation of $\mathbf{f}$ into $\mathbf{g}$, where $\mathbf{f}$ and $\mathbf{g}$ are vectors with entries equal to $p_{B|X}(0|\cdot)$ and $p_{B|Y}(0|\cdot)$, respectively.

For $\mathbb{E}[B] = 1 - a$, the $f$-information between $B$ and $Y$ is given by[3]

$$I_f(B;Y) = \mathbb{E}\left[af\left(\frac{\mathbf{g}_Y}{a}\right) + (1-a)f\left(\frac{1-\mathbf{g}_Y}{1-a}\right)\right].$$

For $0 \leq r, s \leq 1$, we can expand $f\left(\frac{r}{s}\right)$ around 1 as

$$f\left(\frac{r}{s}\right) = \sum_{k=1}^{\infty} \frac{f^{(k)}(1)}{k!}\left(\frac{r-s}{r}\right)^k.$$

Denoting

$$c_k(\alpha) \triangleq \frac{1}{a^{k-1}} + \frac{(-1)^k}{(1-a)^{k-1}},$$

the $f$-information can then be expressed as

$$I_f(B;Y) = \sum_{k=2}^{\infty} \frac{f^{(k)}(1)c_k(a)}{k!}\mathbb{E}\left[(\mathbf{g}_Y - a)^k\right]. \tag{9}$$

Similarly to [9, Chapter 4], for a fixed $\mathbb{E}[B] = 1 - a$, maximizing the principal inertia components between $X$ and $Y$ will always maximize the first term in the expansion (9). To see why this is the case, observe that

$$\begin{aligned}
\mathbb{E}\left[(\mathbf{g}_Y - a)^k\right] &= (\mathbf{g} - a)^T \mathbf{D}_Y (\mathbf{g} - a) \\
&= \mathbf{g}^T \mathbf{D}_Y \mathbf{g} - a^2 \\
&= \mathbf{f}^T \mathbf{D}_X^{1/2} \mathbf{U} \boldsymbol{\Sigma}^2 \mathbf{U}^T \mathbf{D}_x^{1/2} \mathbf{f} - a^2.
\end{aligned} \tag{10}$$

For a fixed $a$ and any $\mathbf{f}$ such that $\mathbf{f}^T \mathbf{1} = a$, (10) is non-decreasing in the diagonal entries of $\boldsymbol{\Sigma}^2$ which, in turn, are exactly the principal inertia components of $X$ and $Y$. Equivalently, (10) is non-decreasing in the $\chi^2$-divergence between $p_{X,Y}$ and $p_X p_Y$.

However, we do note that increasing the principal inertia components **does not** increase the $f$-information between $B$ and $Y$ in general. Indeed, for a fixed $\mathbf{U}$, $\mathbf{V}$ and marginal distributions of $X$ and $Y$, increasing the principal inertia components might not even lead to a valid probability distribution matrix $\mathbf{P}$.

Nevertheless, if $\mathbf{P}$ is conforming and $X$ and $Y$ are uniformly distributed over $[q]$, as shown in Lemma (2), by increasing the principal inertia components we can define a new random variable $\widetilde{Y}$ that results from sending $X$ through a $(\epsilon, q)$-SC, where $\epsilon$ is given in (5). In this case, the $f$-information between $B$ and $Y$ has a simple expression when $B$ is a function of $X$.

---

[3]Note that here we assume that $\mathcal{Y} = [n]$, so there is no ambiguity in indexing $p_{B|Y}(0|Y)$ by $\mathbf{g}_Y$.

**Lemma 3.** *Let $B \to X \to \widetilde{Y}$, where $B = h(X)$ for some $h : [q] \to \{0, 1\}$, $\mathbb{E}[B] = 1 - a$ where $aq$ is an integer, $X$ is uniformly distributed in $[q]$ and $\widetilde{Y}$ is the result of passing $X$ through a $(\epsilon, q)$-SC with $\epsilon \leq (q-1)/q$. Then*

$$I_f(B; \widetilde{Y}) = a^2 f(1 + \sigma_1 c) + 2a(1-a)f(1 - \sigma_1) + (1-a)^2 f(1 + \sigma_1 c^{-1}) \tag{11}$$

