

---

# Probabilistic Dependency Graph Inconsistency, the Natural Loss: Choose your Model, not your Loss Function

---

Anonymous Author  
Anonymous Institution

maybe more  
slower than I'd prefer for abstract.

## Abstract

In a world **blessed** with a great diversity of loss functions, we argue that that choice between them is not a matter of taste (or pragmatics), but of model. Probabilistic dependency graphs (PDGs) are probabilistic models that come equipped with a measure of “inconsistency”. We prove that many standard loss functions arise as the inconsistency of a natural PDG describing the appropriate scenario, and use the same approach to justify a well-known connection between regularizers and priors. We also show that the inconsistency of PDGs captures a large class of statistical divergences, and detail benefits of thinking of them in this way, including an intuitive visual language for deriving inequalities between them. In variational inference, we find that the ELBO, a difficult-to-motivate objective for latent variable models, and variants of it, arise for free out of uncontroversial modeling assumptions—as do simple graphical proofs of their corresponding bounds. Finally, we observe that inconsistency becomes the log partition function (free energy) in the setting where PDGs are factor graphs.

## 1 INTRODUCTION

Many tasks in artificial intelligence have been fruitfully cast as optimization problems, but often the choice of objective is not unique. For instance, a key component of a machine learning system is a loss function which the system minimizes. A wide variety of losses are used in practice. Each implicitly represents different values and results in different behavior, so the

choice of loss can be quite important (Wang et al. 2020; Jadon 2020). Yet, because it’s unclear how to choose a “good” loss function, the choice is usually made by empirics, tradition, and an instinctive calculus acquired through the practice—not explicitly laying out beliefs. Furthermore, there is something to be gained by fiddling with these loss functions: one can add regularization terms, to (dis)incentivize (un)desirable behavior. But the process of tinkering with the objective until it works is often unsatisfying. It can be a tedious game without clear rules or meaning, while results so obtained are arguably overfitted and hard to motivate.

By contrast, a choice of *model* admits more principled discussion, in part because models are testable: it makes sense to ask if a model is accurate. This observation motivates our proposal: instead of specifying a loss function directly, one articulates a situation that gives rise to it, in the (more interpretable) language of probabilistic beliefs and certainties. Concretely, we use the machinery of Probabilistic Dependency Graphs (PDGs), a particularly expressive class of graphical models that can incorporate arbitrary (even inconsistent) probabilistic information in a natural way and comes equipped with a well-motivated measure of inconsistency (Richardson and Halpern 2021).

The goal of this paper is to show that PDGs and their associated inconsistency measure can provide a “universal” model-based loss function. Towards this end, we show that many standard objective functions—cross entropy, square error, many statistical distances, the ELBO, regularizers, and the log partition function—arise naturally as inconsistencies of the appropriate underlying PDG. This is somewhat surprising, since PDGs were not designed with the goal of capturing loss functions at all. Specifying a loss function indirectly like this may be more restrictive, but it is also more intuitive (since it requires no technical familiarity with losses), and admits more epistemically grounded discussion.

For a particularly powerful demonstration, consider the variational autoencoder (VAE), an enormously successful class of generative model that has enabled

breakthroughs in image generation, semantic interpolation, and unsupervised feature learning (Kingma and Welling 2014). Structurally, a VAE for a space  $X$  consists of a (smaller) latent space  $Z$ , a prior distribution  $p(Z)$ , a decoder  $d(X|Z)$ , and an encoder  $e(Z|X)$ . A VAE is not considered a “graphical model” for two reasons. The first is that the encoder  $e(Z|X)$  has the same target variable as  $p(Z)$ , so something like a Bayesian Network cannot simultaneously incorporate them both (besides, they could be inconsistent with one another). The second reason: it is not a VAE’s structure, but rather its *loss function* that makes it tick. A VAE is typically trained by maximizing the “ELBO”, a somewhat difficult-to-motivate function of a sample  $x$ , borrowed from variational analysis. We show that  $-\text{ELBO}(x)$  is also precisely the inconsistency of a PDG containing the probabilistic information of the autoencoder ( $p, d$ , and  $e$ ) and  $x$ . We can form such a PDG precisely because PDGs allow for inconsistency. Thus, PDG semantics simultaneously legitimize the strange structure of the VAE, and also justify its loss function, which can be thought of as a property of the model itself (its inconsistency), rather than some mysterious construction borrowed from physics.

Representing objectives as model inconsistencies, in addition to providing a principled way of selecting an objective, also has beneficial pedagogical side effects, because of the relationships between the underlying models. For instance, we will be able to use the structure of the PDG to get simple, intuitive proofs of technical results, such as the variational inequalities that traditionally motivate the ELBO, and the monotonicity of Rényi divergence.

In the coming sections, we show in more detail how this concept of inconsistency, beyond simply providing a permissive and intuitive modeling framework, reduces exactly to many standard objectives used in machine learning and to measures of statistical distance. We demonstrate that this framework clarifies the relationships between them, by providing simple and clear derivations of otherwise opaque inequalities.

## 2 PRELIMINARIES

We generally use capital letters for variables, and lower case letters for their values. For variables  $X$  and  $Y$ , a conditional probability distribution (cpd)  $p$  on  $Y$  given  $X$ , written  $p(Y|X)$ , consists of a probability distribution on  $Y$  (denoted  $p(Y|X=x)$  or  $p(Y|x)$  for short), for each possible value  $x$  of  $X$ . If  $\mu$  is a probability on outcomes that determine  $X$  and  $Y$ , then  $\mu(X)$  denotes the marginal of  $\mu$  on  $X$ , and  $\mu(Y|X)$  denotes the conditional marginal of  $\mu$  on  $Y$  given  $X$ . Depending on which we find clearer in context, we write ei-

ther  $\mathbb{E}_\mu f$  or  $\mathbb{E}_{\omega \sim \mu} f(\omega)$  for expectation of  $f : \Omega \rightarrow \mathbb{R}$  over a distribution  $\mu$  with outcomes  $\Omega$ . We write  $D(\mu \parallel \nu) = \mathbb{E}_\mu \log \frac{\mu}{\nu}$  for the relative entropy (KL Divergence) of  $\nu$  with respect to  $\mu$ , and for finitely supported  $\mu$ , we write  $H(\mu) := \mathbb{E}_\mu \log \frac{1}{\mu}$  for the entropy of  $\mu$ ,  $H_\mu(X) := H(\mu(X))$  for the marginal entropy on a variable  $X$ , and  $H_\mu(Y|X) := \mathbb{E}_\mu \log 1/\mu(Y|X)$  for the conditional entropy of  $Y$  given  $X$ .

A *probabilistic dependency graph* (PDG) (Richardson and Halpern 2021), like a Bayesian Network (BN), is a directed graph with cpds attached to it. While this data is attached to the nodes of a BN, it is attached to the edges of a PDG. For instance, a PDG of shape  $X \rightarrow Y \leftarrow Z$  contains both a cpd  $p(Y|X)$  and (separately) a cpd  $q(Y|Z)$ , while a BN of the same shape has a single cpd  $\Pr(Y|X, Z)$  on  $Y$  given joint values of  $X$  and  $Z$ . The first interpretation is more expressive, and can encode joint dependence with an extra variable and pair of edges. We now restate the formal definition.

