

The Double-Sided Information-Bottleneck Function

Michael Dikshstein
Technion
Email: michaeldic@campus.technion.ac.il

Or Ordentlich
Hebrew University of Jerusalem
Email: or.ordentlich@mail.huji.ac.il

Shlomo Shamai (Shitz)
Technion
Email: sshlomo@ee.technion.ac.il

Abstract—We consider a two-terminal variant (double-sided) of the information bottleneck problem, which is related to biclustering. In our setup, X and Y are dependent random variables and the problem is to find two independent channels $P_{U|X}$ and $P_{V|Y}$ (setting the Markovian structure $U \rightarrow X \rightarrow Y \rightarrow V$) that maximize $I(U; V)$ subject to constraints on the relevant mutual information expressions: $I(U; X)$ and $I(V; Y)$. For jointly Gaussian X and Y , we show that Gaussian channels are optimal in the low-SNR regime, but not for general SNR. Similarly, it is shown that for a doubly symmetric binary source, binary symmetric channels are optimal when the correlation is low, and are suboptimal for high correlation. We conjecture that Z and S channels are optimal when the correlation is 1 (i.e., $X = Y$), and provide supporting numerical evidence.

I. INTRODUCTION

Let (X, Y) be a bivariate source characterized by a fixed joint probability law P_{XY} and consider all Markov chains $U \rightarrow X \rightarrow Y \rightarrow V$. The Double-Sided Information-Bottleneck (DSIB) function is defined as [1]

$$R_{P_{XY}}(C_u, C_v) \triangleq \max I(U; V), \quad (1)$$

where the maximization is over all $P_{U|X}$ and $P_{V|Y}$ satisfying $I(U; X) \leq C_u$ and $I(V; Y) \leq C_v$. This problem is illustrated in Figure 1. In our study we aim to determine the maximum value and the achieving conditional distributions ($P_{U|X}, P_{V|Y}$) (test channels) of $I(U; V)$ for various fixed sources P_{XY} , and constraints C_u, C_v .

The problem we consider originates from the domain of clustering. Clustering is applied to organize similar entities in unsupervised learning [2]. It has numerous practical applications in data science, such as joint word-document clustering, gene expression [3], and pattern recognition. The data in those applications is arranged as a contingency table. Usually, clustering is performed on one dimension of the table, but sometimes it is useful to apply clustering on both dimensions of the contingency table [4]. For example, when there is a

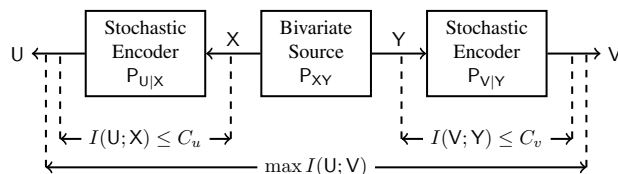


Fig. 1. Block diagram of the Double-Sided Information-Bottleneck function.

strong correlation between the rows and the columns of the table or when high-dimensional sparse structures are handled.

An information-theoretic approach to biclustering was initiated by Dhillon et al. [5]. They have regarded the normalized nonnegative contingency table as a joint probability distribution matrix of two random variables. Mutual information was proposed as a measure for optimal co-clustering. An optimization algorithm was presented that intertwines both row and column clustering at all stages. Distributed Clustering from a formal information-theoretic perspective was first explicitly considered by Pichler et al. [1], [6]. Consider a bivariate memoryless source with joint law P_{XY} . This source generates n i.i.d. copies (X^n, Y^n) , of (X, Y) . Each sequence is observed at two different encoders, and each encoder generates a description of the observed sequence, $f(X^n)$, and respectively, $g(Y^n)$. The objective is to construct the mappings f and g such that the normalized mutual information between the descriptions would be maximal, while the description coding has bounded rate constraints. Single-letter inner and outer bounds for a general P_{XY} were derived. An example of a Doubly Symmetric Binary Source (DSBS) source was given, and several converse results were established. Furthermore, connections to the standard Information Bottleneck [7] and the Multiple Description CEO problems [8] were given. It was also shown that an information-theoretic biclustering problem is equivalent to hypothesis testing against independence with multiterminal data compression [9] and a pattern recognition problem [10].

The DSIB problem addressed in our paper is, in fact, a single letter version of the Distributed Clustering problem. The inner bound in [1] coincides with our definition of the problem. Moreover, if the Markov condition $U \rightarrow X \rightarrow Y \rightarrow Z$ is imposed on the multiletter variant, then those problems coincide.

A similar setting with a maximal correlation criterion between the reconstructed random variables has been considered in [11], [12]. Furthermore, it is sometimes the case that the optimal biclustering problem is easier to solve than the standard single-sided clustering problem. For example, Courtade-Kumar conjecture [13] for the standard single-sided clustering setting was ultimately proved for the biclustering setting [14]. A particular case, where (X, Y) are drawn from DSBS distribution, and the mappings f and g are restricted to be boolean functions, was addressed in [14]. The bound $I(f(X^n); g(Y^n)) \leq I(X; Y)$ was established, which is tight if and only if f and g are dictator functions. There are many

other variations of multi-user IB in the literature [15], [16], [17], [18], [19], [20], [21], [22], [23], [24], [25].

A recent comprehensive tutorial on the Information Bottleneck method and related problems is given in [26]. Applications of Information Bottleneck problem in Machine Learning are detailed in [17], [27], [28], [29], [30], [31], [32].

II. PROBLEM FORMULATION AND BASIC PROPERTIES

The DSIB function is closely related to the standard Single-Sided Information Bottleneck (SSIB) [7].