*where $\sigma_1 = \rho_m(X; \widetilde{Y}) = 1 - \epsilon q (q-1)^{-1}$ and $c \triangleq (1-a)a^{-1}$. In particular, for $f(x) = x \log x$, then $I_f(X; \widetilde{Y}) = I(X; \widetilde{Y})$, and for $\sigma_1 = 1 - 2\delta$*

$$I(B; \widetilde{Y}) = h_b(a) - \alpha H_b(2\delta(1-a)) - (1-a)H_b(2\delta a) \tag{12}$$

$$\leq 1 - H_b(\delta), \tag{13}$$

*where $H_b(x) \triangleq -x \log(x) - (1-x) \log(1-x)$ is the binary entropy function.*

*Proof:* Since $B$ is a deterministic function of $X$ and $aq$ is an integer, $\mathbf{f}$ is a vector with $aq$ entries equal to 1 and $(1-a)q$ entries equal to 0. It follows from (6) that

$$\begin{aligned} I_f(B; \widetilde{Y}) &= \frac{1}{q} \sum_{i=1}^q af\left(\frac{(1-\sigma_1)a + \mathbf{f}_i \sigma_1}{a}\right) + (1-a)f\left(\frac{1 - (1-\sigma_1)a - \mathbf{f}_i \sigma_i}{1-a}\right) \\ &= a^2 f\left(1 + \sigma_1 \frac{1-a}{a}\right) + 2a(1-a)f(1-\sigma_1) + (1-a)^2 f\left(1 + \sigma_1 \frac{a}{1-a}\right). \end{aligned}$$

Letting $f(x) = x \log x$, (12) follows immediately. Since (12) is concave in $a$ and symmetric around $a = 1/2$, it is maximized at $a = 1/2$, resulting in (13). ∎

### C. On the "Most Informative Bit" Conjecture

We now return to channels with additive binary noise, analyzed is Section IV-A. Let $X^n$ be a uniformly distributed binary string of length $n$ ($\mathcal{X} = \{-1, 1\}$), and $Y^n$ the result of passing $X^n$ through a memoryless binary symmetric channel with crossover probability $\delta \leq 1/2$. Kumar and Courtade conjectured [7] that for all binary $B$ and $B \to X^n \to Y^n$ we have

$$I(B; Y^n) \leq 1 - H_b(\delta). \quad \text{(conjecture)} \tag{14}$$

It is sufficient to consider $B$ a function of $X^n$, denoted by $B = h(X^n)$, $h : \{-1, 1\}^n \to \{0, 1\}$, and we make this assumption henceforth.

From the discussion in Section IV-A, for the memoryless binary symmetric channel $Y^n = X^n \oplus Z^n$, where $Z^n$ is an i.i.d. string with $\Pr\{Z_i = 1\} = 1 - \delta$, and any $\mathcal{S} \in [n]$,

$$\begin{aligned} \mathbb{E}[\chi_{\mathcal{S}}(Y^n) | X^n] &= \chi_{\mathcal{S}}(X^n)\left(\Pr\{\chi_{\mathcal{S}}(Z^n) = 1\} - \Pr\{\chi_{\mathcal{S}}(Z^n) = -1\}\right) \\ &= \chi_{\mathcal{S}}(X^n)\left(2\Pr\{\chi_{\mathcal{S}}(Z^n) = 1\} - 1\right) \\ &= \chi_{\mathcal{S}}(X^n)(1 - 2\delta)^{|\mathcal{S}|}. \end{aligned}$$

It follows directly that $c_{\mathcal{S}} = (1 - 2\delta)^{|\mathcal{S}|}$ for all $\mathcal{S} \subseteq [n]$. Consequently, from Theorem 2, the principal inertia components of $X^n$ and $Y^n$ are of the form $(1 - 2\delta)^{2|\mathcal{S}|}$ for some $\mathcal{S} \subseteq [n]$. Observe that the principal inertia components act as a low pass filter on the vector of conditional probabilities $\mathbf{f}$ given in (7).