**Definition 1.** A Probabilistic Dependency Graph (PDG) is a tuple  $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ , where

- $\mathcal{N}$  is a set of nodes, corresponding to variables;
- $\mathcal{V}$  associates each node  $X \in \mathcal{N}$  with a set  $\mathcal{V}(X)$  of possible values that the variable  $X$  can take;
- $\mathcal{E}$  is a set of labeled edges  $\{X \xrightarrow{L} Y\}$ , each with a source  $X$  and target  $Y$  from  $\mathcal{N}$ ;
- $\mathbf{p}$  associates a cpd  $\mathbf{p}_L(Y|X)$  to each edge  $X \xrightarrow{L} Y \in \mathcal{E}$ ;
- $\boldsymbol{\alpha}$  associates to each edge  $X \xrightarrow{L} Y$  a non-negative number  $\alpha_L$  representing the modeler’s confidence in the functional dependence of  $Y$  on  $X$ ;
- $\boldsymbol{\beta}$  associates to each edge  $L$  a real number  $\beta_L$ , the modeler’s subjective confidence in the reliability of the cpd  $\mathbf{p}_L$ .  $\square$

We conflate a cpd’s symbol as the edge label, so we draw the PDG with one edge attached to the cpd  $f(Y|X)$  as  $\boxed{X} \xrightarrow{f} \boxed{Y}$ . Definition 1 is equivalent to one in which edge sources and targets are both sets of variables. This allows us to indicate joint dependence with multi-tailed arrows, joint distributions with multi-headed arrows, and unconditional distributions with nothing at the tail. For instance, we draw

$$p(Y|X, Z) \text{ as } \boxed{Z} \xrightarrow{p} \boxed{Y} \text{, and } q(A, B) \text{ as } \boxed{A} \xrightarrow{q} \boxed{B}.$$

Like other graphical models, PDGs have semantics in terms of joint distributions  $\mu$  over all variables. Most directly, a PDG  $\mathcal{M}$  determines two scoring functions on joint distributions  $\mu$ . For the purposes of this paper, the more important of the two is the *incompatibility* of  $\mu$  with respect to  $\mathcal{M}$ , which measures the quantitative

has been repeated a few times

if you add a diagram in, even in appendix, that would be really helpful

this definition is maybe a little intimidating



Essentially the only choices we’ve made in specifying the PDG of [Proposition 3](#) are the confidences. But  $\text{CrossEntropy}(\text{Pr}_{\mathcal{D}}, p)$  is the expected code length per sample from  $\text{Pr}_{\mathcal{D}}$ , when using codes optimized for the (incorrect) distribution  $p$ . So implicitly, a modeler using cross-entropy has already articulated a belief the data distribution is the “true one”. To get the same effect from a PDG, the modeler must make this belief explicit by placing infinite certainty in  $\text{Pr}_{\mathcal{D}}$ .

Now consider an orthogonal generalization of [Proposition 2](#), in which the sample  $x$  is only **a partial observation** of a joint model  $p(X, Z)$ .

**Proposition 4.** *If  $p(X, Z)$  is a joint distribution, the information content of the partial observation  $X = x$  is given by*

$$I_p[X=x] = \left\langle \left\langle \begin{array}{c} \text{Z} \xleftarrow{p} \text{X} \xleftarrow{x} \end{array} \right\rangle \right\rangle. \quad (2)$$

Intuitively, the inconsistency of the PDG on the right of (2) is localized to  $X$ , where the observation conflicts with the distribution; other variables don’t make a difference. The multi-sample partial-observation generalization also holds; see [Appendix B.3](#).

So far we have considered models of an unconditional distribution  $p(X)$ . Because they are unconditional, such models must describe how to generate a complete sample  $X$  without input, and so are called *generative*; the process of training them is called *unsupervised learning* ([Hastie, Tibshirani, and Friedman 2009](#)). In the (more common) *supervised* setting, we train *discriminative* models to predict  $Y$  from  $X$ , from labeled samples  $\{(x_i, y_i)\}_i$ . There, cross entropy loss is perhaps even more dominant—and it is essentially the inconsistency of a PDG consisting of the predictor  $h(Y|X)$  together with high-confidence data.

**Proposition 5** (Cross Entropy, Supervised). *The inconsistency of the PDG comprising a probabilistic predictor  $h(Y|X)$ , and a high-confidence empirical distribution  $\text{Pr}_{\mathcal{D}}$  of a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^m$  is equal to the cross-entropy loss (minus the empirical uncertainty in  $Y$  given  $X$ , a constant that depends only on  $\mathcal{D}$ ). That is,*

$$\left\langle \left\langle \begin{array}{c} \text{Pr}_{\mathcal{D}} \xrightarrow{(\infty)} \text{X} \xrightarrow{h} \text{Y} \end{array} \right\rangle \right\rangle = \frac{1}{m} \sum_{i=1}^m \log \frac{1}{h(y_i | x_i)} - H_{\text{Pr}_{\mathcal{D}}}(Y|X).$$

Simple evaluation metrics, such as the accuracy of a classifier, or the mean squared error of a regressor, also arise naturally as inconsistencies.

**Proposition 6** (Log Accuracy as Inconsistency). *Consider functions  $f, h : X \rightarrow Y$  from inputs to labels,*

*where  $h$  is a predictor and  $f$  generates the true labels. The inconsistency of believing  $f$  and  $h$  (with any confidences), and a distribution  $D(X)$  with confidence  $\beta$ , is  $\beta$  times the log accuracy of  $h$ . That is,*

$$\left\langle \left\langle \begin{array}{c} D \xrightarrow{(\beta)} \text{X} \xrightarrow{\begin{array}{c} h \text{ (r)} \\ f \text{ (s)} \end{array}} \text{Y} \end{array} \right\rangle \right\rangle = -\beta \log_{x \sim D} \Pr_{x \sim D}(f(x) = h(x)) = \beta I_D[f = h]. \quad (3)$$

One often speaks of the accuracy of a hypothesis  $h$  (leaving the true labels  $f$  and empirical distribution  $D$  implicit). But in some sense,  $D(X)$  plays a more primary role: the inconsistency in (3) is scaled by the confidence in  $D$ , and does not depend at all on the confidences in  $h$  or  $f$ . Why is this the case? Because  $f$  is deterministic, codes optimized for it cannot express a sample  $(x, y)$  such that  $y \neq f(x)$ , and so a joint distribution  $\mu$  incurs infinite cost if  $\mu(x, y) > 0$ . The same is true for  $h$ , so we can only consider  $\mu$  such that  $\mu(f=h)=1$ , a restriction which in turn generates inconsistency equal to  $D$ ’s surprisal that  $h$  is correct. In other words, the optimal distribution  $\mu^*$  throws out incorrect samples, so its conditional  $\mu^*(Y|x)$  is undefined unless  $h(x)$  is already correct. This illustrates why accuracy gives no gradient information for training  $h$  (only  $D$ ). Note that this is precisely the opposite of how cross entropy played out in [Proposition 5](#): there we were unwilling to budge on either the true labels or input distribution, and the optimal distribution told us how to modify  $h$ . Even these simple observations—the lack of gradient information, and relationships to other metrics—are clarified by an underlying model.

When  $Y \cong \mathbb{R}^n$ , an estimator  $h(Y|X)$  is referred to as a regressor instead of a classifier. In this setting, most answers are incorrect, but some more so than others. A common way of measuring incorrectness is with mean squared error (MSE):  $\mathbb{E}|f(X) - Y|^2$ . MSE is also the inconsistency of believing that the labels and predictor have Gaussian noise—often a reasonable assumption due to the central limit theorem.

**Proposition 7** (MSE as Inconsistency).

$$\left\langle \left\langle \begin{array}{c} D \xrightarrow{(\infty)} \text{X} \xrightarrow{\begin{array}{c} f \rightarrow \mu_f \mathcal{N}_1 \\ h \rightarrow \mu_h \mathcal{N}_1 \end{array}} \text{Y} \end{array} \right\rangle \right\rangle = \mathbb{E}_D |f(X) - h(X)|^2 =: \text{MSE}_D(f, h),$$

where  $\mathcal{N}_1(Y|\mu)$  is a unit Gaussian on  $Y$  with mean  $\mu$ .

In the appendix, we treat general Gaussian predictors, with arbitrary variances and confidences.

## 4 REGULARIZERS AND PRIORS

Regularizers are extra terms added to loss functions, which provide a source of inductive bias on model

not sure what a partial observation is?

link to proof

link to proof

link to proof

link to proof



parameters. There is a well-known correspondence between using a regularizer and doing maximum *a posteriori* inference with a prior,<sup>3</sup> in which L2 regularization corresponds to a Gaussian prior (Rennie 2003), while L1 regularization corresponds to a Laplacian prior (Williams 1995). Note that the ability to make principled modeling choices about regularizers is a primary benefit of this correspondence. Our approach provides another justification of it.

