

# THE DUAL INFORMATION BOTTLENECK

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The Information-Bottleneck (IB) framework suggests a general characterization of optimal representations in learning, and deep learning in particular. It is based on the optimal trade off between the representation complexity and accuracy, both of which are quantified by mutual information. The problem is solved by alternating projections between the *encoder* and *decoder* of the representation, which can be performed locally at each representation level. The framework, however, has practical drawbacks, in that mutual information is notoriously difficult to handle at high dimension, and only has closed form solutions in special cases. Further, because it aims to extract representations which are minimal sufficient statistics of the data with respect to the desired label, it does not necessarily optimize the actual prediction of unseen labels. Here we present a formal dual problem to the IB which has several interesting properties. By switching the order in the KL-divergence between the representation decoder and data, the optimal decoder becomes the geometric rather than the arithmetic mean of the input points. While providing a good approximation to the original IB, it also preserves the form of exponential families, and optimizes the mutual information on the predicted label rather than the desired one. We also analyze the critical points of the dual IB and discuss their importance for the quality of this approach.

## 1 INTRODUCTION

### 1.1 THE INFORMATION BOTTLENECK METHOD

The Information Bottleneck (IB) method (Tishby et al., 1999), is an information-theoretic framework for describing efficient representations of an “input” random variable  $X$  (input patterns), for predicting an “output” variable  $Y$  (desired label). In this setting the joint distribution of  $X$  and  $Y$ ,  $p(x, y)$  defines the problem, or rule, and the training data are a finite sample from this distribution. In general, we assume that  $p(y | x)$  is strictly stochastic, and hence bounded away from  $\{0, 1\}$ <sup>1</sup>. The representation variable  $\hat{X}$  is in general a stochastic function of  $X$  which forms a Markov chain  $Y \rightarrow X \rightarrow \hat{X}$ , and only depends on  $Y$  through the input  $X$ . We call the map  $p(\hat{x} | x)$  the *encoder* of the representation and denote by  $p(y | \hat{x})$  the *Bayes optimal decoder* for this representation; i.e., the best possible prediction of the *desired label*  $Y$  from the representation  $\hat{X}$ .

The IB trade off between the encoder and decoder mutual information values is defined by the minimization of the Lagrangian:

$$\mathcal{F}[p_\beta(\hat{x} | x); p_\beta(\hat{x}); p_\beta(y | \hat{x})] = I(X; \hat{X}) - \beta I(Y; \hat{X}), \quad (1)$$

independently over the convex sets of the normalized distributions,  $\{p_\beta(\hat{x} | x)\}$ ,  $\{p_\beta(\hat{x})\}$  and  $\{p_\beta(y | \hat{x})\}$ , given a positive Lagrange multiplier  $\beta$ . As shown in (Tishby et al., 1999; Shamir et al., 2010), this is a natural generalization of the classical concept of *Minimal Sufficient Statistics* (Cover & Thomas, 2006), where the estimated parameter is replaced by the output variable  $Y$  and *exact* statistical sufficiency is characterized by the mutual information equality:  $I(\hat{X}; Y) = I(X; Y)$ . The minimality of the statistics is captured by the minimization of  $I(X; \hat{X})$ , due to the Data Processing Inequality (DPI). However, non-trivial minimal sufficient statistics only exist for very special parametric distributions known as exponential families (Brown, 1986). Thus in general, the IB relaxes

<sup>1</sup>This may seem like a limitation of the method, but it can be shown that in the deterministic limit the IB is well defined. This stochastic assumption simplifies the analysis but does not pose any real limitation.

the minimal sufficiency problem to a continuous family of representations  $\hat{X}$  which are characterized by the trade off between compression,  $I(X; \hat{X}) \equiv I_X$ , and accuracy,  $I(Y; \hat{X}) \equiv I_Y$ , along a convex line in the *Information-Plane* ( $I_Y$  vs.  $I_X$ ). When the rule  $p(x, y)$  is strictly stochastic, the convex optimal line is smooth and each point along the line is uniquely characterized by the value of  $\beta$ . We can then consider the optimal representations  $\hat{x} = \hat{x}(\beta)$  as encoder-decoder pairs:  $(p_\beta(x|\hat{x}), p_\beta(y|\hat{x}))^2$  - a point in the continuous manifold defined by the Cartesian product of these distribution simplexes. We also consider a small variation of these representations,  $\delta\hat{x}$ , as an infinitesimal change in this (encoder-decoder) continuous manifold (not necessarily on the optimal line(s)).

## 1.2 IB AND RATE-DISTORTION THEORY

The IB optimization trade off can be considered as a generalized rate-distortion problem (Cover & Thomas, 2006) with the distortion function between a data point,  $x$  and a representation point  $\hat{x}$  taken as the KL-divergence between their predictions of the desired label  $y$ :

$$d_{\text{IB}}(x, \hat{x}) = D[p(y|x) || p_\beta(y|\hat{x})] = \sum_y p(y|x) \log \frac{p(y|x)}{p_\beta(y|\hat{x})}. \quad (2)$$

The expected distortion  $\langle d_{\text{IB}}(x, \hat{x}) \rangle_{p_\beta(x, \hat{x})}$  for the optimal decoder is simply the label-information loss:  $I(X; Y) - I(\hat{X}; Y)$ , using the Markov chain condition. Thus minimizing the expected IB distortion is equivalent to maximizing  $I(\hat{X}; Y)$ , or minimizing equation 1. Minimizing this distortion is equivalent to minimizing the log-loss or the cross-entropy loss, as done in most deep learning applications, and it upper-bounds other loss functions such as the  $\mathcal{L}_1$ -loss (due to the Pinsker inequality, or (Painsky & Wornell, 2018)). The Pinsker inequality shows that the relative entropy gives a symmetric upper bound to the  $\mathcal{L}_1$ -loss, or variation distance,  $D[p||q] \geq \frac{1}{2 \log 2} \|p - q\|_1^2$ .

## 1.3 THE IB EQUATIONS

For discrete  $X$  and  $Y$ , a necessary condition for the IB (local) minimization is given by the three self-consistent equations for the optimal encoder-decoder pairs, known as *the IB equations*:

$$\begin{cases} (i) & p_\beta(\hat{x}|x) = \frac{p(\hat{x})}{Z(x; \beta)} e^{-\beta D[p(y|x) || p_\beta(y|\hat{x})]} \\ (ii) & p_\beta(\hat{x}) = \sum_x p_\beta(\hat{x}|x) p(x) \\ (iii) & p_\beta(y|\hat{x}) = \sum_x p(y|x) p_\beta(x|\hat{x}) \end{cases}. \quad (3)$$

Iterating these equations is a generalized, Blahut-Arimoto, alternating projection algorithm (Tusnady & Csiszar, 1984; Cover & Thomas, 2006) and it converges to a stationary point of the Lagrangian, equation 1 (Tishby et al., 1999). Notice that the minimizing decoder, (equation 3-(iii)), is precisely the *Bayes optimal decoder* for the representation  $\hat{x}(\beta)$ , given the Markov chain conditions.

## 1.4 CRITICAL POINTS AND CRITICAL SLOWING-DOWN

One of the most interesting aspects of the IB equations is the existence of critical points along the optimal line of solutions in the *Information-Plane* (i.e. the information curve). At these points the representations change topology and cardinality (number of clusters) (Zaslavsky & Tishby, 2019; Parker et al., 2003) and they form the skeleton of the information curve and representation space. Under the strict stochastic assumption, the information-curve is a smooth function of the Lagrange multiplier  $\beta$ , but its derivative may not be smooth. Critical points are bifurcations of the solutions, which are values of  $\beta$  for which two different solutions (representations) co-exist. To identify such points we perform a perturbation analysis of the IB equations, as in (Zaslavsky & Tishby, 2019). Taking a small perturbation of the representation, denoted for brevity by  $\delta\hat{x}$ , the changes in the log encoder and log decoder that satisfy equation 3 for a given  $\beta$  can be determined through the nonlinear eigenvalues problems:

$$[I - \beta C_{xx'}^{\text{IB}}(\hat{x}, \beta)] \frac{\partial \log p_\beta(x'| \hat{x})}{\partial \hat{x}} = 0, \quad [I - \beta C_{yy'}^{\text{IB}}(\hat{x}, \beta)] \frac{\partial \log p_\beta(y' | \hat{x})}{\partial \hat{x}} = 0, \quad (4)$$

<sup>2</sup>Here we use the *inverse encoder*, which is in the fixed dimension simplex of distributions over  $X$ .

with the two square matrices  $C$  defined by:

$$C_{xx'}^{\text{IB}}(\hat{x}, \beta) = \sum_y p(y | x) \frac{p_\beta(x' | \hat{x})}{p_\beta(y | \hat{x})} p(y | x'), \quad C_{yy'}^{\text{IB}}(\hat{x}, \beta) = \sum_x p(y | x) \frac{p_\beta(x | \hat{x})}{p_\beta(y | \hat{x})} p(y' | x). \quad (5)$$

As shown in (Zaslavsky & Tishby, 2019), these two matrices have the same eigenvalues and have non-trivial eigenvectors (i.e., different co-existing optimal representations) precisely at the critical values of  $\beta$ , the bifurcation points of the IB solution. At these points the cardinality of the representation  $\hat{X}$  (the number of “IB-clusters”) changes due to splits of clusters, resulting in topological phase transitions in the encoder. These critical points form the “skeleton” of the topology of the optimal representations. Between critical points the optimal representations change continuously (with  $\beta$ ).