Can the noise distribution be modified so that the principal inertia components act as an all-pass filter? More specifically, what happens when $\widetilde{Y}^n = X^n \oplus W^n$, where $W^n$ is such that the principal inertia components between $X^n$ and $\widetilde{Y}^n$ satisfy $\sigma_i = 1 - 2\delta$? Then, from Lemma 2, $\widetilde{Y}^n$ is the result of sending $X^n$ through a $(\epsilon, 2^n)$-SC with

$\epsilon = 2\delta(1 - 2^{-n})$. Therefore, from (13),

$$I(B; \widetilde{Y}^n) \leq 1 - H_b(\delta).$$

For any function $h : \{-1, 1\}^n \to \{0, 1\}$ such that $B = h(X^n)$, from standard results in Fourier analysis of Boolean functions [17, Prop. 1.1], $h(X^n)$ can be expanded as

$$h(X^n) = \sum_{\mathcal{S} \subseteq [n]} \hat{h}_{\mathcal{S}} \chi_{\mathcal{S}}(X^n).$$

The value of $B$ is uniquely determined by the action of $h$ on $\chi_{\mathcal{S}}(X^n)$. Consequently, for a fixed function $h$, one could expect that $\widetilde{Y}^n$ should be more informative about $B$ than $Y^n$, since the parity bits $\chi_{\mathcal{S}}(X^n)$ are more reliably estimated from $\widetilde{Y}^n$ than from $Y^n$. Indeed, the memoryless binary symmetric channel attenuates $\chi_{\mathcal{S}}(X^n)$ exponentially in $|\mathcal{S}|$, acting (as argued previously) as a low-pass filter. In addition, if one could prove that for any fixed $h$ the inequality $I(B; Y^n) \leq I(B; \widetilde{Y}^n)$ holds, then (14) would be proven true. This motivates the following conjecture.

**Conjecture 1.** *For all $h : \{-1, 1\}^n \to \{0, 1\}$ and $B = h(X^n)$*

$$I(B; Y^n) \leq I(B; \widetilde{Y}^n).$$

We note that Conjecture 1 is not true in general if $B$ is not a deterministic function of $X^n$. In the next section, we provide further evidence for this conjecture by investigating information metrics between $B$ and an estimate $\widehat{B}$ derived from $Y^n$.

## V. ONE-BIT ESTIMATORS

Let $B \to X \to Y \to \widehat{B}$, where $B$ and $\widehat{B}$ are binary random variables with $\mathbb{E}[B] = 1 - a$ and $\mathbb{E}[\widehat{B}] = 1 - b$. We denote by $\mathbf{x} \in \mathbb{R}^m$ and $\mathbf{y} \in \mathbb{R}^n$ the column vectors with entries $\mathbf{x}_i = p_{B|X}(0|i)$ and $\mathbf{y}_j = p_{\widehat{B}|Y}(0|j)$. The joint distribution matrix of $B$ and $\widehat{B}$ is given by

$$\mathbf{P}_{B,\widehat{B}} = \begin{pmatrix} z & a - z \\ b - z & 1 - a - b + z \end{pmatrix}, \tag{15}$$

where $z = \mathbf{x}^T \mathbf{P} \mathbf{y} = \Pr\{B = \widehat{B} = 0\}$. For fixed values of $a$ and $b$, the joint distribution of $B$ and $\widehat{B}$ only depends on $z$.

Let $f : \mathcal{P}_{2 \times 2} \to \mathbb{R}$, and, with a slight abuse of notation, we also denote $f$ as a function of the entries of the $2 \times 2$ matrix as $f(a, b, z)$. If $f$ is convex in $z$ for a fixed $a$ and $b$, then $f$ is maximized at one of the extreme values of $z$. Examples of such functions $f$ include mutual information and expected error probability. Therefore, characterizing the maximum and minimum values of $z$ is equivalent to characterizing the maximum value of $f$ over all possible mappings $X \to B$ and $Y \to \widehat{B}$. This leads to the following definition.

**Definition 5.** For a fixed $\mathbf{P}$, the minimum and maximum values of $z$ over all possible mappings $X \to B$ and $Y \to \widehat{B}$ where $\mathbb{E}[B] = 1 - a$ and $\mathbb{E}[\widehat{B}] = 1 - b$ is defined as

$$z_l^*(a, b, \mathbf{P}) \triangleq \min_{\substack{\mathbf{x} \in \mathcal{C}^m(a, \mathbf{P}^T) \\ \mathbf{y} \in \mathcal{C}^n(b, \mathbf{P})}} \mathbf{x}^T \mathbf{P} \mathbf{y} \quad \text{and} \quad z_u^*(a, b, \mathbf{P}) \triangleq \max_{\substack{\mathbf{x} \in \mathcal{C}^m(a, \mathbf{P}^T) \\ \mathbf{y} \in \mathcal{C}^n(b, \mathbf{P})}} \mathbf{x}^T \mathbf{P} \mathbf{y},$$

respectively, and $\mathcal{C}^n(a, \mathbf{P})$ is defined in (1).