Definition 1 (SSIB): Let (X, V) be a pair of random variables with $|\mathcal{X}| = n$, $|\mathcal{V}| = m$ and fixed P_{XV} . Denote by \mathbf{q}_x the marginal probability vector of X , and let T be the transition matrix from X to V , i.e.,

$$T_{ij} \triangleq P(V = i | X = j), \quad 1 \leq i \leq m, 1 \leq j \leq n.$$

Consider all random variables U satisfying the Markov chain $U \rightarrow X \rightarrow V$. The SSIB function is defined as

$$\hat{R}_T(\mathbf{q}_x, C) \triangleq \begin{aligned} & \underset{P_{U|X}}{\text{maximize}} && I(U; V) \\ & \text{subject to} && I(X; U) \leq C. \end{aligned} \quad (2)$$

The following tight cardinality bound in the single-sided counterpart of our problem was established in [33].

Lemma 1 ([33, Th. 9]): The optimization over U in (2) can be restricted to $|\mathcal{U}| \leq n$.

This bound was actually already proved for the corresponding dual problem, namely, the IB Lagrangian, in [34]. But since $\hat{R}_T(\mathbf{q}_x, C)$ is generally not a strictly convex function of C , it cannot be directly applied for the primal problem (2).

Definition 2 (DSIB): Let (X, Y) be a pair of random variables with $|\mathcal{X}| = n$, $|\mathcal{Y}| = m$ and fixed P_{XY} . Consider all the random variables U and V satisfying the Markov chain $U \rightarrow X \rightarrow Y \rightarrow V$. The DSIB function, $R : [0, H(X)] \times [0, H(Y)] \rightarrow \mathbb{R}_+$ is defined as

$$R_{P_{XY}}(C_u, C_v) \triangleq \begin{aligned} & \underset{P_{U|X}, P_{V|Y}}{\text{maximize}} && I(U; V) \\ & \text{subject to} && I(X; U) \leq C_u, I(Y; V) \leq C_v. \end{aligned} \quad (3)$$

Occasionally we will drop the subscript denoting the particular choice of the bivariate source P_{XY} .

An equivalent form of (3) is

$$R(C_u, C_v) \triangleq \begin{aligned} & \underset{P_{V|Y}}{\text{maximize}} && \underset{P_{U|X}}{\text{maximize}} && I(U; V). \\ & \text{subject to} && \text{subject to} && \\ & && I(Y; V) \leq C_v && I(X; U) \leq C_u \end{aligned} \quad (4)$$

Fix $P_{V|Y}$ that satisfies $I(Y; V) \leq C_v$. Denote by T_v the transition matrix from Y to V and respectively by W the transition matrix from X to Y . Since $P_{V|Y}$ and P_{XY} are fixed, $P_{XV} = \sum_y P_{V|Y}(\cdot|y)P_{XY}(\cdot, y)$ is also fixed. Denote by $T \triangleq T_v W$ the transition matrix from X to V . Therefore, the inner maximization term in (4) is just the SSIB function with parameters T and C_u , namely, $\hat{R}_T(\mathbf{q}_x, C_u)$. Hence, our

problem can also be defined in the following two equivalent ways:

$$R(C_u, C_v) \triangleq \begin{aligned} & \underset{T_v}{\text{maximize}} && \hat{R}_{T_v W}(\mathbf{q}_x, C_u) \\ & \text{subject to} && I(\mathbf{q}_y, T_v) \leq C_v, \end{aligned} \quad (5)$$

or

$$R(C_u, C_v) \triangleq \begin{aligned} & \underset{T_u}{\text{maximize}} && \hat{R}_{T_u \bar{W}}(\mathbf{q}_y, C_v) \\ & \text{subject to} && I(\mathbf{q}_x, T_u) \leq C_u. \end{aligned} \quad (6)$$

where T_u is the transition matrix from X to U , and \bar{W} is the transition matrix from Y to X . This representation gives us a different perspective on our problem as an optimal channel for the SSIB setting.

The bound from Lemma 1 can be utilized to give cardinality bound for the double sided problem.

Proposition 2.1: For the DSIB optimization problem defined in (3) it suffices to consider random variables U and V with cardinalities $|\mathcal{U}| \leq n$ and $|\mathcal{V}| \leq m$.

Proof. Let T_u and T_v be two arbitrary transition matrices. By Lemma 1, there exists \tilde{T}_u with $|\tilde{\mathcal{U}}| \leq n$ such that $I(\tilde{U}; V) \geq I(U; V)$ and $I(X; \tilde{U}) \leq C_u$. Similarly, T_v can be replaced with \tilde{T}_v , $|\tilde{\mathcal{V}}| \leq m$ such that

$$I(\tilde{U}; \tilde{V}) \geq I(\tilde{U}; V) \geq I(U; V),$$

and $I(Y; \tilde{V}) \leq C_v$. Therefore, there exists an optimal solution with $|\mathcal{U}| = n$ and $|\mathcal{V}| = m$. ■

III. MAIN RESULTS

In this section we will present the main analytical outcomes of our study. First we consider the scenario where our bivariate source is binary, and specifically DSBS. Then we treat the case where X and Y are jointly Gaussian.

A. Doubly Symmetric Binary Source (DSBS)

Let (X, Y) be a DSBS with parameter p , i.e., $P_{XY}(x, y) = \frac{1}{2}(p \cdot \mathbb{1}(x \neq y) + (1-p)\mathbb{1}(x = y))$. Here we emphasize the dependence of the DSIB on the parameter p as $R(C_u, C_v, p)$.