The important computational consequence of critical points is known as *critical slowing down* (Tredicce et al., 2004). For binary  $Y$ , near a critical point the convergence time,  $\tau_\beta$ , of the iterations of equation 3 scales like:  $\tau_\beta \sim 1/(1 - \beta\lambda_2)$ , where  $\lambda_2$  is the second eigenvalue of either  $C_{yy'}^{\text{IB}}$  or  $C_{xx'}^{\text{IB}}$ . At criticality,  $\lambda_2(\hat{x}) = \beta^{-1}$  and the number of iterations diverges. This phenomenon dominates any local minimization of equation 3 which is based on alternate encoder-decoder optimization.

The discussed phenomenons of the IB are graphically demonstrated in Figure 1. The figure shows the bifurcation diagram, the non-trivial eigenvalues of  $C_{yy'}^{\text{IB}}$  and convergence time of the IB iterations, for a simple problem with 4 critical points.

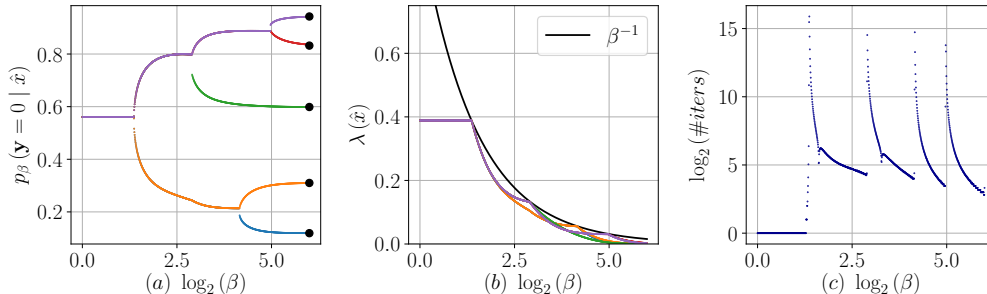


Figure 1: An example of IB solutions presenting the critical points and the *critical slowing down* of the alternating projection algorithm. (a) The algorithm’s solution to  $p_\beta(y = 0 | \hat{x})$  as a function of  $\beta$ . The black dots depict the input distribution  $p(y = 0 | x)$ . (b) The second eigenvalue of  $C_{yy'}(\hat{x}; \beta)$ ,  $\lambda_2(\hat{x})$ , along with  $\beta^{-1}$  as a function of  $\beta$ . (c) Convergence time of the algorithm as a function of  $\beta$ .

### 1.5 CONTRIBUTION OF THIS WORK

Supervised learning is generally separated into two phases: the training phase, where the internal representations are formed from the training data, and the prediction phase, where these representations are used to predict labels of new input patterns. In our Markov chain description we need to add another variable,  $\hat{Y}$ , the *predicted label* which obtains the same values as  $Y$ , but distributed differently:

$$\overbrace{Y \rightarrow X \rightarrow \hat{X}_\beta}^{\text{training}} \rightarrow \underbrace{\hat{Y}}_{\text{prediction}}. \quad (6)$$

The left-hand part of this chain describes the representation training, and the right-hand part is the Maximum Likelihood (ML) prediction using these representations (Slonim et al., 2006). So far the prediction variable  $\hat{Y}$  has not been part of the IB optimization problem. It has been implicitly assumed that the *Bayes optimal decoder*,  $p_\beta(y | \hat{x})$ , which minimizes the IB distortion at a given  $\beta$ , is also the best choice for making predictions of  $\hat{Y}$  from the representation  $\hat{X}_\beta$  through the right-hand Markov chain. That is, denoting  $p_\beta(\hat{y} | \hat{x}) \equiv p_\beta(y | \hat{x})$  the ML decoder is the mixture over the

internal representations:

$$p_\beta(\hat{y} | \mathbf{x}) \equiv \sum_{\hat{\mathbf{x}}} p_\beta(\hat{y} | \hat{\mathbf{x}}) p_\beta(\hat{\mathbf{x}} | \mathbf{x}). \quad (7)$$

This, however, is not necessarily optimal. For example due to finite sample, in which it can be very different from the one obtained from the full distribution (Shamir et al., 2010).

Here we show that by merely switching the order of the arguments in the KL-divergence of the IB distortion, namely:

$$d_{\text{dualIB}}(\mathbf{x}, \hat{\mathbf{x}}) = D[p_\beta(y | \hat{\mathbf{x}}) || p(y | \mathbf{x})] = \sum_y p_\beta(y | \hat{\mathbf{x}}) \log \frac{p_\beta(y | \hat{\mathbf{x}})}{p(y | \mathbf{x})}, \quad (8)$$

which in geometric terms is known as the *dual* distortion problem (Felice & Ay, 2019), we obtain several interesting results: (i) The decoder changes from the arithmetic to the geometric mean of the cluster data distributions; (ii) it preserves exponential form of the original data distribution, if one exists, for all values of  $\beta$ ; (iii) it preserves the low dimensional sufficient statistics of the data, making the scaling to large problems much easier; (iv) it (variationally) optimizes the predicted label information  $I(\mathbf{X}; \hat{\mathbf{Y}})$ , and with it the desired and predicted label information  $I(\mathbf{Y}; \hat{\mathbf{Y}})$ .

The remainder of the paper is organized as follows. We solve the dualIB problem in §2. We discuss the dualIB’s critical points §2.1 and provide a comparison to the original IB in §2.2. Next, in §3 we focus on the special case of data from exponential families. We conclude in §4 with further extensions and possible applications to deep learning.

## 2 THE DUALIB: MAXIMIZING THE PREDICTION INFORMATION

In the dualIB framework we consider the “full” learning Markov chain, equation 6. That is, given the “input” random variable  $\mathbf{X}$  (input patterns), the “output” random variable  $\mathbf{Y}$  (desired label) and the “representation”  $\hat{\mathbf{X}}$ , the ML optimal “predicted label”  $\hat{\mathbf{Y}}$  is given by the mixture distribution, equation 7.<sup>3</sup>

The dualIB optimization can be written as the following rate-distortion problem:

$$\mathcal{F}^*[p_\beta(\hat{\mathbf{x}} | \mathbf{x}); p_\beta(\hat{\mathbf{x}}); p_\beta(y | \hat{\mathbf{x}})] = I(\mathbf{X}; \hat{\mathbf{X}}) + \beta \langle d_{\text{dualIB}}(\mathbf{x}, \hat{\mathbf{x}}) \rangle_{p(\mathbf{x}, \hat{\mathbf{x}})}, \quad (9)$$

with the average distortion given in terms of mutual information on  $\hat{\mathbf{Y}}$ ,  $I(\hat{\mathbf{X}}; \hat{\mathbf{Y}})$  and  $I(\mathbf{X}; \hat{\mathbf{Y}})$ :

$$\begin{aligned} \langle d_{\text{dualIB}}(\mathbf{x}, \hat{\mathbf{x}}) \rangle_{p(\mathbf{x}, \hat{\mathbf{x}})} &= \underbrace{I(\hat{\mathbf{X}}; \hat{\mathbf{Y}}) - I(\mathbf{X}; \hat{\mathbf{Y}})}_{(a)} + \underbrace{\langle D[p_\beta(\hat{y} | \mathbf{x}) || p(\hat{y} | \mathbf{x})] \rangle_{p(\mathbf{x})}}_{(b)} \\ &\geq I(\hat{\mathbf{X}}; \hat{\mathbf{Y}}) - I(\mathbf{X}; \hat{\mathbf{Y}}). \end{aligned} \quad (10)$$

This is similar to the known IB relation:  $\langle d_{\text{IB}}(\mathbf{x}, \hat{\mathbf{x}}) \rangle_{p(\mathbf{x}, \hat{\mathbf{x}})} = I(\mathbf{Y}; \mathbf{X}) - I(\mathbf{Y}; \hat{\mathbf{X}})$  with an extra positive term (b).