The next lemma provides a simple upper-bound for $z_u^*(a, b, \mathbf{P})$ in terms of the largest principal inertia components or, equivalently, the maximal correlation between $X$ and $Y$.

**Lemma 4.** $z_u^*(a, b, \mathbf{P}) \leq ab + \rho_m(X; Y)\sqrt{a(1-a)b(1-b)}$.

**Remark 4.** An analogous result was derived by Witsenhausen [11, Thm. 2] for bounding the probability of agreement of a common bit derived from two correlated sources.

*Proof:* Let $\mathbf{x} \in \mathcal{C}^m(a, \mathbf{P}^T)$ and $\mathbf{y} \in \mathcal{C}^n(b, \mathbf{P})$. Then, for $\mathbf{P}$ decomposed as in (4) and $\mathbf{\Sigma}^- = \text{diag}(0, \sigma_1, \ldots, \sigma_d)$,

$$\mathbf{x}^T \mathbf{P}\mathbf{y} = ab + \mathbf{x}^T \mathbf{D}_X^{1/2} \mathbf{U} \mathbf{\Sigma}^- \mathbf{V}^T \mathbf{D}_Y^{1/2} \mathbf{y}$$
$$= ab + \hat{\mathbf{x}}^T \mathbf{\Sigma}^- \hat{\mathbf{y}}, \tag{16}$$

where $\hat{\mathbf{x}} \triangleq \mathbf{U}^T \mathbf{D}_X^{1/2} \mathbf{x}$ and $\hat{\mathbf{y}} \triangleq \mathbf{V}^T \mathbf{D}_Y^{1/2} \mathbf{y}$. Since $\hat{\mathbf{x}}_1 = \|\hat{\mathbf{x}}\|_2 = a$ and $\hat{\mathbf{y}}_1 = \|\hat{\mathbf{y}}\|_2 = b$, then

$$\hat{\mathbf{x}}^T \mathbf{\Sigma}^- \hat{\mathbf{y}} = \sum_{i=2}^{d+1} \sigma_{i-1} \hat{\mathbf{x}}_i \hat{\mathbf{y}}_i$$
$$\leq \sigma_1 \sqrt{\left(\|\hat{\mathbf{x}}\|_2^2 - \hat{\mathbf{x}}_1^2\right)\left(\|\hat{\mathbf{y}}\|_2^2 - \hat{\mathbf{y}}_1^2\right)}$$
$$= \sigma_1 \sqrt{(a - a^2)(b - b^2)}.$$

The result follows by noting that $\sigma_1 = \rho_m(X; Y)$. ■

We will focus in the rest of this section on functions and corresponding estimators that are (i) unbiased ($a = b$) and (ii) satisfy $z = \Pr\{\hat{B} = B = 0\} \geq a^2$. The set of all such mappings is given by

$$\mathcal{H}(a, \mathbf{P}) \triangleq \left\{(\mathbf{x}, \mathbf{y}) \mid \mathbf{x} \in \mathcal{C}^m(a, \mathbf{P}^T), \mathbf{y} \in \mathcal{C}^n(a, \mathbf{P}), \mathbf{x}^T \mathbf{P}\mathbf{y} \geq a^2\right\}.$$

The next results provide upper and lower bounds for $z$ for the mappings in $\mathcal{H}(a, \mathbf{P})$.

**Lemma 5.** *Let $0 \leq a \leq 1/2$ and $\mathbf{P}$ be fixed. For any $(\mathbf{x}, \mathbf{y}) \in \mathcal{H}(a, \mathbf{P})$*

$$a^2 \leq z \leq a^2 + \rho_m(X; Y)a(1-a), \tag{17}$$

*where $z = \mathbf{x}^T \mathbf{P}\mathbf{y}$.*

*Proof:* The lower bound for $z$ follows directly from the definition of $\mathcal{H}(a, \mathbf{P})$, and the upper bound follows from Lemma 4. ■

The previous lemma allows us to provide an upper bound over the mappings in $\mathcal{H}(a, \mathbf{P})$ for the $f$-information between $B$ and $\widehat{B}$ when $I_f$ is non-negative.