**Proposition 8.** *Suppose you have a parameterized model  $p(Y|\Theta)$ , a prior  $q(\Theta)$ , and a trusted distribution  $D(Y)$ . The inconsistency of also believing  $\Theta = \theta$  is the cross entropy loss, plus the regularizer:  $\log \frac{1}{q(\theta)}$  times your confidence in  $q$ . That is,*

$$\left\langle \left( \begin{array}{c} q \\ \beta \end{array} \right) \rightarrow \Theta \xrightarrow{p} Y \right\rangle_{D \uparrow_{(\infty)}} = \mathbb{E}_{y \sim D} \log \frac{1}{p(y|\theta)} + \beta \log \frac{1}{q(\theta)} - H(D) \quad (4)$$

If our prior is  $q(\theta) = \frac{1}{k} \exp(-\frac{1}{2}\theta^2)$ , a (discretized) unit gaussian, then the right hand side of (4) becomes

$$\underbrace{\mathbb{E}_D \log \frac{1}{p(Y|\theta)}}_{\text{Cross entropy loss (data-fit cost of } \theta)} + \underbrace{\frac{\beta}{2}\theta_0}_{\text{L2 regularizer (complexity cost of } \theta)} + \underbrace{\beta \log k - H(D)}_{\text{constant in } p \text{ and } \theta},$$

which is the L2 regularized version of Proposition 3. Moreover, the regularization strength corresponds exactly to the confidence  $\beta$ . What about other priors? It is not difficult to see that if we use a (discretized) unit Laplacian prior,  $q(\theta) \propto \exp(-|\theta|)$ , the second term instead becomes  $\beta|\theta_0|$ , which is L1 regularization. More generally, to consider a complexity measure  $U(\theta)$ , we need only include the Gibbs distribution  $\Pr_U(\theta) \propto \exp(-U(\theta))$  into our PDG. We remark that nothing here is specific to cross entropy; any of the objectives we describe can be regularized in this way.

## 5 STATISTICAL DISTANCES AS INCONSISTENCIES

Suppose you are concerned with a single variable  $X$ . One friend has told you that it is distributed according to  $p(X)$ ; another has told you that it follows  $q(X)$ . You adopt both beliefs. Your mental state will be inconsistent if (and only if)  $p \neq q$ , with more inconsistency the more  $p$  and  $q$  differ. Thus the inconsistency of a PDG comprising  $p$  and  $q$  is a measure of divergence. Recall that a PDG also allows us to specify the confidences  $\beta_p$  and  $\beta_q$  of each cpd, so we can form a PDG divergence  $D_{(r,s)}^{\text{PDG}}(p||q)$  for every setting  $(r, s)$  of  $(\beta_p, \beta_q)$ . It turns out that a large class of statistical divergences arise in this way. We start with a familiar one.

<sup>3</sup>A full account can be found in the appendix.

**Proposition 9** (KL Divergence as Inconsistency). *The inconsistency of believing  $p$  with complete certainty, and also  $q$  with some finite certainty  $\beta$ , is  $\beta$  times the KL Divergence (or relative entropy) of  $q$  with respect to  $p$ . That is,*

$$\left\langle \left( \begin{array}{c} p \\ (\infty) \end{array} \right) \rightarrow X \leftarrow \left( \begin{array}{c} q \\ \beta \end{array} \right) \right\rangle = \beta D(p || q).$$

This result gives us an intuitive interpretation of the asymmetry of relative entropy / KL divergence, and a prescription about when it makes sense to use it.  $D(p || q)$  is the inconsistency of a mental state containing both  $p$  and  $q$ , when absolutely certain of  $p$  (and not willing to budge on it). This concords with the standard intuition that  $D(p || q)$  reflects the amount of information required to change  $q$  into  $p$ , which is why we call it the relative entropy “from  $q$  to  $p$ ”.

We now consider the general case of a PDG comprising  $p(X)$  and  $q(X)$  with arbitrary confidences.

**Lemma 10.** *The inconsistency of a PDG comprising  $p(X)$  with confidence  $r$  and  $q(X)$  with confidence  $s$  is given in closed form by*

$$\left\langle \left( \begin{array}{c} p \\ (r) \end{array} \right) \rightarrow X \leftarrow \left( \begin{array}{c} q \\ (s) \end{array} \right) \right\rangle = -(r+s) \log \sum_x \left( p(x)^r q(x)^s \right)^{\frac{1}{r+s}}.$$

Of the many generalizations of KL divergence, Rényi divergences, first characterized by Alfréd Rényi 1961 are perhaps the most significant, as few others have found either application or an interpretation in terms of coding theory (Van Erven and Harremos 2014). The Rényi divergence of order  $\alpha$  between two distributions  $p(X)$  and  $q(X)$  is given by

$$D_\alpha(p || q) := \frac{1}{1-\alpha} \log \sum_{x \in \mathcal{V}(X)} p(x)^\alpha q(x)^{1-\alpha}. \quad (5)$$

Rényi introduced this measure in the same paper as the more general class of  $f$ -divergences, but directs his attention towards those of the form (5), because they satisfy a natural weakening of standard postulates for Shannon entropy due to Fadeev (1957). Concretely, every symmetric, continuous measure that additively separates over independent events, and with a certain “mean-value property”, up to scaling, is of the form (5) for some  $\alpha$  (Rényi 1961). It follows from Lemma 10 that carves out almost the same class: every Rényi divergence is a PDG divergence, and every (non-limiting) PDG divergence is a (scaled) Rényi divergence.

**Corollary 10.1** (Rényi Divergences).

$$\left\langle \left( \begin{array}{c} p \\ (r) \end{array} \right) \rightarrow X \leftarrow \left( \begin{array}{c} q \\ (s) \end{array} \right) \right\rangle = s \cdot D_{\frac{r}{r+s}}(p || q)$$

and  $D_\alpha(p || q) = \left\langle \left( \begin{array}{c} p \\ (\frac{1}{1-\alpha}) \end{array} \right) \rightarrow X \leftarrow \left( \begin{array}{c} q \\ 1 \end{array} \right) \right\rangle$

[link to proof]

I really really like this diagram. Is there any way you can get it earlier? It feels like a key to some few contributions really.

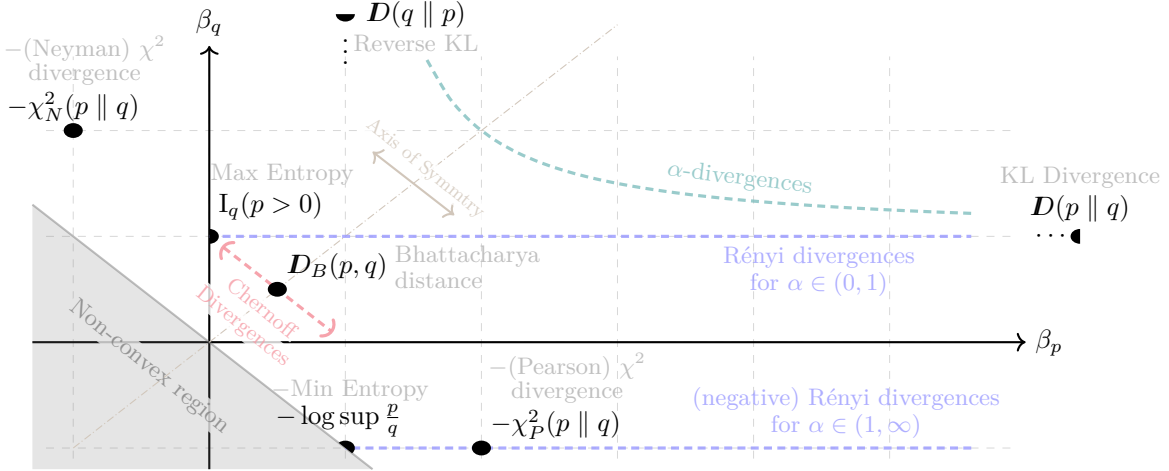


Figure 1: A map of the inconsistency of the PDG comprising  $p(X)$  and  $q(X)$ , as we vary their respective confidences  $\beta_p$  and  $\beta_q$ . Solid circles indicate well-known named measures, semicircles indicate limiting values, and the heavily dashed lines are well-established classes.