Both terms, (a) and (b), vanish precisely when  $\hat{\mathbf{X}}$  is a sufficient statistic for  $\mathbf{X}$  with respect to  $\hat{\mathbf{Y}}$ , since we can then reverse the order of  $\mathbf{X}$  and  $\hat{\mathbf{X}}$  in the Markov chain (equation 6). This replaces the roles of  $\mathbf{Y}$  and  $\hat{\mathbf{Y}}$  as the variable for which  $\hat{\mathbf{X}}_\beta$  are approximately minimally sufficient statistics. In that sense the dualIB shifts the emphasis from the training phase to the prediction phase. This implies that minimizing the dualIB functional maximizes the mutual information between  $\mathbf{Y}$  and  $\hat{\mathbf{Y}}$ ,  $I(\mathbf{Y}; \hat{\mathbf{Y}})$ , as well as the mutual information  $I(\mathbf{X}; \hat{\mathbf{Y}})$ . This is illustrated in figure 4(a), (b).

The next theorem states the form of the solutions of the *Dual Information Bottleneck*:

**Theorem 1.** *The minima of equation 9, can be obtained by generalized Blahut-Arimoto iterations between the encoder and the decoder as in the original IB, with the following modifications: (i) Replace the distortion by its dual in the encoder update; (ii) Update the decoder by the encoder’s “geometric” mean of the data distributions  $p(y | \mathbf{x})$ .*

<sup>3</sup>We abuse the notation  $p(\hat{y} | \mathbf{x}) \equiv p(y | \mathbf{x})$ , when there is no  $\beta$  subscript.

The proof is given in §A.2.

The alternating projections between the encoder and decoder, which converge to a solution of the dualIB at a given value of the Lagrange multiplier  $\beta$ , are implemented by the following iterative algorithm<sup>4</sup>:

---

**Algorithm 1** dualIB iterative algorithm

---

- 1: **while**  $\left| p_{\beta}^{t+1}(y | \hat{x}) - p_{\beta}^t(y | \hat{x}) \right| > \epsilon$  **do**
  - 2:  $Z_{\hat{x}|\mathbf{x}}^{t+1}(\mathbf{x}; \beta) = \sum_{\hat{x}} p_{\beta}^t(\hat{x}) e^{-\beta D[p_{\beta}^t(y|\hat{x}) \| p(y|\mathbf{x})]}$
  - 3:  $p_{\beta}^{t+1}(\hat{x} | \mathbf{x}) = \frac{p_{\beta}^t(\hat{x})}{Z_{\hat{x}|\mathbf{x}}^{t+1}(\mathbf{x}; \beta)} e^{-\beta D[p_{\beta}^t(y|\hat{x}) \| p(y|\mathbf{x})]}$
  - 4:  $p_{\beta}^{t+1}(\hat{x}) = \sum_{\mathbf{x}} p_{\beta}^{t+1}(\hat{x} | \mathbf{x}) p(\mathbf{x})$
  - 5:  $Z_{y|\hat{x}}^{t+1}(\hat{x}; \beta) = \sum_{y} e^{\sum_{\mathbf{x}} p_{\beta}^{t+1}(\mathbf{x}|\hat{x}) \log p(y|\mathbf{x})}$
  - 6:  $p_{\beta}^{t+1}(y | \hat{x}) = \frac{1}{Z_{y|\hat{x}}^{t+1}(\hat{x}; \beta)} e^{\sum_{\mathbf{x}} p_{\beta}^{t+1}(\mathbf{x}|\hat{x}) \log p(y|\mathbf{x})}$
  - 7: **return**  $p_{\beta}(\hat{x}), p_{\beta}(\hat{x} | \mathbf{x}), p_{\beta}(y | \hat{x})$
- 

Here the  $Z$ 's are standard normalization factors (partition functions).

As in the IB, the **encoder update** (row 3) and **decoder update** (row 6) are the core of the algorithm. Figure 2 illustrates the properties of the dualIB solutions on the same small problem as in Figure 1.

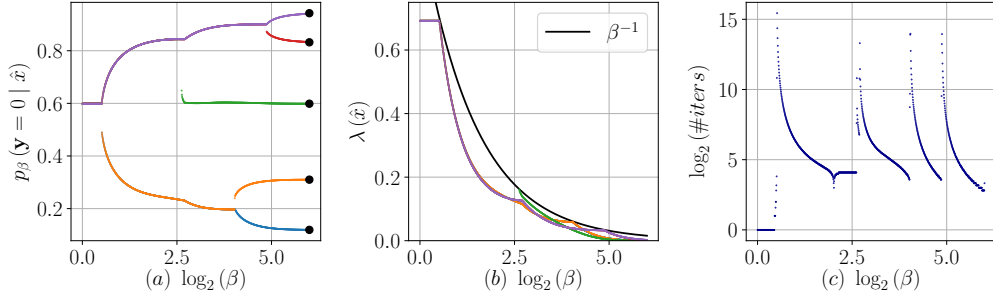


Figure 2: Same as Figure 1 for the dualIB solutions with the  $C_{yy'}^{\text{dualIB}}$  matrix. (a) The representations bifurcate at the critical values of  $\beta$ . (b) Critical points appear when an eigenvalue (color) crosses the  $\beta^{-1}$  (black) line. (c) Near these points there is *critical slowing down* and the numbers of iterations of Algorithm 1 diverge.

## 2.1 THE CRITICAL POINTS OF THE dualIB

As discussed in §1.4 the skeleton of the IB optimal bound (the information curve) is constituted by the critical points in which the topology (cardinality) of the representation changes. A similar stability analysis of the dualIB equations reveals similar conditions for the critical points.

**Theorem 2.** *The dualIB critical points are detected by non-trivial solutions of the nonlinear eigenvalue problem:*

$$[I - \beta C_{xx'}^{\text{dualIB}}(\hat{x}, \beta)] \frac{\partial \log p_{\beta}(x' | \hat{x})}{\partial \hat{x}} = 0, \quad [I - \beta C_{yy'}^{\text{dualIB}}(\hat{x}, \beta)] \frac{\partial \log p_{\beta}(y' | \hat{x})}{\partial \hat{x}} = 0, \quad (11)$$

<sup>4</sup>Unless stated otherwise, we use the convention  $\log = \ln \equiv \log_e$

with the matrices  $C^{\text{dualIB}}$  given by:

$$\begin{aligned} C_{xx'}^{\text{dualIB}}(\hat{x}; \beta) &= \sum_{y, \tilde{y}, \tilde{x}} p_{\beta}(y | \hat{x}) p_{\beta}(\tilde{x} | \hat{x}) \log \frac{p(y | x)}{p(y | \tilde{x})} \cdot p_{\beta}(x' | \hat{x}) p_{\beta}(\tilde{y} | \hat{x}) \log \frac{p(y | x')}{p(\tilde{y} | x')} \\ C_{yy'}^{\text{dualIB}}(\hat{x}; \beta) &= \sum_{x, \tilde{x}, \tilde{y}} p_{\beta}(x | \hat{x}) p_{\beta}(\tilde{y} | \hat{x}) \log \frac{p(y | x)}{p(\tilde{y} | x)} \cdot p_{\beta}(y' | \hat{x}) p_{\beta}(\tilde{x} | \hat{x}) \log \frac{p(y' | x)}{p(y' | \tilde{x})}. \end{aligned} \quad (12)$$

The proof to *theorem 2* is given in §A.3.

**Lemma 3.** *The matrices  $C_{xx'}^{\text{dualIB}}(\hat{x}; \beta)$ ,  $C_{yy'}^{\text{dualIB}}(\hat{x}; \beta)$  have the same eigenvalues  $\{\lambda_i\}$ , with  $\lambda_1(\hat{x}) = 0$ . With binary  $Y$ , the critical points are obtained at  $\lambda_2(\hat{x}) = \beta^{-1}$ .*

The proof of *lemma 3* given in section §A.3.1.

As in the IB, at the critical points,  $\beta_c^{\text{dualIB}}$ , the partial derivatives of the encoder and decoder with respect to  $\beta$ ,  $\partial \log p_{\beta}(x | \hat{x}) / \partial \beta$ ,  $\partial \log p_{\beta}(y | \hat{x}) / \partial \beta$ , have multiple (at least two) values. This results in discontinuities (cusps) in the encoder and decoder mutual information values as functions of  $\beta$  along the optimal line, with an undefined second derivative.

## 2.2 THE INFORMATION PLANE OF THE dualIB

The *Information-Plane*,  $I_X = I(\hat{X}; X)$  vs.  $I_Y = I(\hat{X}; Y)$ , is the standard visualization of the compression-prediction trade off of the IB. It can be defined for *any* encoder once the decoder is the Bayes optimal (equation 3-(iii)), for which the  $I_Y$  is the actual information of the representation on the desired label (Tishby et al., 1999).