**Theorem 3.** *For any non-negative $I_f$ and fixed $a$ and $\mathbf{P}$,*

$$\sup_{(\mathbf{x}, \mathbf{y}) \in \mathcal{H}(a, \mathbf{P})} I_f(B; \hat{B}) \leq a^2 f(1 + \sigma_1 c) + 2a(1-a)f(1 - \sigma_1) + (1-a)^2 f(1 + \sigma_1 c^{-1}) \tag{18}$$

*where here $\sigma_1 = \rho_m(X; \widetilde{Y})$ and $c \triangleq (1-a)a^{-1}$. In particular, for $a = 1/2$,*

$$\sup_{(\mathbf{x}, \mathbf{y}) \in \mathcal{H}(1/2, \mathbf{P})} I_f(B; \hat{B}) \leq \frac{1}{2}\left(f(1 - \sigma_1) + f(1 + \sigma_1)\right). \tag{19}$$

*Proof:* Using the matrix form of the joint distribution between $B$ and $\widehat{B}$ given in (15), for $\mathbb{E}[B] = \mathbb{E}[\widehat{B}] = 1 - a$, the $f$ information is given by

$$I_f(B; \hat{B}) = a^2 f\left(\frac{z}{a^2}\right) + 2a(1-a)f\left(\frac{a - z}{a(1-a)}\right) + (1-a)^2 f\left(\frac{1 - 2a + z}{(1-a)^2}\right). \tag{20}$$

Consequently, (20) is convex in $z$. For $(\mathbf{x}, \mathbf{y}) \in \mathcal{H}(a, \mathbf{P})$, it follows from Lemma 5 that $z$ is restricted to the interval

in (17). Since $I_f(B; \hat{B})$ is non-negative by assumption, $I_f(B; \hat{B}) = 0$ for $z = a^2$ and (20) is convex in $z$, then $I_f(B; \hat{B})$ is non-decreasing in $z$ for $z$ in (17). Substituting $z = a^2 + \rho_m(X; Y)a(1-a)$ in (20), inequality (18) follows. ∎

**Remark 5.** Note that the right-hand side of (18) matches the right-hand side of (11), and provides further evidence for Conjecture 1. This result indicates that, for conforming probability distributions, the information between a binary function and its corresponding unbiased estimate is maximized when all the principal inertia components have the same value.

Following the same approach from Lemma 3, we find the next bound for the mutual information between $B$ and $\hat{B}$.

**Corollary 1.** *For $a$ fixed and $\rho_m(X; Y) = 1 - 2\delta$*

$$\sup_{(\mathbf{x}, \mathbf{y}) \in \mathcal{H}(a, \mathbf{P})} I(B; \hat{B}) \leq 1 - H_b(\delta).$$

We now provide a few application examples for the results derived in this section.

### A. Lower Bounding the Estimation Error Probability

For $z$ given in (15), the average estimation error probability is given by $\Pr\{B \neq \hat{B}\} = a + b - 2z$, which is a convex (linear) function of $z$. If $a$ and $b$ are fixed, then the error probability is minimized when $z$ is maximized. Therefore

$$\Pr\{B \neq \hat{B}\} \geq a + b - 2z_u^*(a, b).$$

Using the bound from Lemma 4, it follows that

$$\Pr\{B \neq \hat{B}\} \geq a + b - 2ab - 2\rho_m(X; Y)\sqrt{a(1-a)b(1-b)}. \tag{21}$$

The bound (21) is exactly the bound derived by Witsenhausen in [11, Thm 2.]. Furthermore, minimizing the right-hand side of (21) over $0 \leq b \leq 1/2$, we arrive at

$$\Pr\{B \neq \hat{B}\} \geq \frac{1}{2}\left(1 - \sqrt{1 - 4a(1-a)(1 - \rho_m(X; Y)^2)}\right), \tag{22}$$

which is a particular form of the bound derived by Calmon *et al.* [15, Thm. 3].

### B. Memoryless Binary Symmetric Channels with Uniform Inputs

We now turn our attention back to the setting considered in Section IV-A. Let $Y^n$ be the result of passing $X^n$ through a memoryless binary symmetric channel with crossover probability $\delta$, $X^n$ uniformly distributed, and $B \to X^n \to Y^n \to \hat{B}$. Then $\rho_m(X^n; Y^n) = 1 - 2\delta$ and, from (22), when $\mathbb{E}[B] = 1/2$,

$$\Pr\{B \neq \hat{B}\} \geq \delta.$$

Consequently, inferring any unbiased one-bit function of the input of a binary symmetric channel is at least as hard (in terms of error probability) as inferring a single output from a single input.