However, the two classes are not identical, because the PDG divergences admit extra limit points. One of the biggest differences is that the reverse KL divergence  $D(q \parallel p)$  is not a Rényi entropy of the form  $D_\alpha(p \parallel q)$  for any value (or limit) of  $\alpha$ . This lack of symmetry has led others (e.g., [Cichocki and Amari 2010](#)) to instead work with a re-scaled symmetric version of the Rényi entropy, called  $\alpha$ -divergence, which as an additional factor of  $\frac{1}{\alpha}$ . The relationships between these quantities can be seen in [Figure 1](#).

The Chernoff divergence measures the tightest possible exponential bound on probability of error ([Nielsen 2011](#)) in Bayesian hypothesis testing. It also happens to be the smallest possible inconsistency of simultaneously believing  $p$  and  $q$ , with total confidence 1.

**Corollary 10.2.** *The Chernoff Divergence between  $p$  and  $q$  equals*

$$\inf_{\beta \in (0,1)} \left\langle \left\langle \frac{p}{(\beta)} \rightarrow \boxed{X} \leftarrow \frac{q}{(1-\beta)} \right\rangle \right\rangle.$$

One significant consequence of representing divergences as inconsistencies is that we can use [Lemma 1](#) to derive relationships between them. The following facts follow directly from [Figure 1](#), by inspection.

**Corollary 10.3.** 1. Rényi entropy is monotonic in its parameter  $\alpha$ .  
2.  $D(p \parallel q) \geq 2D_B(p, q) \leq D(q \parallel p)$ .  
3. If  $q(p > 0) < 1$  (i.e.,  $q \not\ll p$ ), then  $D(q \parallel p) = \infty$ .

These divergences correspond to PDGs with only two edges and one variable. What about more complex graphs? For a start, the usual notion of a conditional divergence  $D_{r,s}^{\text{PDG}}(p(Y|X) \parallel q(Y|X) | r(X)) := \mathbb{E}_{x \sim r} D_{r,s}^{\text{PDG}}(p(Y|x) \parallel q(Y|x))$  falls out of PDGs of the form

$$\frac{r}{(\infty)} \rightarrow \boxed{X} \begin{matrix} \xrightarrow{p(r)} \\ \xleftarrow{q(s)} \end{matrix} \boxed{Y}.$$

[Lemma 1](#), together with minor structural manipulation, gives visual proofs of standard divergence properties. Such a proof of the Data Processing Inequality can be seen in [Figure 2](#). And in general, PDG inconsistency can be viewed as a vast generalization of these divergences to arbitrary structured objects.

## 6 VARIATIONAL OBJECTIVES AND BOUNDS

The fact that the incompatibility of  $\mathcal{M}$  with a *specific* joint distribution  $\mu$  is an upper bound on the inconsistency is not a deep one, but it is of a variational flavor. Here, we focus on the more surprising converse: PDG semantics capture general aspects of variational inference and provide a graphical proof language for it.

### 6.1 PDGs and Variational Approximations

We begin by recounting the standard development of the ‘Evidence Lower Bound’ (ELBO), a standard objective for training latent variable models ([Blei, Kucukelbir, and McAuliffe 2017](#), §2.2). Suppose we have a model  $p(X, Z)$ , but only have access to observations of  $x$ . In service of adjusting  $p(X, Z)$  to make our observations more likely, we would like to maximize  $\log p(X=x)$ , the “evidence” of  $x$  ([Proposition 4](#)). Unfortunately, computing  $p(x) = \sum_z p(x, z)$  requires summing over all of  $Z$ , which can be intractable. The variational approach is as follows: fix a family of distributions  $\mathcal{Q}$  that is easy to sample from, choose some

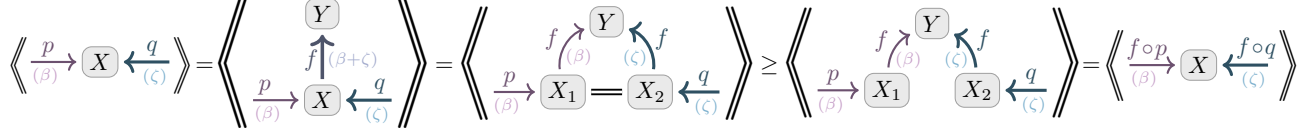


Figure 2: A visual proof of the data-processing inequality:  $D_{(\beta, \zeta)}^{\text{PDG}}(p \parallel q) \geq D_{(\beta, \zeta)}^{\text{PDG}}(f \circ p \parallel f \circ q)$ . In words: the cpd  $f(Y|X)$  can always be satisfied, so adds no inconsistency. It is then equivalent to split  $f$  and the variable  $X$  into  $X_1$  and  $X_2$  with edges enforcing  $X_1 = X_2$ . But removing such edges can only decrease inconsistency. Finally, compose the remaining cpds to give the result. A full justification can be found in the appendix.

$q(Z) \in \mathcal{Q}$ , and define  $\text{ELBO}_{p,q}(x) := \mathbb{E}_{z \sim q} \log \frac{p(x,z)}{q(z)}$ . This is something we can estimate, since we can sample from  $q$ . By Jensen's inequality,

$$\text{ELBO}_{p,q}(x) = \mathbb{E}_q \log \frac{p(x, Z)}{q(Z)} \leq \log \left[ \mathbb{E}_q \frac{p(x, Z)}{q(Z)} \right] = \log p(X),$$

with equality if  $q(Z) = p(Z)$ . So to maximize  $p(X)$ , it suffices to adjust  $p$  and  $q$  to maximize  $\text{ELBO}_{p,q}(x)$ ,<sup>4</sup> provided  $\mathcal{Q}$  is expressive enough.

The formula for the ELBO is somewhat difficult to make sense of.<sup>5</sup> Nevertheless, it arises naturally as the inconsistency of the appropriate PDG.

**Proposition 11.** *The negative ELBO of  $x$  is the inconsistency of the PDG containing  $p, q$ , and  $X = x$ , with high confidence in  $q$ . That is,*

$$-\text{ELBO}_{p,q}(x) = \left\langle \frac{q}{(\infty)} \rightarrow Z \xrightarrow{p} X \leftarrow x \right\rangle.$$

Owing to its structure, a PDG is often more intuitive and easier to work with than the formula for its inconsistency. To illustrate, we now give a simple and visually intuitive proof of the bound traditionally used to motivate the ELBO, via Lemma 1:

$$\log \frac{1}{p(x)} = \left\langle \frac{p}{Z} \rightarrow X \right\rangle \leq \left\langle \frac{q}{Z} \rightarrow X \right\rangle = -\text{ELBO}_{p,q}(x).$$

The first and last equalities are Propositions 4 and 11 respectively. Now to reap some pedagogical benefits. The second PDG has more edges so it is clearly at least as inconsistent. Furthermore, it's easy to see that equality holds when  $q(Z) = p(Z)$ : the best distribution for the left PDG has marginal  $p(Z)$  anyway, so insisting on it incurs no further cost.

## 6.2 Variational Auto-Encoders and PDGs

An autoencoder is a probabilistic model intended to compress a variable  $X$  (e.g., an image) to a compact

latent representation  $Z$ . Its structure is given by two conditional distributions: an encoder  $e(Z|X)$ , and a decoder  $d(X|Z)$ . Of course, not all pairs of cpds fill this role equally well. One important consideration is the *reconstruction error* (6): when we decode an encoded image, we would like it to resemble the original.

$$\text{Rec}(x) := \mathbb{E}_{z \sim e(Z|x)} \underbrace{I_{d(X|z)}(x)}_{\left( \begin{array}{l} \text{additional bits required to} \\ \text{decode } x \text{ from its encoding } z \end{array} \right)} = \sum_z e(z|x) \log \frac{1}{d(x|z)} \quad (6)$$

There are other desiderata as well. Perhaps good latent representations  $Z$  have uncorrelated components, and are normally distributed. We encode such wishful thinking as a belief  $p(Z)$ , known as a variational prior.