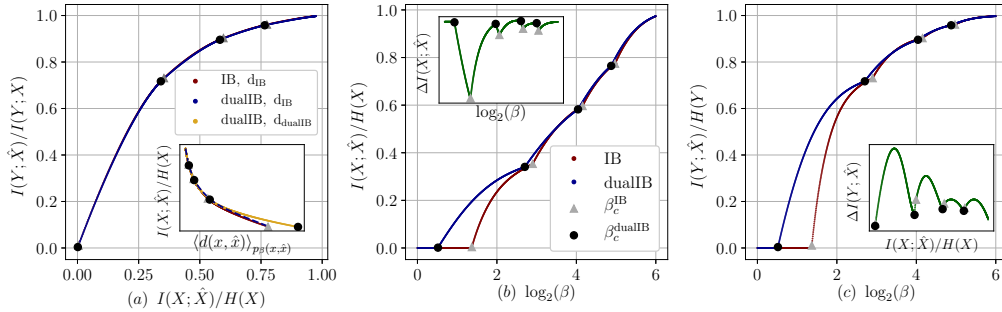


Figure 3: The IB’s and dualIB’s *Information Plane*. (a)  $I_Y$  vs.  $I_X$  for the two algorithms. The black dots are the dualIB critical points,  $\beta_c^{\text{dualIB}}$ , and the grey triangles are the IB critical points,  $\beta_c^{\text{IB}}$ . The corresponding distortion functions are shown in the inset. (b) The functions  $I_X^{\text{IB}}(\beta)$  and  $I_X^{\text{dualIB}}(\beta)$ . Both curves are monotonic and concave between the critical points. The inset indicates the relative difference between the curves, where the alternating order of the critical points is clearly observed. (c) Similarly,  $I_Y^{\text{IB}}(\beta)$  and  $I_Y^{\text{dualIB}}(\beta)$  are monotonic and piece-wise concave. The relative discrepancy between the information curves is clearly minimized at the dualIB critical points (inset). The functions approach each other for large  $\beta$ .

Comparing the dualIB information curve to the IB curve shows the quality of this approximation. Figure 3(a) depicts this comparison. While we know that  $I_Y^{\text{IB}}(\beta)$  is always higher, the two curves are almost indistinguishable. To better understand the relationship between these two curves, we look at the values of  $I_X$  and  $I_Y$  as functions of the corresponding  $\beta$  (Figure 3(b),(c)). The important role of the critical points is revealed as the corresponding cusps along these curves. As we argue below, the IB information values are strictly below those of the dualIB, but the distance between them is minimized precisely at the dual critical points.

**Lemma 4.**  *$I_X(\beta)$  and  $I_Y(\beta)$ , along the optimal lines, are non-decreasing piece-wise concave, functions of  $\beta$ . When their second derivative (with respect to  $\beta$ ) is defined, it is strictly negative.*

**Lemma 5.** For any sub-optimal information curve  $(I_X, I_Y)$ ,  $I_X^{\text{IB}}(\beta) \leq I_X(\beta)$  and  $I_Y^{\text{IB}}(\beta) \leq I_Y(\beta)$ , for all values of  $\beta$ .

Proofs of the above are given in §A.4.

The information plane properties are summarized by the following theorem.

**Theorem 6.** (i) The critical points of the two algorithms alternate: for each critical point,  $\beta_c^{\text{dualIB}} \leq \beta_c^{\text{IB}}$ . (ii) The distance between the two information curves is minimized precisely at the dualIB critical points  $\beta_c^{\text{dualIB}}$ . (iii) The two curves approach each other as  $\beta \rightarrow \infty$ .

*Proof.* The proof follows from lemmas 4 and 5, together with the critical points analysis above, and is only sketched here. As the encoder and decoder at the critical points,  $\beta_c^{\text{IB}}$  and  $\beta_c^{\text{dualIB}}$ , have different left and right derivatives, they form cusps in the curves of the mutual information  $(I_X$  and  $I_Y)$  as functions of  $\beta$ . These cusps can only be consistent with the optimality of the IB curves if  $\beta_c^{\text{dualIB}} < \beta_c^{\text{IB}}$  (this is true for any sub-optimal distortion), otherwise the curves intersect.

Moreover, at the dualIB critical points, the distance between the curves is minimized due to the strict concavity of the functions segments between the critical points. As the critical points imply discontinuity in the derivative, this results in a "jump" in the information values. Therefore, at any  $\beta_c^{\text{dualIB}}$  the distance between the curves has a (local) minimum. This is depicted in Figure 3, comparing  $I_X(\beta)$  and  $I_Y(\beta)$  and their differences for the two algorithms.

The two curves approach each other for large  $\beta$  since the two distortion functions become close in the low distortion limit (as long as  $p(y | x)$  is bounded away from 0).  $\square$

### 3 THE dualIB FOR EXPONENTIAL FAMILIES

Distributions of exponential families form the elegant theoretical core of parametric statistics and often emerge as maximum entropy (Jaynes, 1957) or stochastic equilibrium distributions, subject to observed constraints. They also form the class of parametric distributions for which exact, finite dimensional and additive, Minimal Sufficient Statistics exist (Kullback, 1959). One of the key properties of the dualIB is that it preserves the exponential form of the original data distribution,  $p(x, y)$  for all values of  $\beta$ .

Assuming that the data distribution is of the form:

$$p(y | x) = e^{-\sum_{r=0}^d \lambda^r(y) A_r(x)} = e^{-\lambda(y) \cdot \mathbf{A}(x) - \log Z_{y|x}(x)} = \prod_{r=0}^d e^{-\lambda^r(y) A_r(x)}, \quad (13)$$

where  $A_r(x)$  are  $d$  linearly independent functions of the input  $x$  and  $\lambda^r(y)$  are functions of the label  $y$ , or the parameters of this exponential family. The  $\lambda(y)$  can also be considered Lagrange multipliers associated with the constraints conditional expectations  $\langle A_r(x) \rangle_{p(x|y)}$  in entropy maximization.

The normalization factors,  $Z_{y|x}(x)$ , are written, for brevity, as  $\lambda_x^0 \equiv \log(\sum_y \prod_{r=1}^d e^{-\lambda^r(y) A_r(x)})$  with  $A_0(x) \equiv 1$ . We do not constrain the marginal  $p(x)$ .

The important fact about the exponential form is that all the mutual information,  $I(X; Y)$ , is fully captured by the  $d$  conditional expectations,  $\langle A_r(x) \rangle_{p(x|y)}$ , since these are the (minimal) sufficient statistics for the parameters. This means that all the relevant information (in the training sample) is captured solely by  $d$ -dimensional empirical expectations. This can lead to a huge reduction in computational complexity (from  $\dim(\mathbf{X})$  to  $d$ ).

We next show that for the dualIB, this dimension reduction is preserved or improved along the dual information curve, for all values of  $\beta$ .

**Theorem 7.** For data from an exponential family, equation 13, the optimal decoder of the dualIB, at a given  $\beta$ , is given by:

$$p_\beta(y | \hat{x}) = e^{-\sum_{r=1}^d \lambda^r(y) A_{r,\beta}(\hat{x}) - \lambda_\beta^0(\hat{x})}, \quad \lambda_\beta^0(\hat{x}) = \log\left(\sum_y e^{-\sum_{r=1}^d \lambda^r(y) A_{r,\beta}(\hat{x})}\right), \quad (14)$$

and the respective encoder:

$$p_\beta(\hat{x} | x) = \frac{p_\beta(\hat{x})e^{\beta\lambda_\beta^0(\hat{x})}}{Z_{\hat{x}|x}(x; \beta)} e^{-\beta \sum_{r=1}^d \lambda_\beta^r(\hat{x})[A_r(x) - A_{r,\beta}(\hat{x})]}, \quad (15)$$

with the constraints and multipliers expectations,

$$A_{r,\beta}(\hat{x}) \equiv \sum_x p_\beta(x | \hat{x}) A_r(x), \quad \lambda_\beta^r(\hat{x}) \equiv \sum_y p_\beta(y | \hat{x}) \lambda^r(y), \quad 1 \leq r \leq d. \quad (16)$$

The complete derivations for this section are given in §A.5.

This defines a simplified iterative algorithm to solve the dualExpIB problem, since we can replace the decoders' update rule at each iteration in the dualIB algorithm (Algorithm 1, row 6) with the simplified expression given in equation 14. The dualExpIB algorithm is more efficient because decoders' update amounts to estimating the  $d$  dimensional *constraints expectation*  $A_{r,\beta}(\hat{x})$ .

One would expect that as we decrease  $\beta$  the dimensionality of the *constraints expectation*,  $A_\beta(\hat{x})$ , should reduce as well. This is the most natural implication of the dimensionality reduction of  $\hat{x}$ . Figure 4(c) indeed shows that this is the case.  $A_\beta(\hat{x})$  follows the same bifurcation pattern as  $p_\beta(y | \hat{x})$ .