Using the result from Corollary 1, it follows that when $\mathbb{E}[B] = \mathbb{E}[\hat{B}] = a$ and $\Pr\{B = \hat{B} = 0\} \geq a^2$, then

$$I(B; \hat{B}) \leq 1 - H_b(\delta). \tag{23}$$

**Remark 6.** Anantharam *et al.* presented in [8] a computer aided proof that the upper bound (23) holds for any $B \to X^n \to Y^n \to \widehat{B}$. However, we highlight that the methods introduced here allowed an analytical derivation of the inequality (23), which, in turn, is a particular case of the more general setting studied by Anantharam *et al.*

## ACKNOWLEDGEMENT

## REFERENCES

[1] M. Greenacre, Ed., *Theory and Applications of Correspondence Analysis*. Academic Press, Mar. 1984.

[2] M. Greenacre and T. Hastie, "The geometric interpretation of correspondence analysis," *J. Am. Stat. Assoc.*, vol. 82, no. 398, pp. 437–447, Jun. 1987.

[3] H. O. Hirschfeld, "A connection between correlation and contingency," in *Proc. Camb. Philos. Soc.*, vol. 31, 1935, pp. 520–524.

[4] H. Gebelein, "Das statistische problem der korrelation als variations- und eigenwertproblem und sein zusammenhang mit der ausgleichsrechnung," *ZAMM - Z. Angew. Math. Mech.*, vol. 21, no. 6, pp. 364–379, 1941.

[5] O. Sarmanov, "Maximum correlation coefficient (nonsymmetric case)," *Selected Translations in Mathematical Statistics and Probability*, vol. 2, pp. 207–210, 1962.

[6] A. Rényi, "On measures of dependence," *Acta Math. Acad. Sci. H.*, vol. 10, no. 3-4, pp. 441–451, Sep. 1959.

[7] G. R. Kumar and T. A. Courtade, "Which boolean functions are most informative?" in *Proc. 2013 IEEE Int. Symp. on Inf. Theory*, 2013, pp. 226–230.

[8] V. Anantharam, A. Gohari, S. Kamath, and C. Nair, "On hypercontractivity and the mutual information between boolean functions," in *Proc. 51st Ann. Allerton Conf. Commun., Contr., and Comput.*, Oct. 2013, pp. 13–19.

[9] I. Csiszár, *Information Theory And Statistics: A Tutorial*. Now Publishers Inc, 2004.

[10] Y. Polyanskiy and S. Verdú, "Arimoto channel coding converse and Rényi divergence," in *Proc. 48th Ann. Allerton Conf. Commun., Contr., and Comput.*, 2010, pp. 1327–1333.

[11] H. S. Witsenhausen, "On sequences of pairs of dependent random variables," *SIAM J. on Appl. Math.*, vol. 28, no. 1, pp. 100–113, Jan. 1975.

[12] R. Ahlswede and P. Gács, "Spreading of sets in product spaces and hypercontraction of the markov operator," *Ann. Probab.*, vol. 4, no. 6, pp. 925–939, Dec. 1976.

[13] V. Anantharam, A. Gohari, S. Kamath, and C. Nair, "On Maximal Correlation, Hypercontractivity, and the Data Processing Inequality studied by Erkip and Cover," *arXiv:1304.6133 [cs.IT]*, Apr. 2013.

[14] Y. Polyanskiy, "Hypothesis testing via a comparator," in *Proc. 2012 IEEE Int. Symp. on Inf. Theory*, Jul. 2012, pp. 2206–2210.

[15] F. P. Calmon, M. Varia, M. Médard, M. Christiansen, K. Duffy, and S. Tessaro, "Bounds on inference," in *Proc. 51st Ann. Allerton Conf. Commun., Contr., and Comput.*, Oct. 2013, pp. 567–574.

[16] W. Kang and S. Ulukus, "A new data processing inequality and its applications in distributed source and channel coding," *IEEE Trans. Inform. Theory*, vol. 57, no. 1, pp. 56–69, 2011.

[17] R. O'Donnell, "Some topics in analysis of boolean functions," in *Proc. 40th ACM Symp. on Theory of Computing*, 2008, pp. 569–578.

[18] M. Raginsky, J. G. Silva, S. Lazebnik, and R. Willett, "A recursive procedure for density estimation on the binary hypercube," *Electron. J. Statist.*, vol. 7, pp. 820–858, 2013.