The data of a *variational auto-encoder* (Kingma and Welling 2014) consists of  $e(Z|X)$ ,  $d(X|Z)$ , and  $p(Z)$ . The encoder  $e(Z|X)$  can be used as a variational approximation of  $Z$ , differing from  $q(Z)$  of Section 6.1 only in that it can depend on  $X$ . Here, the analogue of the ELBO becomes

$$\text{ELBO}_{p,e,d}(x) := \mathbb{E}_{z \sim e(Z|x)} \left[ \log \frac{p(z)d(x|z)}{e(z|x)} \right] = -\text{Rec}(x) - D(e(Z|x) \parallel p).$$

This gives us the following analog of Proposition 11.

**Proposition 12.** *The VAE loss of a sample  $x$  is the inconsistency of the PDG comprising the encoder (with high confidence, as it defines the encoding), decoder  $d$ , prior  $p$ , and  $x$ . That is,*

$$-\text{ELBO}_{p,e,d}(x) = \left\langle \frac{p}{Z} \xrightarrow{d} X \xleftarrow{e} x \right\rangle.$$

Propositions 12 and 16 can be used to derive another variational bound. Once again, the addition of the edge  $e$  cannot decrease the inconsistency (Lemma 1), but believing it with high confidence does make it possible to know the encoding of a sample, making inference tractable. This results in the following simple visual proof:

$$\log \frac{1}{\Pr(x)} = \left\langle \frac{p}{Z} \xrightarrow{d} X \right\rangle \leq \left\langle \frac{p}{Z} \xrightarrow{d} X \xleftarrow{e} x \right\rangle = -\text{ELBO}_{p,e,d}(x).$$

<sup>4</sup>For many iid samples:  $\max_{p,q} \sum_{x \in \mathcal{D}} \text{ELBO}_{p,q}(x)$ .

<sup>5</sup>Especially if  $p$  and  $q$  are densities. See Appendix A

[link to proof]

[link to proof]

Here  $\Pr(X)$  is the marginal of  $p(Z)d(X|Z)$  on  $X$ . See the appendix for multi-sample analogs of the bound and both propositions.

### 6.3 $\beta$ -VAE Objective

The ELBO is not the only objective that has been used to train networks with a VAE structure. In the most common variant, Higgins et al. (2016) argue that one might want to weight the reconstruction error (6) and the ‘KL term’ differently. They suggest an objective of the form

$$\beta\text{-ELBO}_{p,e,d}(x) := \text{Rec}(x) - \beta \mathbf{D}(e(Z|x) \parallel p),$$

which when  $\beta = 1$  is the ELBO as before. The authors view  $\beta$  as a regularization strength, and argue that it sometimes helps to have a stronger prior. Sure enough, the  $\beta$ -VAE objective is the inconsistency of the same PDG as before, but with confidence  $\beta$  in  $p(Z)$ .<sup>6</sup>

## 7 FREE ENERGY AND INCONSISTENCY

A weighted factor graph  $\Psi = (\phi_J, \theta_J)_{J \in \mathcal{J}}$ , where each  $\theta_J$  is a real-valued weight,  $J$  is associated with a subset of variables  $\mathbf{X}_J$ , and  $\phi_J : \mathcal{V}(\mathbf{X}_J) \rightarrow \mathbb{R}$ , determines a distribution by

$$\Pr_{\Psi}(\mathbf{x}) = \frac{1}{Z_{\Psi}} \prod_{J \in \mathcal{J}} \phi_J(\mathbf{x}_J)^{\theta_J}.$$

$Z_{\Psi}$  is the constant  $\sum_{\mathbf{x}} \prod_{J \in \mathcal{J}} \phi_J(\mathbf{x}_J)^{\theta_J}$  required to normalize the distribution, and is known as the *partition function*. Computing  $\log Z_{\Psi}$  is intimately related to much of probabilistic inference in factor graphs (Ma et al. 2013). Following Richardson and Halpern (2021), let  $\mathbf{m}_{\Psi}$  be the PDG with edges  $\{ \overset{J}{\rightarrow} \mathbf{X}_J \}_{J \in \mathcal{J}}$ , cpds  $p_J(\mathbf{X}_J) \propto \phi_J(\mathbf{X}_J)$ , and weights  $\alpha_J, \beta_J := \theta_J$ .

If the factors are normalized and all variables are edge targets, then  $Z_{\Psi} \leq 1$ , so  $\log \frac{1}{Z_{\Psi}} \geq 0$  measures how far the product of factors is from being a probability distribution. So in a sense, it measures  $\Psi$ ’s inconsistency.

**Proposition 13.** *For any weighted factor graph  $\Psi$  we have  $\langle \mathbf{m}_{\Psi} \rangle_1 = -\log Z_{\Psi}$ .*

The exponential families generated by weighted factor graphs are a cornerstone of statistical mechanics, where  $-\log Z_{\Psi}$  is known as the (Heimholz) free energy. The principle of free-energy minimization has been enormously successful in describing of not only chemical and biological systems (Chipot and Pohorille 2007), but also cognitive ones (Friston 2009).

<sup>6</sup>Both names originate in thermodynamic coldness  $\beta$ .

## 8 REVERSE-ENGINEERING A LOSS

Given an arbitrary loss function  $\ell(X)$ , can we find a PDG that gives rise to it? To a first approximation, the answer seems to be yes. Without affecting its semantics, one may add the variable  $T$  that takes values  $\{\mathbf{t}, \mathbf{f}\}$ , and the event  $T = \mathbf{t}$ , to any PDG. Now, given a non-negative cost function, define the cpd  $\hat{c}(T|X)$  by  $\hat{c}(\mathbf{t}|x) := \exp(-c(x))$ . By threatening to generate the falsehood  $\mathbf{f}$ , with probability dependent on the cost of  $X$ ,  $\hat{c}$  ties the value of  $X$  to inconsistency.

**Proposition 14.**  $\left\langle \left\langle \frac{p}{(\infty)} \rightarrow X \right\rangle \hat{c} \rightarrow T \leftarrow \mathbf{t} \right\rangle = \mathbb{E}_{x \sim p} c(x)$ .

[\[link to proof\]](#)

Setting  $\beta_p := \infty$  is not realistic since we’re still training the model  $p$ , but it is necessary to recover the formula. Any bridge we build to generate inconsistency based on the value of  $X$  will also work in reverse: the PDG squirms and contorts the probability of  $X$  to disperse the inconsistency. One cannot simply “emit inconsistency” and expect it collect without affecting the optimal probabilities. With more complexity, even setting every  $\beta := \infty$  may not be enough to prevent the squirming. To illustrate, consider a model of the supervised learning setting, with data  $\mathcal{D}$ , model  $p$ , and an arbitrary loss function  $\ell$ . Define:

$$\mathcal{S} := \begin{array}{c} \Pr_{\mathcal{D}} \rightarrow Y \\ \downarrow \text{Pr}_{\mathcal{D}} \\ X \xrightarrow{p} Y' \end{array} \quad \begin{array}{c} \hat{\ell} \\ \uparrow \\ T \end{array} \quad \text{and} \quad L := \mathbb{E}_{\substack{(x,y) \sim \Pr_{\mathcal{D}} \\ y' \sim p(Y'|x)}} [\ell(y, y')].$$

Given Proposition 14, one might imagine  $\langle \mathcal{S} \rangle = L$ , but this is not so. But  $\langle \mathcal{S} \rangle$  is arguably nicer. The optimal  $p(Y|X)$  according to  $L$  is a point mass on the label  $y_X^*$  minimizing expected loss, but the optimal  $p(Y|X)$  according to  $\langle \mathcal{S} \rangle$  is  $\Pr_{\mathcal{D}}(Y|X)$ , which is *calibrated* (CITATION). If we also strictly enforce the causal picture, the squirming finally stops, as we arrive at  $\lim_{\gamma \rightarrow \infty} \langle \mathcal{S} \rangle_{\gamma} = L$ .