A special case where we have a complete analytical solution is when p tends to $1/2$. Let $h_2(p) : [0, 1] \rightarrow [0, 1]$ be the binary entropy function and $h_2^{-1}(\cdot)$ its inverse, restricted to $[0, 1/2]$. Throughout this paper all logarithms are taken to base 2 unless stated otherwise.

Theorem 1: Suppose $p = \frac{1}{2} - \epsilon$, and consider $\epsilon \rightarrow 0$. Then
$$R(C_u, C_v, \epsilon) = 2\epsilon^2 \log_e(1 - 2h_2^{-1}(1 - C_u))^2 (1 - 2h_2^{-1}(1 - C_v))^2 + o(\epsilon^2). \quad (7)$$

This theorem is proved in Appendix A. For the lower bound we take $P_{U|X}$ and $P_{V|Y}$ to be Binary Symmetric Channels (BSCs). In the following section we will give a numerical evidence that BSC test-channels are in fact optimal provided that p is sufficiently large. However, for small p this is no longer the case and we believe the following holds.

Conjecture 1: Let $X = Y$, i.e., $p = 0$. The test channels $P_{U|X}$ and $P_{V|X}$ that achieve $R(C_u, C_v, 0)$ are Z-channel and S-channel respectively.

Remark 1: Our results in the numerical section strongly support this conjecture. In fact they prove it within the resolution of the experiments, i.e., for optimizing over a dense set of test-channels rather than all test-channels. Nevertheless, we were not able to find an analytical proof for this result.

Remark 2: Suppose $X = Y$, $I(X; U) = C_u$, and $I(X; V) = C_v$. Since $I(U; V) = I(U; X) + I(V; X) - I(X; U, V)$ (as $U \rightarrow X \rightarrow Y \rightarrow V$ form a Markov chain in this order) then maximizing $I(U; V)$ is equivalent to minimizing $I(X; U, V)$, namely, minimizing information combining as in [35], [22]. Therefore, Conjecture 1 is equivalent to the conjecture that among all channels with $I(X; U) \geq C_u$ and $I(Y; V) \geq C_v$, Z and S are the worst channels for information combining.

This observation leads us the following additional conjecture.

Conjecture 2: The test channels $P_{U|X}$ and $P_{V|X}$ that maximize $I(X; U, V)$ are both Z channels.

Remark 3: Suppose now that p is arbitrary and assume that one of the channels $P_{U|X}$ or $P_{V|Y}$ is restricted to be a Binary Memoryless Symmetric (BMS) channel [36, Ch. 4], then the maximal $I(U; V)$ is attained by BSC channels, as those are the symmetric channels minimizing $I(X; U, V)$ [35]. It is not surprising that once the BMS constraint is removed, symmetric channels are no longer optimal (see the discussion in [35, Sec. VI.C]).

B. Gaussian Double-Sided Information Bottleneck (GDSIB)

In this subsection we consider a specific setting where (X, Y) is a Gaussian bivariate source, namely,

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right). \quad (8)$$

We establish achievability schemes and show that Gaussian test-channels $P_{U|X}$ and $P_{V|Y}$ are optimal for vanishing SNR. Furthermore we show an elegant representation of the problem through *Hermite polynomials*. We denote the Gaussian DSIB function with explicit dependency on ρ as $R(C_u, C_v, \rho)$.

Proposition 3.1: Let $H_n(\cdot)$ be the n th order probabilistic Hermite polynomial, then,

$$I(U; V) = \mathbb{E}_{U, V} \left[\log \left(\sum_{n=0}^{\infty} \frac{\rho^n}{n!} \mathbb{E} [H_n(X)|U] \mathbb{E} [H_n(Y)|V] \right) \right]. \quad (9)$$

This representation follows by considering $I(U; V) = D(P_{U, V} || P_U \cdot P_V)$ and expressing $\frac{P_{U, V}}{P_U \cdot P_V}$ using Mehler Kernel [37]. Mehler Kernel decomposition is a special case of a much richer family of Lancaster distributions [38].

Now we give two lower bounds on $R(C_u, C_v, \rho)$. Our first lower bound is established by choosing $P_{U|X}$ and $P_{V|Y}$ to be additive Gaussian channels, satisfying the mutual information (MI) constraints with equality.

Proposition 3.2: A lower bound on $R(C_u, C_v, \rho)$ is given by

$$R(C_u, C_v, \rho) \geq -\frac{1}{2} \log(1 - \rho^2 (1 - 2^{-2C_u}) (1 - 2^{-2C_v})). \quad (10)$$

Although it was shown in [17] that choosing the test channel to be Gaussian is optimal for the single-sided variant, it is not the case for its double-sided extension. We will show this by

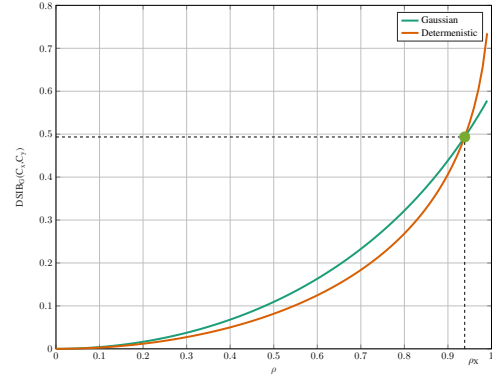


Fig. 2. Comparison of the lower bounds from Proposition 3.2 and Proposition 3.3.

examining a specific set of values for the rates constraints, $(C_u, C_v) = (1, 1)$. Furthermore, we choose the test channels $P_{U|X}$ and $P_{V|Y}$ to be deterministic quantizers.