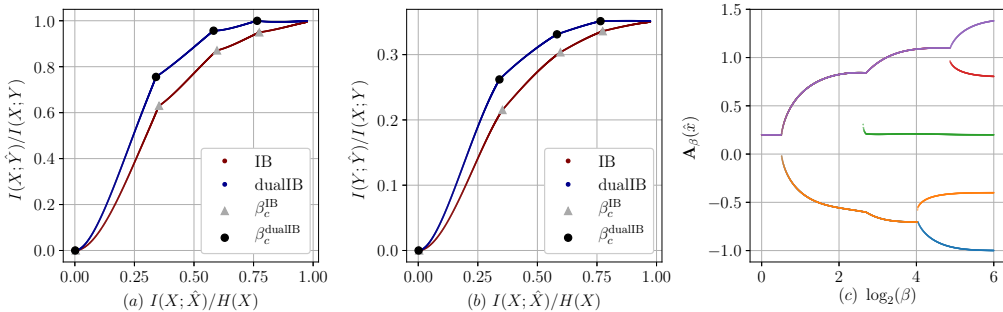


Figure 4: The two left plots present the *Information Planes* with respect to  $\hat{Y}$ , in both the dualIB appears above the IB. (a)  $I(X; \hat{Y})$  vs.  $I(X; \hat{X})$ . (b)  $I(Y; \hat{Y})$  vs.  $I(X; \hat{X})$ . (c) The *constraint expectation*,  $A_\beta(\hat{x})$  as a function of  $\beta$  for a problem with a single constraint ( $d = 1$ ). The relevant dimension of the representation decreases at the critical points.

## 4 CONCLUSION

We presented a new, dual formulation of the Information Bottleneck framework, based on switching the arguments in the original IB distortion function. This simple change has several interesting consequences: (i) it provides a good approximation to the original IB while keeping the algorithm in the low relevant dimension of the original data. This can significantly reduce the complexity of finding good IB representations; (ii) it optimizes the information between the representation and the *predicted label* rather than the desired label as in the original IB. This can improve the generalization error when trained on small samples since the predicted label is the one used in practice. (iii) It preserves the exponential form of data from exponential families, while reducing the dimensionality of the compressed representations. This important property was known to be satisfied by the Gaussian IB (Chechik et al., 2005) but not known for other distributions. The Gaussian case is self-dual in that sense. Generalizing this property to other exponential families was an open problem for many years. (iv) The exponential form of the optimal encoder-decoder pairs allows for the application to distributions with special symmetries, which can be naturally expressed in this form.

The topology of the reduced representations is determined by the critical points, where the cardinality of the representation changes. The critical points form the skeleton of the optimal solutions and most of the computation time is spent near the critical points. We analyzed the critical points of the dualIB and their analytic relations to those of the original IB.

This paves the way for completely new applications of the IB to representation learning, in particular deep learning of data with low internal dimensionality or special symmetries.



## REFERENCES

- Lawrence D. Brown. Fundamentals of statistical exponential families with applications in statistical decision theory. *Lecture Notes-Monograph Series*, 9:i–279, 1986. ISSN 07492170. URL <http://www.jstor.org/stable/4355554>.
- Gal Chechik, Amir Globerson, Naftali Tishby, and Yair Weiss. Information bottleneck for gaussian variables. *J. Mach. Learn. Res.*, 6:165–188, December 2005. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1046920.1046926>.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience, New York, NY, USA, 2006. ISBN 0471241954.
- Domenico Felice and Nihat Ay. Divergence functions in information geometry. In Frank Nielsen and Frédéric Barbaresco (eds.), *Geometric Science of Information - 4th International Conference, GSI 2019, Toulouse, France, August 27-29, 2019, Proceedings*, volume 11712 of *Lecture Notes in Computer Science*, pp. 433–442. Springer, 2019. ISBN 978-3-030-26979-1. doi: 10.1007/978-3-030-26980-7\_45. URL [https://doi.org/10.1007/978-3-030-26980-7\\_45](https://doi.org/10.1007/978-3-030-26980-7_45).
- Ran Gilad-bachrach, Amir Navot, and Naftali Tishby. An information theoretic tradeoff between complexity and accuracy. In *In Proceedings of the COLT*, pp. 595–609. Springer, 2003.
- E. T. Jaynes. Information theory and statistical mechanics. *Phys. Rev.*, 106:620–630, May 1957. doi: 10.1103/PhysRev.106.620. URL <https://link.aps.org/doi/10.1103/PhysRev.106.620>.
- S. Kullback. *Information Theory and Statistics*. Wiley, New York, 1959.
- Amichai Painsky and Gregory W. Wornell. Bregman Divergence Bounds and the Universality of the Logarithmic Loss. *arXiv e-prints*, art. arXiv:1810.07014, Oct 2018.
- Albert E. Parker, Tomáš Gedeon, and Alexander G. Dimitrov. Annealing and the rate distortion problem. In S. Becker, S. Thrun, and K. Obermayer (eds.), *Advances in Neural Information Processing Systems 15*, pp. 993–976. MIT Press, 2003. URL <http://papers.nips.cc/paper/2264-annealing-and-the-rate-distortion-problem.pdf>.
- Ohad Shamir, Sivan Sabato, and Naftali Tishby. Learning and generalization with the information bottleneck. *Theor. Comput. Sci.*, 411:2696–2711, 2010.
- Noam Slonim, Nir Friedman, and Naftali Tishby. Multivariate information bottleneck. *Neural Computation*, 18(8):1739–1789, 2006. doi: 10.1162/neco.2006.18.8.1739. URL <https://doi.org/10.1162/neco.2006.18.8.1739>. PMID: 16771652.
- Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method. In *Proc. of the 37-th Annual Allerton Conference on Communication, Control and Computing*, pp. 368–377, 1999.
- J. R. Tredicce, G. L. Lippi, Paul Mandel, B. Charasse, A. Chevalier, and B. Picqu. Critical slowing down at a bifurcation. *American Journal of Physics*, 72(6):799–809, 2004. doi: 10.1119/1.1688783. URL <https://doi.org/10.1119/1.1688783>.
- G. Tusnady and I. Csiszar. Information geometry and alternating minimization procedures. *Statistics & Decisions: Supplement Issues*, 1:205–237, 1984.
- Noga Zaslavsky and Naftali Tishby. Deterministic annealing and the evolution of optimal information bottleneck representations. *Preprint*, 2019.

## A APPENDIX

### A.1 DUALIB MATHEMATICAL FORMULATION

As presented in §1.5 the dualIB is solved with respect to the full Markov chain (equation 6) in which we introduce the new variable,  $\hat{Y}$ , the *predicted label*. Thus, in analogy to the IB we want to write the optimization problem in term of  $\hat{Y}$ .

Developing the expected distortion we find:

$$\begin{aligned} \langle d_{\text{dualIB}}(\mathbf{x}, \hat{\mathbf{x}}) \rangle_{p_{\beta}(\mathbf{x}, \hat{\mathbf{x}})} &= \sum_{\mathbf{x}, \hat{\mathbf{x}}} p_{\beta}(\mathbf{x}, \hat{\mathbf{x}}) \sum_{\hat{\mathbf{y}}} p_{\beta}(\hat{\mathbf{y}} | \hat{\mathbf{x}}) \log \frac{p_{\beta}(\hat{\mathbf{y}} | \hat{\mathbf{x}})}{p(\hat{\mathbf{y}} | \mathbf{x})} \\ &= \sum_{\hat{\mathbf{x}}, \hat{\mathbf{y}}} p_{\beta}(\hat{\mathbf{x}}) p_{\beta}(\hat{\mathbf{y}} | \hat{\mathbf{x}}) \log \frac{p_{\beta}(\hat{\mathbf{y}} | \hat{\mathbf{x}})}{p_{\beta}(\hat{\mathbf{y}})} - \sum_{\mathbf{x}, \hat{\mathbf{y}}} p(\mathbf{x}) p_{\beta}(\hat{\mathbf{y}} | \mathbf{x}) \log \frac{p_{\beta}(\hat{\mathbf{y}} | \mathbf{x})}{p_{\beta}(\hat{\mathbf{y}})} \\ &\quad + \sum_{\mathbf{x}, \hat{\mathbf{y}}} p(\mathbf{x}) p_{\beta}(\hat{\mathbf{y}} | \mathbf{x}) \log \frac{p_{\beta}(\hat{\mathbf{y}} | \mathbf{x})}{p(\hat{\mathbf{y}} | \mathbf{x})} \\ &= I(\hat{\mathbf{X}}; \hat{\mathbf{Y}}) - I(\mathbf{X}; \hat{\mathbf{Y}}) + \langle D[p_{\beta}(\hat{\mathbf{y}} | \mathbf{x}) \| p(\hat{\mathbf{y}} | \mathbf{x})] \rangle_{p(\mathbf{x})}. \end{aligned}$$

Allowing the dual optimization problem to be written as:

$$\mathcal{F}^* [p(\hat{\mathbf{x}} | \mathbf{x}); p(\hat{\mathbf{x}}); p(\mathbf{y} | \hat{\mathbf{x}})] = I(\mathbf{X}; \hat{\mathbf{X}}) - \beta \left\{ I(\mathbf{X}; \hat{\mathbf{Y}}) - I(\hat{\mathbf{X}}; \hat{\mathbf{Y}}) - \langle D[p_{\beta}(\hat{\mathbf{y}} | \mathbf{x}) \| p(\hat{\mathbf{y}} | \mathbf{x})] \rangle_{p(\mathbf{x})} \right\}.$$

### A.2 THE DUALIB SOLUTIONS

To prove *theorem 1* we want to obtain the normalized distributions minimizing the dualIB rate-distortion problem.