To summarize, working backwards requires making strange choices, and indicating absolute certainty in them. In the end, we must confront our modeling choices; good loss functions come from good models.

## 9 FINAL REMARKS

We have now seen that PDG semantics not only capture structured objects such as Bayesian Networks and Factor Graphs as in Richardson and Halpern (2021), but in the same stroke also generate many standard loss functions, including some non-trivial ones. In each case, the appropriate loss is a simple consequence of carefully articulating modeling assumptions. Viewing

[\[link to proof\]](#)



loss functions in this way also has beneficial side effects, including an intuitive visual proof language for reasoning about the relationships between them.

Our “universal loss function”, which blurs the line between model and objective function, may be of particular interest to the AI safety community.

## Acknowledgements

## References

- Blei, David M, Alp Kucukelbir, and Jon D McAuliffe (2017). “Variational Inference: A Review for Statisticians.” In: *Journal of the American statistical Association* 112.518, pp. 859–877.
- Chipot, Christophe and Andrew Pohorille (2007). “Free Energy Calculations.” In: *Springer Series in Chemical Physics* 86, pp. 159–184.
- Cichocki, Andrzej and Shun-ichi Amari (2010). “Families of Alpha Beta and Gamma Divergences: Flexible and Robust Measures of Similarities.” In: *Entropy* 12.6, pp. 1532–1568.
- Fadeev, DK (1957). “Zum Begriff der Entropie einer endlichen Wahrscheinlichkeitsschemas.” In: *Arbeiten zur Informationstheorie I. Deutscher Verlag der Wissenschaften*, pp. 85–90.
- Friston, Karl (2009). “The Free-Energy Principle: a Rough Guide to the Brain?” In: *Trends in Cognitive Sciences* 13.7, pp. 293–301.
- Grover, Aditya and Stefano Ermon (2018). *Lecture notes in Deep Generative Models*. [deepgenerativemodels.github.io/notes/](https://deepgenerativemodels.github.io/notes/).
- Hastie, Trevor, Robert Tibshirani, and Jerome Friedman (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Second. Springer.
- Higgins, Irina et al. (2016). “Beta-VAE: Learning Basic visual concepts with a constrained variational framework.” In:
- Jadon, Shruti (2020). “A Survey of Loss Functions for Semantic Segmentation.” In: *2020 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB)*. IEEE, pp. 1–7.
- Kingma, Diederik P and Max Welling (2014). “Auto-Encoding Variational Bayes.” In: *Proceedings of the International Conference on Learning Representations (ICLR)*. arXiv: [1312.6114](https://arxiv.org/abs/1312.6114) [stat.ML].
- Ma, Jianzhu et al. (2013). “Estimating the Partition Function of Graphical Models using Langevin Importance Sampling.” In: *Artificial Intelligence and Statistics*. PMLR, pp. 433–441.
- MacKay, David (2003). *Information Theory, Inference and Learning Algorithms*. Cambridge University Press.
- Myung, In Jae (2003). “Tutorial on Maximum Likelihood Estimation.” In: *Journal of mathematical Psychology* 47.1, pp. 90–100.
- Nielsen, Frank (2011). “Chernoff Information of Exponential Families.” In: *arXiv preprint arXiv:1102.2684*.
- Rennie, Jason (2003). “On l2-norm regularization and the Gaussian prior.” In:
- Rényi, Alfréd (1961). “On Measures of Entropy and Information.” In: *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*. University of California Press, pp. 547–561.
- Richardson, Oliver and Joseph Y Halpern (2021). “Probabilistic Dependency Graphs.” In: *AAAI ’21*. arXiv: [2012.10800](https://arxiv.org/abs/2012.10800) [cs.AI].
- Tribus, Myron (1961). “Information Theory as the Basis for Thermostatistics and Thermodynamics.” In:
- Van Erven, Tim and Peter Harremoës (2014). “Rényi Divergence and Kullback-Leibler divergence.” In: *IEEE Transactions on Information Theory* 60.7, pp. 3797–3820.
- Wang, Qi et al. (2020). “A Comprehensive Survey of Loss Functions in Machine Learning.” In: *Annals of Data Science*, pp. 1–26.
- Williams, Peter M (1995). “Bayesian regularization and pruning using a Laplace prior.” In: *Neural Computation* 7.1, pp. 117–143.

Name	$p$	Formula
Harmonic	$(p = -1):$	$\text{HM}_w(\mathbf{r}) = 1 / (\sum_{i=1}^n w_i / r_i)$
Geometric	$(\lim p \rightarrow 0):$	$\text{GM}_w(\mathbf{r}) = \prod_{i=1}^n r_i^{w_i}$
Arithmetic	$(p = 1):$	$\text{AM}_w(\mathbf{r}) = \sum_{i=1}^n w_i r_i$
Quadratic	$(p = 2):$	$\text{QM}_w(\mathbf{r}) = \sqrt{\sum_{i=1}^n w_i r_i^2}$

Table 1: special cases of the  $p$ -power mean  $\text{M}_p^w(\mathbf{r})$

## A THE FINE PRINT FOR PROBABILITY DENSITIES

Many of our results ([Propositions 2 to 5](#), [11](#), [12](#), [16](#), [18](#) and [19](#)) technically require the distribution to be represented with a mass function (not a density function (pdf)). A PDG containing both pdf and a finitely supported distribution on the same variable will typically have infinite inconsistency. But this is not just a quirk of the PDG formalism.

Probability density is not dimensionless (like probability mass), but rather has inverse  $X$ -units (e.g., probability per meter), so depends on an arbitrary choice of scale (the pdf for probability per meter and per centimeter will yield different numbers). In places where the objective does not have units that cancel before we take a logarithm, the use of a probability density  $p(X)$  becomes sensitive to this arbitrary choice of parameterization. For instance, the analog of surprisal,  $-\log p(x)$  for a pdf  $p$ , or its expectation, called differential entropy, both suffer from this problem. On the other hand, this choice of scale ultimately amounts to an additive constant.

Moreover, beyond a certain point, decreasing the discretization size  $k$  of a discretized approximation  $\tilde{p}_k(X)$  *also* contributes a constant that depends only on  $k$ . But such constants are irrelevant for optimization, justifying the use of the continuous analogues as loss functions.

The bottom line is that these results hold for EVERY discretization — but in the limit as the discretization becomes smaller, the relevant quantity may diverge. However, this divergence stems from an additive constant, which is irrelevant for optimization. Therefore, using the densities results in a morally equivalent, but finite loss function.

## B FURTHER RESULTS AND GENERALIZATIONS

### B.1 Full Characterization of Gaussian Predictors

The more general result, which is most clearly articulated in terms of a power mean.

**Definition 2.** The weighted power mean  $\text{M}_p^w(\mathbf{r})$  of the collection of real numbers  $\mathbf{r} = r_1, \dots, r_n$  with respect to the convex weights  $w = w_1, \dots, w_n$  satisfying  $\sum_i w_i = 1$ , is given by

$$\text{M}_p^w(\mathbf{r}) := \left( \sum_{i=1}^n w_i (r_i)^p \right)^{\frac{1}{p}}.$$

We omit the superscript as a shorthand for the uniform weighting  $w_i = 1/N$ . Most standard means, such as those in [Table 1](#), are special cases.  $\square$

It is well known that  $\text{M}_p^w(\mathbf{r})$  is increasing in  $p$ , and strictly so if not all elements of  $\mathbf{r}$  are identical. In particular,  $\text{QM}_w(a, b) > \text{GM}_w(a, b)$  for all  $a \neq b$  and positive weights  $w$ . We now present the result.