Proposition 3.3: Let $(C_u, C_v) = (1, 1)$, then

$$R(1, 1, \rho) \geq 1 - h_2 \left(\frac{\arccos \rho}{\pi} \right). \quad (11)$$

The proof of this bound is developed in Appendix B.

We compare the bounds from Proposition 3.2 and Proposition 3.3 with $(C_u, C_v) = (1, 1)$ in Figure 2. The most unexpected observation here is that the deterministic quantizers lower bound outperform the Gaussian test-channels for high values of ρ . The crossing point of those bounds is

$$\rho_{\text{cross}} = \frac{e}{\sqrt{1+e^2}} \rightarrow \sqrt{SNR_{\text{cross}}} = \frac{\rho_{\text{cross}}}{\sqrt{1-\rho_{\text{cross}}^2}} = e. \quad (12)$$

For $\rho \rightarrow 0$, the exact asymptotic behavior of the Gaussian and deterministic test-channels, respectively, for $C_u = C_v = 1$ bit is

$$\lim_{\rho \rightarrow 0} -\frac{1}{2} \log(1 - \rho^2 (1 - 2^{-2C_u}) (1 - 2^{-2C_v})) = \frac{9 \log e}{32} \rho^2 + o(\rho^2),$$

$$\lim_{\rho \rightarrow 0} 1 - h_2 \left(\frac{\arccos \rho}{\pi} \right) = \frac{2 \log e}{\pi^2} \rho^2 + o(\rho^2).$$

Hence, the Gaussian choice outperforms the second lower bound for vanishing SNR.

Theorem 2: For small ρ , the GDSIB function is given by

$$R(C_u, C_v, \rho) = \frac{\rho^2 \log e}{2} (1 - 2^{-2C_u}) (1 - 2^{-2C_v}) + o(\rho^2). \quad (13)$$

The lower bound follows from Proposition 3.2. The upper bound is established by considering the kernel representation from Proposition 3.1 in the limit of vanishing ρ . The detailed proof is given in Appendix C.

IV. NUMERICAL RESULTS

In this section we consider the case where $P_{X,Y}$ is a DSBS with parameter p . Since V is binary, we can describe $P_{V|Y}$ using the following transition matrix

$$T_v \triangleq \begin{pmatrix} \alpha & \beta \\ 1 - \alpha & 1 - \beta \end{pmatrix}. \quad (14)$$

Now, let \mathcal{F}_v be the set of all feasible $P_{V|Y}$, i.e.,

$$\mathcal{F}_v(C_v) \triangleq \left\{ (\alpha, \beta) : h\left(\frac{\alpha + \beta}{2}\right) - \frac{1}{2}(h(\alpha) + h(\beta)) = C_v \right\}$$

(equality constraint follows since we maximize a convex function over a convex set). This set defines a curve $\beta(\alpha) \triangleq \beta(\alpha, C_v)$. Examples of such curves for various values of C_v are shown in Figure 3.

Fix $(\alpha, \beta(\alpha)) \in \mathcal{F}_v$ and consider $\hat{R}_T(\mathbf{q}, C_u) = \hat{R}(C_u, \alpha, p)$ from (2), where $\bar{x} \triangleq 1 - x$ and

$$T = T_v W = \begin{pmatrix} \alpha\bar{p} + \beta(\alpha)p & \alpha p + \beta(\alpha)\bar{p} \\ \bar{\alpha}\bar{p} + \bar{\beta}(\alpha)p & \bar{\alpha}p + \bar{\beta}(\alpha)\bar{p} \end{pmatrix}. \quad (15)$$

Since SSIB is equivalent to the problem that was studied in [39] (as was shown in [40, Sec. 3]), we can utilize the recipe from [39, Sec. IV] to evaluate $\hat{R}(C_u, \alpha, p)$. With that tool in hand, we can compute $\hat{R}(C_u, \alpha, p)$ for every $\alpha \in \mathcal{F}_v$ and then find α^* that maximizes $\hat{R}(C_u, \alpha, p)$. This is equivalent to

$$R(C_u, C_v, p) \triangleq \max_{\alpha \in \mathcal{F}_v} \hat{R}(C_u, \alpha, p). \quad (16)$$

The graph of $\hat{R}(C_u, \alpha, p)$ for $C_u = C_v = 0.3$ bit and various values of p is shown in Figure 4. For $p = 0$ ($X = Y$) we see that $\alpha^* = 1$, in such case $P_{V|Y}$ and $P_{U|X}$ are Z and S channels, respectively. For small but increasing values of p we observe a smooth transition from Z-channel, $\alpha^*(p = 0)$, points on the curve $\beta(\alpha)$, towards the point $\alpha = 1 - \beta$, which corresponds to a BSC. The graphs of $R(C_u, C_v, \cdot)$, $\alpha^*(\cdot)$, $\beta^*(\alpha^*(\cdot))$ and $\alpha^*(\cdot) + \beta^*(\alpha^*(\cdot))$ are shown in Figure 5. As expected $R(C_u, C_v, \cdot)$ is a decreasing function of p . Furthermore, the channel transition probability has a threshold value, θ , which depends on C_u and C_v . If $p > \theta(C_u, C_v)$ then $P_{V|Y}$ is a BSC (and thus also $P_{U|X}$ [39, Sec. IV.A]).