*Proof.* (i) Given that the problem is formulated as a rate-distortion problem the encoder's update rule must be the known minimizer of the distortion function (Cover & Thomas, 2006). Thus the IB encoder with the dual distortion is plugged in. (ii) For the decoder, by considering a small perturbation in the distortion  $d_{\text{dualIB}}(\mathbf{x}, \hat{\mathbf{x}})$ , with  $\alpha(\hat{\mathbf{x}})$  the normalization Lagrange multiplier, we obtain:

$$\begin{aligned} \delta d_{\text{dualIB}}(\mathbf{x}, \hat{\mathbf{x}}) &= \delta \left( \sum_{\mathbf{y}} p_{\beta}(\mathbf{y} | \hat{\mathbf{x}}) \log \frac{p_{\beta}(\mathbf{y} | \hat{\mathbf{x}})}{p(\mathbf{y} | \mathbf{x})} + \alpha(\hat{\mathbf{x}}) \left( \sum_{\mathbf{y}} p_{\beta}(\mathbf{y} | \hat{\mathbf{x}}) - 1 \right) \right) \\ \frac{\delta d_{\text{dualIB}}(\mathbf{x}, \hat{\mathbf{x}})}{\delta p_{\beta}(\mathbf{y} | \hat{\mathbf{x}})} &= \log \frac{p_{\beta}(\mathbf{y} | \hat{\mathbf{x}})}{p(\mathbf{y} | \mathbf{x})} + 1 + \alpha(\hat{\mathbf{x}}). \end{aligned}$$

Hence, minimizing the expected distortion becomes:

$$\begin{aligned} 0 &= \sum_{\mathbf{x}} p_{\beta}(\mathbf{x} | \hat{\mathbf{x}}) \left[ \log \frac{p_{\beta}(\mathbf{y} | \hat{\mathbf{x}})}{p(\mathbf{y} | \mathbf{x})} + 1 \right] + \alpha(\hat{\mathbf{x}}) \\ &= \log p_{\beta}(\mathbf{y} | \hat{\mathbf{x}}) - \sum_{\mathbf{x}} p_{\beta}(\mathbf{x} | \hat{\mathbf{x}}) \log p(\mathbf{y} | \mathbf{x}) + 1 + \alpha(\hat{\mathbf{x}}), \end{aligned}$$

which yields Algorithm 1, row 6.  $\square$

Considering the dualIB encoder-decoder, Algorithm 1, we find that  $\langle d_{\text{dualIB}}(\mathbf{x}, \hat{\mathbf{x}}) \rangle_{p_{\beta}(\mathbf{x}, \hat{\mathbf{x}})}$  reduces to the expectation of the decoder's log partition function:

$$\begin{aligned} \langle d_{\text{dualIB}}(\mathbf{x}, \hat{\mathbf{x}}) \rangle_{p_{\beta}(\mathbf{x}, \hat{\mathbf{x}})} &= \sum_{\mathbf{x}, \hat{\mathbf{x}}} p_{\beta}(\mathbf{x}, \hat{\mathbf{x}}) \sum_{\mathbf{y}} p_{\beta}(\mathbf{y} | \hat{\mathbf{x}}) \log \frac{p_{\beta}(\mathbf{y} | \hat{\mathbf{x}})}{p(\mathbf{y} | \mathbf{x})} \\ &= -\langle \log Z_{\mathbf{y} | \hat{\mathbf{x}}}(\hat{\mathbf{x}}; \beta) \rangle_{p_{\beta}(\hat{\mathbf{x}})} + \sum_{\hat{\mathbf{x}}, \mathbf{y}} p_{\beta}(\hat{\mathbf{x}}) \left[ \sum_{\mathbf{x}'} p_{\beta}(\mathbf{x}' | \hat{\mathbf{x}}) \log p(\mathbf{y} | \mathbf{x}') - \sum_{\mathbf{x}} p_{\beta}(\mathbf{x} | \hat{\mathbf{x}}) \log p(\mathbf{y} | \mathbf{x}) \right] \\ &= -\langle \log Z_{\mathbf{y} | \hat{\mathbf{x}}}(\hat{\mathbf{x}}; \beta) \rangle_{p_{\beta}(\hat{\mathbf{x}})}. \end{aligned}$$

### A.3 STABILITY ANALYSIS

Here we provide the detailed stability analysis allowing the definition of the matrices  $C_{xx'}^{\text{dualIB}}$ ,  $C_{yy'}^{\text{dualIB}}$  (equation 12) and which allows us to claim that they obey the same rules as the  $C$  matrices (equation 5) in equation 4. Considering a variation in  $\hat{x}$  we get:

$$\begin{aligned} \frac{\partial \log p_\beta(x | \hat{x})}{\partial \hat{x}} &= \beta \sum_y p_\beta(y | \hat{x}) \left( \log \frac{p(y | x)}{p_\beta(y | \hat{x})} - 1 \right) \frac{\partial \log p_\beta(y | \hat{x})}{\partial \hat{x}} \\ &= \beta \sum_y p_\beta(y | \hat{x}) \left[ \log p(y | x) - \sum_{\tilde{x}} p_\beta(\tilde{x} | \hat{x}) \log p(y | \tilde{x}) \right] \frac{\partial \log p_\beta(y | \hat{x})}{\partial \hat{x}} \\ &\quad + \beta \sum_y \log Z_{\mathbf{y}|\hat{x}}(\hat{x}; \beta) \frac{\partial p_\beta(y | \hat{x})}{\partial \hat{x}} \\ &= \beta \sum_{y, \tilde{x}} p_\beta(y | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y | x)}{p(y | \tilde{x})} \frac{\partial \log p_\beta(y | \hat{x})}{\partial \hat{x}}, \end{aligned} \quad (17)$$

$$\begin{aligned} \frac{\partial \log p_\beta(y | \hat{x})}{\partial \hat{x}} &= - \frac{1}{Z_{\mathbf{y}|\hat{x}}(\hat{x}; \beta)} \frac{\partial Z_{\mathbf{y}|\hat{x}}(\hat{x}; \beta)}{\partial \hat{x}} + \sum_x p_\beta(x | \hat{x}) \log p(y | x) \frac{\partial \log p_\beta(x | \hat{x})}{\partial \hat{x}} \\ &= - \sum_{\tilde{y}} p_\beta(\tilde{y} | \hat{x}) \sum_x p_\beta(x | \hat{x}) \log p(\tilde{y} | x) \frac{\partial \log p_\beta(x | \hat{x})}{\partial \hat{x}} \\ &\quad + \sum_x p_\beta(x | \hat{x}) \log p(y | x) \frac{\partial \log p_\beta(x | \hat{x})}{\partial \hat{x}} \\ &= \sum_{x, \tilde{y}} p_\beta(x | \hat{x}) p_\beta(\tilde{y} | \hat{x}) \log \frac{p(y | x)}{p(\tilde{y} | x)} \frac{\partial \log p_\beta(x | \hat{x})}{\partial \hat{x}}. \end{aligned} \quad (18)$$

Substituting equation 18 into equation 17 and vice versa one obtains:

$$\begin{aligned} \frac{\partial \log p_\beta(x | \hat{x})}{\partial \hat{x}} &= \beta \sum_{x', y, \tilde{y}, \tilde{x}} p_\beta(y | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y | x)}{p(y | \tilde{x})} \\ &\quad \cdot p_\beta(x' | \hat{x}) p_\beta(\tilde{y} | \hat{x}) \log \frac{p(y | x')}{p(\tilde{y} | x')} \frac{\partial \log p_\beta(x' | \hat{x})}{\partial \hat{x}} \\ \frac{\partial \log p_\beta(y | \hat{x})}{\partial \hat{x}} &= \beta \sum_{x, y', \tilde{x}, \tilde{y}} p_\beta(x | \hat{x}) p_\beta(\tilde{y} | \hat{x}) \log \frac{p(y | x)}{p(\tilde{y} | x)} \\ &\quad \cdot p_\beta(y' | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y' | x)}{p(y' | \tilde{x})} \frac{\partial \log p_\beta(y' | \hat{x})}{\partial \hat{x}}. \end{aligned}$$