**Proposition 15.** *Consider a PDG containing two (distinct) conditional Gaussian distributions on a variable*

[\[link to proof\]](#)

$Y$ , whose parameters can both depend on a variable  $X$ . Its inconsistency takes the form

$$\left\langle\!\!\left\langle \begin{array}{c} \xrightarrow{D!} X \\ \begin{array}{c} \xrightarrow{f} \mu_1 \\ \xrightarrow{t} \sigma_1 \\ \xrightarrow{s} \sigma_2 \\ \xrightarrow{h} \mu_2 \end{array} \\ \xrightarrow{\mathcal{N}} Y \end{array} \right\rangle\!\!\right\rangle = \frac{1}{2} \mathbb{E}_D \left[ \text{HM}(\beta_1, \beta_2) \frac{1}{2} \left( \frac{\mu_1 - \mu_2}{\text{QM}_{\hat{\beta}}(\sigma_1, \sigma_2)} \right)^2 + \text{AM}(\beta_1, \beta_2) \log \frac{\text{QM}_{\hat{\beta}}(\sigma_1, \sigma_2)}{\text{GM}_{\hat{\beta}}(\sigma_1, \sigma_2)} \right] \quad (7)$$

$$= \mathbb{E}_{x \sim D} \left[ \frac{\beta_1 \beta_2}{2} \frac{(f(x) - h(x))^2}{\beta_2 s(x)^2 + \beta_1 t(x)^2} + \frac{\beta_1 + \beta_2}{2} \log \frac{\beta_2 s(x)^2 + \beta_1 t(x)^2}{(\beta_1 + \beta_2)(s(x)^{\beta_2} t(x)^{\beta_1})^{\frac{1}{\beta_1 + \beta_2}}} \right]$$

where  $\hat{\beta} = (\frac{\beta_2}{\beta_1 + \beta_2}, \frac{\beta_1}{\beta_1 + \beta_2})$  represents the normalized and reversed vector of confidences  $\beta = (\beta_1, \beta_2)$  for the two distributions, and  $\mu_1 = f(X)$ ,  $\mu_2 = g(X)$ ,  $\sigma_1 = s(X)$ ,  $\sigma_2 = t(X)$  are random variables over  $X$ .

Plugging in  $s(x) = t(x) = 1$  and  $\beta_1 = \beta_2 = 1$  proves:

**Proposition 7.**

$$\left\langle\!\!\left\langle \begin{array}{c} \xrightarrow[\infty]{D} X \\ \begin{array}{c} \xrightarrow{f} \mu_f \\ \xrightarrow{h} \mu_h \end{array} \\ \xrightarrow{\mathcal{N}_1} Y \end{array} \right\rangle\!\!\right\rangle = \mathbb{E}_D |f(X) - h(X)|^2 =: \text{MSE}_D(f, h),$$

where  $\mathcal{N}_1(Y|\mu)$  is a unit Gaussian on  $Y$  with mean  $\mu$ .

This PDG is also equal to

$$\left\langle\!\!\left\langle \begin{array}{c} \xrightarrow{D!} X \\ \begin{array}{c} \xrightarrow{\mathcal{N}(f(x), 1)} \\ \xrightarrow{\mathcal{N}(g(x), 1)} \end{array} \\ \xrightarrow{\quad} Y \end{array} \right\rangle\!\!\right\rangle$$

which illustrates an orthogonal point: that PDGs handle composition of functions as one would expect, so that it is equivalent to model an entire process as a single arrow, or to break it into stages, ascribing an arrow to each stage, with one step of randomization.

As a bonus, [Proposition 15](#) also gives a proof of the inequality of the weighted geometric and quadratic means.

**Corollary 15.1.** For all  $\beta, \sigma_1, \sigma_2$ , and  $D$ , we have:

$$\frac{\mathbb{E}}{D} \frac{\text{QM}_{\hat{\beta}}(\sigma_1, \sigma_2)}{\text{GM}_{\hat{\beta}}(\sigma_1, \sigma_2)} > 0.$$

## B.2 Full-Dataset ELBO and Bounds

We now present the promised multi-sample analogs from [Section 6.1](#).

**Proposition 16.** The following analog of [Proposition 12](#) for a whole dataset  $\mathcal{D}$  holds:

[link to proof](#)

$$- \mathbb{E}_{\text{Pr}_{\mathcal{D}}} \text{ELBO}_{p,e,d}(X) = \left\langle\!\!\left\langle \begin{array}{c} \xrightarrow{p} Z \\ \begin{array}{c} \xrightarrow{d} X \\ \xleftarrow{e!} X \end{array} \\ \xleftarrow{\text{Pr}_{\mathcal{D}}!} \end{array} \right\rangle\!\!\right\rangle + H(\text{Pr}_{\mathcal{D}}).$$

[Propositions 3](#) and [16](#) then give us an analog of the visual bounds in the body of the main paper ?? for many

i.i.d. datapoints at once, with only a single application of the inequality:

$$\begin{aligned}
 -\log \Pr(\mathcal{D}) &= -\log \prod_{i=1}^m \left( \Pr(x^{(i)}) \right) = -\frac{1}{m} \sum_{i=1}^m \log \Pr(x^{(i)}) = \\
 &= \mathbb{H}(\Pr_{\mathcal{D}}) + \left\langle \left\langle p \rightarrow Z \xrightarrow{d} X \xleftarrow{\Pr_{\mathcal{D}}!} \right\rangle \right\rangle \leq \left\langle \left\langle p \rightarrow Z \xrightleftharpoons[e!]{d} X \xleftarrow{\Pr_{\mathcal{D}}!} \right\rangle \right\rangle + \mathbb{H}(\Pr_{\mathcal{D}}) \\
 &= -\mathbb{E}_{\Pr_{\mathcal{D}}} \mathbb{ELBO}_{p,e,d}(X)
 \end{aligned}$$

We also have the following formal statement of the claim made in ??.

**Proposition 17.** *The negative  $\beta$ -ELBO objective for a prior  $p(X)$ , encoder  $e(Z | X)$ , decoder  $d(X | Z)$ , at a sample  $x$ , is equal to the inconsistency of the corresponding PDG, where the prior has confidence equal to  $\beta$ . That is,*

$$-\beta\text{-ELBO}_{p,e,d}(x) = \left\langle \left\langle \overset{(\beta)}{p} \rightarrow Z \xrightleftharpoons[e!]{d} X \xleftarrow{x} \right\rangle \right\rangle$$

As a specific case (i.e., effectively by setting  $\beta_p := 0$ ), we get the reconstruction error as the inconsistency of an autoencoder (without a variational prior):

**Corollary 17.1** (reconstruction error as inconsistency).

$$-\text{Rec}_{e,d}(x) := \mathbb{E}_{z \sim e(Z|x)} \mathbb{I}_{d(X|z)}(x) = \left\langle \left\langle Z \xrightleftharpoons[e!]{d} X \xleftarrow{x} \right\rangle \right\rangle$$

### B.3 More Variants of Cross Entropy Results

**Proposition 18.** *Given a model determining a probability distribution with mass function  $p(X)$ , and samples  $\mathcal{D} = \{x_i\}_{i=1}^m$  determining an empirical distribution  $\Pr_{\mathcal{D}}$ , the following are equal, for all  $\gamma \geq 0$ :*

1. The average negative log likelihood  $\ell(p; \mathcal{D}) = -\frac{1}{m} \sum_{i=1}^m \log p(x_i)$
2. The cross entropy of  $p$  relative to  $\Pr_{\mathcal{D}}$
3.  $\llbracket p \rrbracket_{\gamma}(\Pr_{\mathcal{D}}) + (1 + \gamma) \mathbb{H}(\Pr_{\mathcal{D}})$
4.  $\left\langle \left\langle p \rightarrow X \xleftarrow{\Pr_{\mathcal{D}}!} \right\rangle \right\rangle_{\gamma} + (1 + \gamma) \mathbb{H}(\Pr_{\mathcal{D}})$

**Proposition 19.** *The average negative log likelihood  $\ell(p; x) := -\frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \log p(x)$  (which is also the cross entropy) is the inconsistency of the PDG containing  $p$  and the data distribution  $\Pr_{\mathcal{D}}$ , plus the entropy of the data distribution (which is constant in  $p$ ). That is,*

$$\ell(p; \mathcal{D}) = \left\langle \left\langle Z \xleftarrow{p} X \xleftarrow{\Pr_{\mathcal{D}}!} \right\rangle \right\rangle + \mathbb{H}(\Pr_{\mathcal{D}}).$$