V. CONCLUDING REMARKS

In this paper we have considered the Double-Sided Information-Bottleneck problem. Tight cardinality bounds on the auxiliary random variables were obtained for an arbitrary discrete bivariate source. For DSBS we have shown that BSC test-channel are optimal when $p \rightarrow 0.5$. Furthermore, we have shown numerical simulation for arbitrary p , indicating that Z

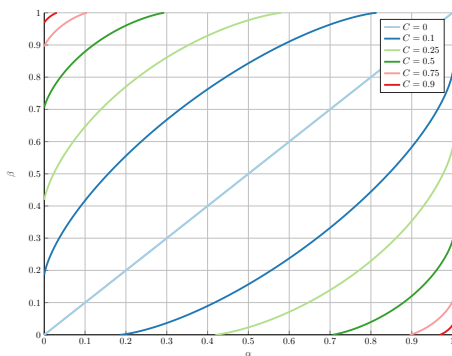


Fig. 3. Set of feasible curves $\beta(\alpha, C)$ for various values of C .

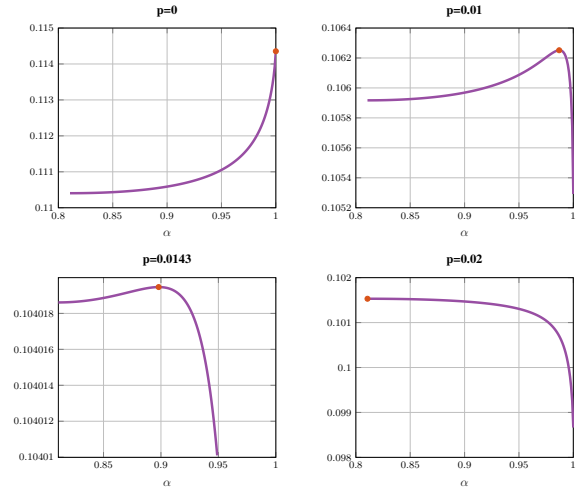


Fig. 4. Evaluation of $\hat{R}(C_u, \alpha, p)$ for $C_u = C_v = 0.3$ bit and various values of p . Recall that $1 - h_2^{-1}(1 - 0.3) = 0.8107$.

and S channels are optimal for $p = 0$. As for the Gaussian bivariate source, representation of $I(U; V)$ utilizing Hermite polynomials was given. Also, the optimality of Gaussian test-channel was demonstrated for vanishing SNR. Moreover, we have constructed a lower bound attained by deterministic quantizers that outperforms the jointly Gaussian choice at high SNR. Note that the solution for the n -letter problem $\max \frac{1}{n} I(U; V)$ for $U \rightarrow X^n \rightarrow Y^n \rightarrow V$ under constraints $I(U; X^n) \leq nC_u$ and $I(V; Y^n) \leq nC_v$, does not tensorize in general. For $X^n = Y^n \sim \text{Ber}^{\otimes n}(0.5)$, we can easily achieve the cut-set bound $I(U; V)/n = \min\{C_u, C_v\}$. In addition, if time sharing is allowed, the results change drastically.

ACKNOWLEDGMENT

The work has been supported by the European Union's Horizon 2020 Research And Innovation Programme, grant

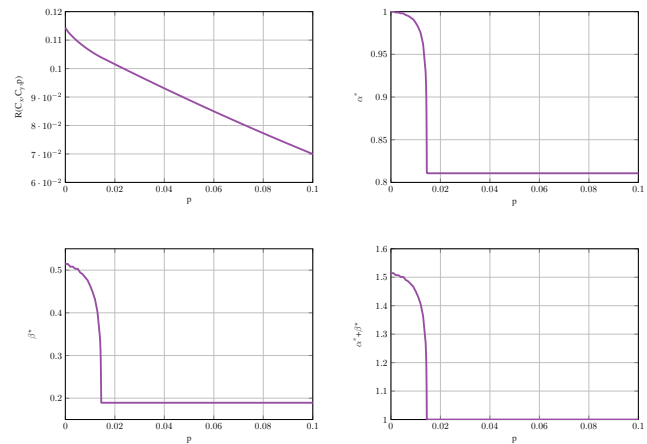


Fig. 5. Evaluation of $R(C_u, C_v, p)$, α^* , $\beta^* \triangleq \beta(\alpha^*)$, and $\alpha^* + \beta^*$ for $C_u = C_v = 0.3$ bit.

agreement no. 694630, by the ISF under Grant 1791/17, and by the WIN consortium via the Israel minister of economy and science.

APPENDIX

A. Proof of Theorem 1

Fix $P_{U|X}$ and $P_{V|Y}$ that satisfy $I(X; U) \leq C_u$ and $I(Y; V) \leq C_v$. Denote $\alpha_u \triangleq P(X = 1|U = u)$, $\beta_v \triangleq P(Y = 1|V = v)$. Using this notation we have

$$I(U; V) = \sum_{u,v} P_U(u)P_V(v)K(u, v, p) \log K(u, v, p), \quad (17)$$

where $K(u, v, p) \triangleq \frac{P_{UV}(u,v)}{P_U(u)P_V(v)}$, and satisfies

$$K(u, v, p) = 2(\bar{\alpha}_u \bar{\beta}_v + \alpha_u p \bar{\beta}_v + \bar{\alpha}_u p \beta_v + \alpha_u \bar{p} \beta_v) = 2\alpha_u * \beta_v * \bar{p}. \quad (18)$$