We now define the  $C^{\text{dualIB}}$  matrices as follows:

$$\begin{aligned} C_{xx'}^{\text{dualIB}}(\hat{x}; \beta) &= \sum_{y, \tilde{y}, \tilde{x}} p_\beta(y | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y | x)}{p(y | \tilde{x})} \cdot p_\beta(x' | \hat{x}) p_\beta(\tilde{y} | \hat{x}) \log \frac{p(y | x')}{p(\tilde{y} | x')} \\ C_{yy'}^{\text{dualIB}}(\hat{x}; \beta) &= \sum_{x, \tilde{x}, \tilde{y}} p_\beta(x | \hat{x}) p_\beta(\tilde{y} | \hat{x}) \log \frac{p(y | x)}{p(\tilde{y} | x)} \cdot p_\beta(y' | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y' | x)}{p(y' | \tilde{x})}. \end{aligned}$$

Using the above definition we have an equivalence to equation 4 in the form of:

$$\left[ I - \beta C_{xx'}^{\text{dualIB}}(\hat{x}, \beta) \right] \frac{\partial \log p_\beta(x' | \hat{x})}{\partial \hat{x}} = 0, \quad \left[ I - \beta C_{yy'}^{\text{dualIB}}(\hat{x}, \beta) \right] \frac{\partial \log p_\beta(y' | \hat{x})}{\partial \hat{x}} = 0.$$

Note that for the binary case, the matrices may be simplified to:

$$\begin{aligned} C_{xx'}^{\text{dualIB}}(\hat{x}; \beta) &= \sum_{y, \tilde{x}} p_\beta(y | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y | x)}{p(y | \tilde{x})} \cdot p_\beta(x' | \hat{x}) (1 - p_\beta(y | \hat{x})) \log \frac{p(y | x')}{1 - p(y | x')} \\ C_{yy'}^{\text{dualIB}}(\hat{x}; \beta) &= \sum_{x, \tilde{x}} p_\beta(x | \hat{x}) (1 - p_\beta(y | \hat{x})) \log \frac{p(y | x)}{1 - p(y | x)} \cdot p_\beta(y' | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y' | x)}{p(y' | \tilde{x})}. \end{aligned}$$

Repeating the above steps considering a small change in  $\beta$  we see that the above obey the same, up to an additive constant, non-linear eigenvalue problem as the  $\hat{x}$  derivatives. Therefore, we conclude that they appear at the same critical points as the variation w.r.t  $\hat{x}$  (the solutions to the non-linear eigenvalue problem depend solely on the derivative coefficients).

### A.3.1 PROOF OF LEMMA 3

We show that the  $C^{\text{dualIB}}$  matrices share the same eigenvalues with  $\lambda_1(\hat{x}) = 0$ .

*Proof.* The matrices,  $C_{xx'}^{\text{dualIB}}(\hat{x}; \beta)$ ,  $C_{yy'}^{\text{dualIB}}(\hat{x}; \beta)$ , are given by:

$$C_{xx'}^{\text{dualIB}}(\hat{x}; \beta) = A_{xy}(\hat{x}; \beta)B_{yx'}(\hat{x}; \beta), \quad C_{yy'}^{\text{dualIB}}(\hat{x}; \beta) = B_{yx}(\hat{x}; \beta)A_{xy'}(\hat{x}; \beta),$$

with:

$$A_{xy}(\hat{x}; \beta) = p_\beta(y | \hat{x}) \sum_{\tilde{x}} p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y | \mathbf{x})}{p(y | \tilde{x})}, \quad B_{yx}(\hat{x}; \beta) = p_\beta(x | \hat{x}) \sum_{\tilde{y}} p_\beta(\tilde{y} | \hat{x}) \log \frac{p(y | \mathbf{x})}{p(\tilde{y} | \mathbf{x})}.$$

Given that the matrices are obtained by the multiplication of the same matrices, it follows that they have the same eigenvalues  $\{\lambda_i(\hat{x}; \beta)\}$ .

To prove that  $\lambda_1(\hat{x}; \beta) = 0$  we show that  $\det(C_{yy'}^{\text{dualIB}}) = 0$ . We present the exact calculation for a binary label  $Y \in \{y_0, y_1\}$  (the argument for general  $Y$  follows by encoding the label as a sequence of bits and discussing the first bit only, as a binary case.):

$$\begin{aligned} \det(C_{yy'}^{\text{dualIB}}(\hat{x}; \beta)) &= \sum_{x, \tilde{x}} p_\beta(x | \hat{x}) p_\beta(y_1 | \hat{x}) \log \frac{p(y_0 | \mathbf{x})}{p(y_1 | \mathbf{x})} \cdot p_\beta(y_0 | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y_0 | \mathbf{x})}{p(y_0 | \tilde{x})} \\ &\quad \cdot \sum_{x', \tilde{x}'} p_\beta(x' | \hat{x}) p_\beta(y_0 | \hat{x}) \log \frac{p(y_1 | \mathbf{x}')}{p(y_0 | \mathbf{x}')} \cdot p_\beta(y_1 | \hat{x}) p_\beta(\tilde{x}' | \hat{x}) \log \frac{p(y_1 | \mathbf{x}')}{p(y_1 | \tilde{x}')} \\ &\quad - \sum_{x, \tilde{x}} p_\beta(x | \hat{x}) p_\beta(y_0 | \hat{x}) \log \frac{p(y_1 | \mathbf{x})}{p(y_0 | \mathbf{x})} \cdot p_\beta(y_0 | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y_0 | \mathbf{x})}{p(y_0 | \tilde{x})} \\ &\quad \cdot \sum_{x', \tilde{x}'} p_\beta(x' | \hat{x}) p_\beta(y_1 | \hat{x}) \log \frac{p(y_0 | \mathbf{x}')}{p(y_1 | \mathbf{x}')} \cdot p_\beta(y_1 | \hat{x}) p_\beta(\tilde{x}' | \hat{x}) \log \frac{p(y_1 | \mathbf{x}')}{p(y_1 | \tilde{x}')} \\ &= \sum_{x, x', \tilde{x}, \tilde{x}'} p_\beta(x | \hat{x}) p_\beta(x' | \hat{x}) p_\beta^2(y_0 | \hat{x}) p_\beta^2(y_1 | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y_0 | \mathbf{x})}{p(y_0 | \tilde{x})} p_\beta(\tilde{x}' | \hat{x}) \log \frac{p(y_1 | \mathbf{x}')}{p(y_1 | \tilde{x}')} \\ &\quad \cdot \left[ \log \frac{p(y_0 | \mathbf{x})}{p(y_1 | \mathbf{x})} \log \frac{p(y_1 | \mathbf{x}')}{p(y_0 | \mathbf{x}')} - \log \frac{p(y_0 | \mathbf{x})}{p(y_1 | \mathbf{x})} \log \frac{p(y_1 | \mathbf{x}')}{p(y_0 | \mathbf{x}')} \right] = 0. \end{aligned}$$

Given that the determinant is 0 implies that  $\lambda_1(\hat{x}) = 0$ .  $\square$

For a binary problem we can describe the non-zero eigenvalue using  $\lambda_2(\hat{x}) = \text{Tr}(C_{yy'}^{\text{dualIB}}(\hat{x}; \beta))$ . That is:

$$\begin{aligned} \lambda_2(\hat{x}) &= \sum_{x, \tilde{x}} p_\beta(x | \hat{x}) p_\beta(y_1 | \hat{x}) \log \frac{p(y_0 | \mathbf{x})}{p(y_1 | \mathbf{x})} \cdot p_\beta(y_0 | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y_0 | \mathbf{x})}{p(y_0 | \tilde{x})} \\ &\quad + \sum_{x, \tilde{x}} p_\beta(x | \hat{x}) p_\beta(y_0 | \hat{x}) \log \frac{p(y_1 | \mathbf{x})}{p(y_0 | \mathbf{x})} \cdot p_\beta(y_1 | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y_1 | \mathbf{x})}{p(y_1 | \tilde{x})} \\ &= p_\beta(y_1 | \hat{x}) p_\beta(y_0 | \hat{x}) \sum_{x, \tilde{x}} p_\beta(x | \hat{x}) p_\beta(\tilde{x} | \hat{x}) \log \frac{p(y_0 | \mathbf{x})}{p(y_1 | \mathbf{x})} \left[ \log \frac{p(y_0 | \mathbf{x})}{p(y_0 | \tilde{x})} - \log \frac{p(y_1 | \mathbf{x})}{p(y_1 | \tilde{x})} \right]. \end{aligned}$$

#### A.4 INFORMATION PLANE ANALYSIS

We rely on known results for the rate-distortion problem and the information plane:

**Lemma 8.**  $I(X; \hat{X})$  is a non-increasing convex function of the distortion  $\langle d(x, \hat{x}) \rangle_{p_\beta(x, \hat{x})}$  with a slope of  $-\beta$ .