Denoting $\epsilon \triangleq \frac{1}{2} - p$ ($p = \frac{1}{2} - \epsilon$), we obtain

$$K(u, v, \epsilon) = 1 + 2\epsilon(1 - 2\alpha_u)(1 - 2\beta_v). \quad (19)$$

Now we rewrite $I(U; V)$ with explicit dependency on ϵ as $I(\epsilon) = \sum_{u,v} P_U(u)P_V(v)K(u, v, \epsilon) \log K(u, v, \epsilon)$. We would like to expand $I(\epsilon)$ with Taylor series around $\epsilon = 0$. Note that $I(0) = 0 = I'(\epsilon)|_{\epsilon=0}$. Furthermore, the second derivative is given by

$$I''(\epsilon)|_{\epsilon=0} = 4 \log e \cdot \left(\sum_u P_U(u)(1-2\alpha_u)^2 \right) \left(\sum_v P_V(v)(1-2\beta_v)^2 \right).$$

Hence,

$$I(\epsilon) = 2\epsilon^2 \log e \cdot \left(\sum_u P_U(u)(1-2\alpha_u)^2 \right) \left(\sum_v P_V(v)(1-2\beta_v)^2 \right) + o(\epsilon^2).$$

Now, note that

$$\alpha_u = \begin{cases} h_2^{-1}(H(X|U=u)) & \alpha_u \leq \frac{1}{2} \\ 1 - h_2^{-1}(H(X|U=u)) & \alpha_u > \frac{1}{2} \end{cases} \quad (20)$$

with similar relation for β_v . Therefore,

$$\begin{aligned} I(\epsilon) &= \frac{2\epsilon^2}{\ln 2} \mathbb{E}_u [1 - 2h_2^{-1}(H(X|U=u))]^2 \mathbb{E}_v [1 - 2h_2^{-1}(H(Y|V=v))]^2 + o(\epsilon^2) \\ &\leq 2\epsilon^2 \log e \cdot (1 - 2h_2^{-1}(H(X|U)))^2 (1 - 2h_2^{-1}(H(Y|V)))^2 + o(\epsilon^2) \\ &\leq 2\epsilon^2 \log e \cdot (1 - 2h_2^{-1}(1 - C_x))^2 (1 - 2h_2^{-1}(1 - C_y))^2 + o(\epsilon^2), \end{aligned}$$

where the first inequality follows since the function $f : x \mapsto (1 - 2h_2^{-1}(x))^2$ is concave and applying Jensen's inequality (the proof is omitted due to space limitation), and the second inequality follows from rate constraints.

B. Proof of Proposition 3.3

We choose U and V to be deterministic functions of X and respectively Y , i.e., $U = \text{sign}(X)$ and $V = \text{sign}(Y)$. In such case the rate constraints are met with equality, namely, $I(U; X) = 1 = I(Y; V)$. We proceed to evaluate the achievable rate,

$$\begin{aligned} I(U; V) &= 1 - P(U=0)h_2(P(V=1|U=0)) - P(U=1)h_2(P(V=0|U=1)) \\ &\stackrel{(a)}{=} 1 - h_2(P(U \neq V)), \end{aligned}$$

where equality in (a) follows since $P(V = 1|U = 0) = P(V = 0|U = 1)$ by symmetry. We therefore obtain the following formula for the "error probability":

$$P(V \neq U) = 1 - P(X < 0, Y < 0) - P(X > 0, Y > 0) \stackrel{(a)}{=} 1 - 2P(X < 0, Y < 0),$$

where (a) also follows from symmetry. Utilizing Sheppard's Formula [41, Ch. 5, p.107], we have $1 - 2P(X < 0, Y < 0) = \frac{\arccos \rho}{\pi}$. This completes the proof of the proposition.

C. Proof of Theorem 2

We assume U and V are continuous RVs. The proof for the discrete case is identical. The joint density $f_{UV}(u, v)$ can be expressed with explicit dependency on ρ as follows:

$$f(u, v; \rho) \triangleq f_U(u)f_V(v) \iint_{\mathbb{R}^2} f_{X|U}(x|u)M(x, y; \rho)f_{Y|V}(y|v)dx dy,$$

where $M(x, y; \rho) = \sum_{n=0}^{\infty} \frac{\rho^n}{n!} H_n(x)H_n(y)$ [37]. Similarly, $I(U; V)$ can also be written with explicit dependency on ρ

$$I(\rho) \triangleq I_\rho(U; V) = \int \int f(u, v; \rho) \log \frac{f(u, v; \rho)}{f_U(u)f_V(v)} du dv.$$

We would like to approximate $I(\rho)$ in the limit $\rho \rightarrow 0$ using a Taylor series up to a second order in ρ . As a first step we evaluate the first two derivatives of $f(u, v; \rho)$ at $\rho = 0$. Note that $M(x, y; 0) = 1$ and

$$\frac{dM}{d\rho} \Big|_{\rho=0} = xy, \quad \frac{d^2M}{d\rho^2} \Big|_{\rho=0} = (x^2 - 1)(y^2 - 1). \quad (21)$$

Thus, $f(u, v; 0) = f_U(u)f_V(v)$,

$$\frac{df}{d\rho} \Big|_{\rho=0} = f_U(u)f_V(v) \mathbb{E}[X|U=u] \mathbb{E}[Y|V=v],$$

and $\frac{d^2f}{d\rho^2} \Big|_{\rho=0} = f_U(u)f_V(v) (\mathbb{E}[X^2|U=u] - 1) (\mathbb{E}[Y^2|V=v] - 1)$.

Expanding $I(\rho)$ in Taylor series around $\rho = 0$ gives us $I(0) = 0 = \frac{dI(\rho)}{d\rho} \Big|_{\rho=0}$ and

$$\frac{d^2I(\rho)}{d\rho^2} \Big|_{\rho=0} = \log e \cdot \mathbb{E} [(\mathbb{E}[X|U])^2] \mathbb{E} [(\mathbb{E}[Y|V])^2].$$

Thus

$$I(\rho) = \frac{\rho^2 \log e}{2} \mathbb{E} [(\mathbb{E}[X|U])^2] \mathbb{E} [(\mathbb{E}[Y|V])^2] + o(\rho^2). \quad (22)$$

Note that $\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|U]]$ and

$$1 = \mathbb{E}[X^2] = \mathbb{E}[\mathbb{E}[X^2|U]] = \mathbb{E}[\text{var}[X|U] + (\mathbb{E}[X|U])^2]. \quad (23)$$

Also, by [42, Corollary to Theorem 8.6.6], $\mathbb{E}[\text{var}[X|U]] \geq \frac{1}{2\pi e} e^{2h(X|U)}$. Moreover, from MI constraint, we have

$$I(X; U) = h(X) - h(X|U) = \frac{1}{2} \log(2\pi e) - h(X|U) \leq C_u,$$

and therefore $h(X|U) \geq \log(2\pi e) - C_u$. Thus we get

$$-C_u \leq \frac{1}{2} \log(\mathbb{E}[\text{var}[X|U]]) \rightarrow \mathbb{E}[\text{var}[X|U]] \geq 2^{-2C_u}. \quad (24)$$

Combining (23) and (24), we obtain $\mathbb{E}[(\mathbb{E}[X|U])^2] \leq 1 - 2^{-2C_u}$. In a very similar method one can show that $\mathbb{E}[(\mathbb{E}[Y|V])^2] \leq 1 - 2^{-2C_v}$. Thus, for $\rho \rightarrow 0$

$$I(\rho) \leq \frac{\rho^2 \log e}{2} (1 - 2^{-2C_u})(1 - 2^{-2C_v}) + o(\rho^2). \quad (25)$$