We emphasize that this is a general result of rate-distortion thus holds for the dualIB as well.

**Lemma 9.** For a fixed encoder  $p_\beta(\hat{x} | x)$  and the Bayes optimal decoder  $p_\beta(y | \hat{x})$ :

$$\langle d_{\text{IB}}(x, \hat{x}) \rangle_{p_\beta(x, \hat{x})} = I(X; Y) - I(\hat{X}; Y).$$

Thus, the information curve,  $I_Y$  vs.  $I_X$ , is a non-decreasing concave function with a positive slope,  $\beta^{-1}$ . The concavity implies that  $\beta$  increases along the curve.

(Cover & Thomas, 2006; Gilad-bachrach et al., 2003).

##### A.4.1 PROOF OF LEMMA 4

In the following section we provide a proof to *lemma 4*, for the IB and dualIB problems.

*Proof.* We want to analyze the behavior of  $I_X(\beta)$ ,  $I_Y(\beta)$ , that is the change in each term as a function of the corresponding  $\beta$ . From *lemma 9*, the concavity of the information curve, we can deduce that both are non-decreasing functions of  $\beta$ . As the two  $\beta$  derivatives are proportional it's enough to discuss the first one.

Next, we focus on their behavior between two critical points. That is, where the cardinality of  $\hat{X}$  is fixed (clusters are "static"). For "static" clusters, the  $\beta$  derivative of  $I_X$ , along the optimal line is given by:

$$\begin{aligned} \frac{\partial I(X; \hat{X})}{\partial \beta} &= -\frac{\partial}{\partial \beta} \left[ \sum_{x, \hat{x}} p_\beta(x, \hat{x}) (\log Z_{\hat{x}|x}(x; \beta) + \beta d(x, \hat{x})) \right] \\ &= -\beta \left\langle d(x, \hat{x}) \frac{\partial \log p_\beta(\hat{x} | x)}{\partial \beta} \right\rangle_{p_\beta(x, \hat{x})} \\ &\approx \beta \left\langle d(x, \hat{x}) \left[ \frac{\partial \log Z_{\hat{x}|x}(x; \beta)}{\partial \beta} + d(x, \hat{x}) \right] \right\rangle_{p_\beta(x, \hat{x})} \\ &\approx \beta \left\langle \underbrace{\langle d^2(x, \hat{x}) \rangle_{p_\beta(\hat{x}|x)}}_{\text{Var}(d(x))} - \underbrace{\langle d(x, \hat{x}) \rangle_{p_\beta(\hat{x}|x)}^2}_{p(x)} \right\rangle. \end{aligned}$$

This first of all reassures that the function is non-decreasing as  $\text{Var}(d(x)) \geq 0$ .

The piece-wise concavity follows from the fact that when the number of clusters is fixed (between the critical points) - increasing  $\beta$  decreases the clusters conditional entropy  $H(\hat{X} | X)$ , as the encoder becomes more deterministic. The mutual information is bounded by  $H(\hat{X})$  and it's  $\beta$  derivative decreases. Further, between the critical points there are no sign changes in the second  $\beta$  derivative.  $\square$

##### A.4.2 PROOF OF LEMMA 5

*Proof.* The information curve has a positive slope,  $\beta^{-1}$ , with  $\beta$  increasing along it, *lemma 9*. That is, given a value of  $\beta$ , there exists a pair  $I_Y^{\text{IB}}(\beta)$ ,  $I_X^{\text{IB}}(\beta)$  such that  $\partial I_Y^{\text{IB}}(\beta) / \partial I_X^{\text{IB}}(\beta) = \beta^{-1}$ . Now, consider a sub-optimal information curve,  $I_Y^*$ ,  $I_X^*$ . There exist values  $\beta'$ ,  $\beta''$  such that:

$$I_X^*(\beta') = I_X^{\text{IB}}(\beta), \quad I_Y^*(\beta'') = I_Y^{\text{IB}}(\beta).$$

The optimality of the IB implies that sub-optimal curves lie below it; i.e, the IB slope is steeper:

$$\beta^{-1} > \beta'^{-1}, \quad \beta^{-1} > \beta''^{-1}.$$

Thus, given that  $\beta$  increases along the information curve, it holds that:

$$I_X^*(\beta) > I_X^*(\beta') = I_X^{\text{IB}}(\beta), \quad I_Y^*(\beta) > I_Y^*(\beta'') = I_Y^{\text{IB}}(\beta).$$

□

### A.5 DERIVATION OF THE DUALEXPIB

We provide elaborate derivations to *theorem 7*; that is, we obtain the dualIB optimal encoder-decoder under the exponential assumption over the data. We use the notations defined in §3.

- The *decoder*, equation 14.  
Substituting the exponential assumption into the dualIB log-decoder yields:

$$\begin{aligned} \log p_\beta(y | \hat{x}) &= \sum_x p_\beta(x | \hat{x}) \log p(y | x) - \log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x}; \beta) \\ &= - \sum_x \sum_{r=0}^d p_\beta(x | \hat{x}) \lambda^r(y) A_r(x) - \log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x}; \beta) \\ &= - \sum_{r=1}^d \lambda^r(y) A_{r,\beta}(\hat{x}) - \langle \lambda_{\mathbf{x}}^0 \rangle_{p_\beta(x|\hat{x})} - \log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x}; \beta). \end{aligned}$$

Taking a closer look at the normalization term:

$$\begin{aligned} Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x}; \beta) &= \sum_y e^{\sum_x p_\beta(x|\hat{x}) \log p(y|x)} = e^{-\langle \lambda_{\mathbf{x}}^0 \rangle_{p_\beta(x|\hat{x})}} \sum_y e^{-\sum_{r=1}^d \lambda^r(y) A_{r,\beta}(\hat{x})} \\ \log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x}; \beta) &= -\langle \lambda_{\mathbf{x}}^0 \rangle_{p_\beta(x|\hat{x})} + \log \left( \sum_y e^{-\sum_{r=1}^d \lambda^r(y) A_{r,\beta}(\hat{x})} \right). \end{aligned}$$

From which it follows that  $\lambda_\beta^0(\hat{x})$  is given by:

$$\lambda_\beta^0(\hat{x}) = \log \left( \sum_y e^{-\sum_{r=1}^d \lambda^r(y) A_{r,\beta}(\hat{x})} \right),$$

and we can conclude that the dualExpIB decoder takes the form:

$$\log p_\beta(y | \hat{x}) = - \sum_{r=0}^d \lambda^r(y) A_{r,\beta}(\hat{x}).$$

- The *encoder*, equation 15.  
The core of the encoder is the dual distortion function which may now be written as:

$$\begin{aligned} d_{\text{dualIB}}(\mathbf{x}, \hat{x}) &= \sum_y p_\beta(y | \hat{x}) \log \frac{p_\beta(y | \hat{x})}{p(y | \mathbf{x})} \\ &= \sum_y p_\beta(y | \hat{x}) \left[ (\lambda_{\mathbf{x}}^0 - \lambda_\beta^0(\hat{x})) + \sum_{r=1}^d \lambda^r(y) (A_r(\mathbf{x}) - A_{r,\beta}(\hat{x})) \right] \\ &= \lambda_{\mathbf{x}}^0 - \lambda_\beta^0(\hat{x}) + \sum_{r=1}^d \lambda_\beta^r(\hat{x}) (A_r(\mathbf{x}) - A_{r,\beta}(\hat{x})), \end{aligned}$$

substituting this into the encoder's definition we obtain:

$$\begin{aligned} p_\beta(\hat{x} | \mathbf{x}) &= \frac{p_\beta(\hat{x})}{Z_{\hat{\mathbf{x}}|\mathbf{x}}(\mathbf{x}; \beta)} e^{-\beta [\lambda_{\mathbf{x}}^0 - \lambda_\beta^0(\hat{x}) + \sum_{r=1}^d \lambda_\beta^r(\hat{x}) (A_r(\mathbf{x}) - A_{r,\beta}(\hat{x}))]} \\ &= \frac{p_\beta(\hat{x}) e^{\beta \lambda_\beta^0(\hat{x})}}{Z_{\hat{\mathbf{x}}|\mathbf{x}}(\mathbf{x}; \beta)} e^{-\beta \sum_{r=1}^d \lambda_\beta^r(\hat{x}) (A_r(\mathbf{x}) - A_{r,\beta}(\hat{x}))}. \end{aligned}$$