REFERENCES

- [1] G. Pichler, P. Piantanida, and G. Matz, "Distributed Information-Theoretic Biclustering," in *Proc. 2016 IEEE Int. Symp. Inf. Theory*, July 2016, pp. 1083–1087.
- [2] A. K. Jain and R. C. Dubes, *Algorithms for Clustering Data*. Prentice-Hall, 1988.
- [3] N. Gupta and S. Aggarwal, "Modeling Biclustering as an optimization problem using Mutual Information," in *2009 Proceeding of International Conference on Methods and Models in Computer Science (ICM2CS)*, 2009, pp. 1–5.
- [4] J. Hartigan, "Direct Clustering of a Data Matrix," *Journal of the American Statistical Association*, vol. 67, no. 337, pp. 123–129, 1972.
- [5] I. S. Dhillon, S. Mallela, and D. S. Modha, "Information-Theoretic Co-clustering," in *Proc. The Ninth ACM SIGKDD Int. Conf. Knowledge Discovery Data Mining*, ser. (KDD '03), 2003.
- [6] G. Pichler, P. Piantanida, and G. Matz, "Distributed Information-Theoretic Clustering," *arXiv*, vol. abs/1602.04605, 2020.
- [7] N. Tishby, F. C. N. Pereira, and W. Bialek, "The information bottleneck method," in *Proc. 37th Annu. Allerton Conf. Commun. Control Comput.*, Sep. 1999, pp. 368–377.
- [8] T. A. Courtade and G. R. Kumar, "Which Boolean Functions Maximize Mutual Information on Noisy Inputs?" *IEEE Trans. Inf. Theory*, vol. 60, no. 8, pp. 4515–4525, 2014.
- [9] T. S. Han, "Hypothesis Testing with Multiterminal Data Compression," *IEEE Trans. Inf. Theory*, vol. 33, no. 6, pp. 759–772, 1987.
- [10] M. B. Westover and J. A. O'Sullivan, "Achievable Rates for Pattern Recognition," *IEEE Trans. Inf. Theory*, vol. 54, no. 1, pp. 299–320, 2008.
- [11] A. Painsky, M. Feder, and N. Tishby, "An information-theoretic framework for non-linear canonical correlation analysis," *CoRR*, vol. abs/1810.13259, 2018.
- [12] A. R. Williamson, "The Impacts of Additive Noise and 1-bit Quantization on the Correlation Coefficient in the Low-SNR Regime," in *Proc. 57th Annu. Allerton Conf. Commun. Control Comput.*, Sep. 2019, pp. 631–638.
- [13] T. A. Courtade and T. Weissman, "Multiterminal Source Coding Under Logarithmic Loss," *IEEE Trans. Inf. Theory*, vol. 60, no. 1, pp. 740–761, 2014.
- [14] G. Pichler, P. Piantanida, and G. Matz, "Dictator Functions Maximize Mutual Information," *Ann. Appl. Prob.*, vol. 28, no. 5, pp. 3094–3101, 2018.
- [15] I. Estella and A. Zaidi, "Distributed Information Bottleneck Method for Discrete and Gaussian Sources," in *Proc. Int. Zurich Seminar Inf. Commun. (IZS)*. ETH Zurich, 2018, pp. 35–39.
- [16] T. Berger, Zhen Zhang, and H. Viswanathan, "The CEO Problem," *IEEE Trans. Inf. Theory*, vol. 42, no. 3, pp. 887–902, 1996.
- [17] G. Chechik, A. Globerson, N. Tishby, and Y. Weiss, "Information Bottleneck for Gaussian Variables," *J. Mach. Learn. Res.*, vol. 6, pp. 165–188, Dec. 2005.
- [18] T. A. Courtade and J. Jiao, "An Extremal Inequality for Long Markov Chains," in *Proc. 52nd Annu. Allerton Conf. Commun., Control Comput.*, Oct. 2014, pp. 763–770.
- [19] E. Erkip and T. M. Cover, "The Efficiency of Investment Information," *IEEE Trans. Inf. Theory*, vol. 44, no. 3, pp. 1026–1040, May 1998.
- [20] I. E. Aguerri and A. Zaidi, "Distributed Variational Representation Learning," *IEEE Trans. Pattern Anal.*, vol. 43, no. 1, pp. 120–138, 2021.
- [21] P. Gács and J. Körner, "Common information is far less than mutual information," *Probl. Contr. Inform. Theory*, vol. 2, no. 2, pp. 149–162, 1973.
- [22] I. Land and J. Huber, "Information Combining," *Found. Trends Commun. Inf. Theory*, vol. 3, no. 3, pp. 227–330, 2006.
- [23] Y. Ugur, I. E. Aguerri, and A. Zaidi, "Vector Gaussian CEO Problem Under Logarithmic Loss and Applications," *IEEE Trans. Inf. Theory*, vol. 66, no. 7, pp. 4183–4202, 2020.
- [24] M. Vera, L. Rey Vega, and P. Piantanida, "Collaborative Information Bottleneck," *IEEE Trans. Inf. Theory*, vol. 65, no. 2, pp. 787–815, 2019.
- [25] Q. Yang, P. Piantanida, and D. Gündüz, "The Multi-layer Information Bottleneck Problem," in *Proc. IEEE Inf. Theory Workshop (ITW)*, Nov. 2017, pp. 404–408.
- [26] A. Zaidi, I. E. Aguerri, and S. S. (Shitz), "On the Information Bottleneck Problems: Models, Connections, Applications and Information Theoretic Views," *Entropy*, vol. 22, no. 2, p. 151, 2020.
- [27] N. Tishby and N. Zaslavsky, "Deep Learning and the Information Bottleneck Principle," in *Proc. Inf. Theory Workshop (ITW)*, Apr. 2015, pp. 1–5.
- [28] A. Alemi, I. Fischer, J. Dillon, and K. Murphy, "Deep Variational Information Bottleneck," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2017.
- [29] P. Farajiparvar, A. Beirami, and M. Nokleby, "Information Bottleneck Methods for Distributed Learning," in *Proc. 56th Annu. Allerton Conf. Commun., Control Comput.*, 2018, pp. 24–31.
- [30] R. A. Amjad and B. C. Geiger, "Learning Representations for Neural Network-Based Classification Using the Information Bottleneck Principle," *IEEE Trans. Pattern Anal.*, vol. 42, no. 9, pp. 2225–2239, 2020.
- [31] A. M. Saxe, Y. Bansal, J. Dapello, M. Advani, A. Kolchinsky, B. D. Tracey, and D. D. Cox, "On the information bottleneck theory of deep learning," *J. Stat. Mech. Theory Exp.*, vol. 2019, no. 12, pp. 1–34, Dec. 2019.
- [32] Z. Goldfeld and Y. Polyanskiy, "The Information Bottleneck Problem and its Applications in Machine Learning," *IEEE J. Sel. Areas Inf. Theory*, vol. 1, no. 1, pp. 19–38, 2020.
- [33] S. Asoodeh and F. P. Calmon, "Bottleneck Problems: An Information and Estimation-Theoretic View," *Entropy*, vol. 22, no. 11, p. 1325, 2020.
- [34] P. Harremoës and N. Tishby, "The Information Bottleneck Revisited or How to Choose a Good Distortion Measure," in *Proc. 2007 IEEE Int. Symp. Inf. Theory*, Jun. 2007, pp. 566–570.
- [35] I. Sutskever, S. Shamaï, and J. Ziv, "Extremes of information combining," *IEEE Trans. Inf. Theory*, vol. 51, no. 4, pp. 1313–1325, Apr. 2005.
- [36] T. Richardson and R. Urbanke, *Modern Coding Theory*. Cambridge, U.K.: Cambridge Univ. Press, 2008.
- [37] F. G. Mehler, "Ueber die Entwicklung einer Function von beliebig vielen Variablen nach Laplaceschen Functionen höherer Ordnung," *J. Reine Angew. Math.*, vol. 66, pp. 161–176, 1866.
- [38] H. O. Lancaster, "The Structure of Bivariate Distributions," *Ann. Math. Statist.*, vol. 29, no. 3, pp. 719–736, 1958.
- [39] H. S. Witsenhausen and A. D. Wyner, "A Conditional Entropy Bound for a Pair of Discrete Random Variables," *IEEE Trans. Inf. Theory*, vol. 21, no. 5, pp. 493–501, Sep. 1975.
- [40] R. Gilad-Bachrach, A. Navot, and N. Tishby, "An information theoretic tradeoff between complexity and accuracy," in *Proc. 16th Conf. Comput. Theory (COLT)*, ser. Lecture Notes in Computer Science, B. Schölkopf and M. K. Warmuth, Eds., vol. 2777. Springer, 2003, pp. 595–609.
- [41] R. O'Donnell, *Analysis of Boolean Functions*, 1st ed. New York, NY, USA: Cambridge Univ. Press, Jun. 2014.
- [42] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. Hoboken, NJ, USA: Wiley, 2006.