## C PROOFS

**Lemma 1.** Suppose PDGs  $\mathcal{M}$  and  $\mathcal{M}'$  differ only in their edges (resp.  $\mathcal{E}$  and  $\mathcal{E}'$ ) and confidences (resp.  $\beta$  and  $\beta'$ ). If  $\mathcal{E} \subseteq \mathcal{E}'$  and  $\beta_L \leq \beta'_L$  for all  $L \in \mathcal{E}$ , then  $\llbracket \mathcal{M} \rrbracket_\gamma \leq \llbracket \mathcal{M}' \rrbracket_\gamma$  for all  $\gamma$ .<sup>7</sup>

*Proof.* For every  $\mu$ , adding more edges only adds non-negative terms to (1). Thus,  $\llbracket \mathcal{M} + \mathcal{M}' \rrbracket_\gamma(\mu) \geq \llbracket \mathcal{M} \rrbracket_\gamma(\mu)$  for all  $\gamma$  and  $\mu$ . So it also holds when we take an infimum over  $\mu$ , yielding  $\llbracket \mathcal{M} + \mathcal{M}' \rrbracket_\gamma \geq \llbracket \mathcal{M} \rrbracket_\gamma$ . Analogously, increasing  $\beta$  results in larger coefficients on the (non-negative) terms of (1) so  $\llbracket \mathcal{M} \rrbracket_\gamma(\mu) \geq \llbracket \mathcal{M}' \rrbracket_\gamma(\mu)$  for all  $\gamma$  and  $\mu$ , so  $\llbracket \mathcal{M} \rrbracket_\gamma \geq \llbracket \mathcal{M}' \rrbracket_\gamma$ .  $\square$

**Proposition 2.** Consider a distribution  $p(X)$ . The inconsistency of the PDG comprising  $p$  and  $X=x$  equals the surprisal  $I_p[X=x]$ . That is,

$$I_p[X=x] = \left\langle \xrightarrow{p} \boxed{X} \xleftarrow{x} \right\rangle.$$

(Recall that  $\llbracket \mathcal{M} \rrbracket$  is the inconsistency of the PDG  $\mathcal{M}$ .)

*Proof.* Any distribution  $\mu(X)$  that places mass on some  $x' \neq x$  will have infinite KL divergence from the point mass on  $x$ . Thus, the only possibility for a finite consistency arises when  $\mu = \delta_x$ , and so

$$\left\langle \xrightarrow{p} \boxed{X} \xleftarrow{x} \right\rangle = \left[ \xrightarrow{p} \boxed{X} \xleftarrow{x} \right](\delta_x) = D(\delta_x \parallel p) = \log \frac{1}{p(x)} = I_p(x).$$

$\square$

**Proposition 18.** Given a model determining a probability distribution with mass function  $p(X)$ , and samples  $\mathcal{D} = \{x_i\}_{i=1}^m$  determining an empirical distribution  $\text{Pr}_{\mathcal{D}}$ , the following are equal, for all  $\gamma \geq 0$ :

1. The average negative log likelihood  $\ell(p; \mathcal{D}) = -\frac{1}{m} \sum_{i=1}^m \log p(x_i)$
2. The cross entropy of  $p$  relative to  $\text{Pr}_{\mathcal{D}}$
3.  $\llbracket p \rrbracket_\gamma(\text{Pr}_{\mathcal{D}}) + (1 + \gamma) H(\text{Pr}_{\mathcal{D}})$
4.  $\left\langle \xrightarrow{p} \boxed{X} \xleftarrow{\text{Pr}_{\mathcal{D}}} \right\rangle_\gamma + (1 + \gamma) H(\text{Pr}_{\mathcal{D}})$

*Proof.* The equality of 1 and 2 is standard. The equality of 3 and 4 can be seen by the fact that in the limit of infinite confidence on  $\text{Pr}_{\mathcal{D}}$ , the optimal distribution must also equal  $\text{Pr}_{\mathcal{D}}$ , so the least inconsistency is attained at this value. Finally it remains to show that the first two and last two are equal:

$$\begin{aligned} \llbracket p \rrbracket_\gamma(\text{Pr}_{\mathcal{D}}) + (1 + \gamma) H(\text{Pr}_{\mathcal{D}}) &= D(\text{Pr}_{\mathcal{D}} \parallel p) - \gamma H(\text{Pr}_{\mathcal{D}}) + (1 + \gamma) H(\text{Pr}_{\mathcal{D}}) \\ &= D(\text{Pr}_{\mathcal{D}} \parallel p) + H(\text{Pr}_{\mathcal{D}}) \\ &= \mathbb{E}_{\text{Pr}_{\mathcal{D}}} \left[ \log \frac{\text{Pr}_{\mathcal{D}}}{p} + \log \frac{1}{\text{Pr}_{\mathcal{D}}} \right] = \mathbb{E}_{\text{Pr}_{\mathcal{D}}} \left[ \log \frac{1}{p} \right], \end{aligned}$$

which is the cross entropy, as desired.  $\square$

**Proposition 4.** If  $p(X, Z)$  is a joint distribution, the information content of the partial observation  $X = x$  is given by

$$I_p[X=x] = \left\langle \boxed{Z} \xrightarrow{p} \boxed{X} \xleftarrow{x} \right\rangle. \quad (2)$$

<sup>7</sup>All proofs can be found in [Appendix C](#).

*Proof.* As before, all mass of  $\mu$  must be on  $x$  for it to have a finite score. Thus it suffices to consider joint distributions of the form  $\mu(X, Z) = \delta_x(X)\mu(Z)$ . We have

$$\begin{aligned}
 \left\langle \begin{array}{c} \text{Z} \xleftarrow{p} \text{X} \xleftarrow{x} \end{array} \right\rangle &= \inf_{\mu(Z)} \left[ \begin{array}{c} \text{Z} \xleftarrow{p} \text{X} \xleftarrow{x} \end{array} \right] \left( \delta_x(X)\mu(Z) \right) \\
 &= \inf_{\mu(Z)} D\left( \delta_x(X)\mu(Z) \parallel p(X, Z) \right) \\
 &= \inf_{\mu(Z)} \mathbb{E}_{z \sim \mu} \log \frac{\mu(z)}{p(x, z)} = \inf_{\mu(Z)} \mathbb{E}_{z \sim \mu} \log \frac{\mu(z)}{p(x, z)} \frac{p(x)}{p(x)} \\
 &= \inf_{\mu(Z)} \mathbb{E}_{z \sim \mu} \left[ \log \frac{\mu(z)}{p(z | x)} + \log \frac{1}{p(x)} \right] \\
 &= \inf_{\mu(Z)} \left[ D(\mu(Z) \parallel p(Z | x)) \right] + \log \frac{1}{p(x)} \\
 &= \log \frac{1}{p(x)} = I_p(x) \quad \text{[Gibbs Inequality]}
 \end{aligned}$$

□

**Proposition 19.** *The average negative log likelihood  $\ell(p; x) := -\frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \log p(x)$  (which is also the cross entropy) is the inconsistency of the PDG containing  $p$  and the data distribution  $\text{Pr}_{\mathcal{D}}$ , plus the entropy of the data distribution (which is constant in  $p$ ). That is,*

$$\ell(p; \mathcal{D}) = \left\langle \begin{array}{c} \text{Z} \xleftarrow{p} \text{X} \xleftarrow{\text{Pr}_{\mathcal{D}}!} \end{array} \right\rangle + H(\text{Pr}_{\mathcal{D}}).$$

*Proof.* The same idea as in Proposition 4, but a little more complicated.

$$\begin{aligned}
 \left\langle \begin{array}{c} \text{Z} \xleftarrow{p} \text{X} \xleftarrow{\text{Pr}_{\mathcal{D}}!} \end{array} \right\rangle &= \inf_{\mu(Z|X)} \left[ \begin{array}{c} \text{Z} \xleftarrow{p} \text{X} \xleftarrow{\text{Pr}_{\mathcal{D}}!} \end{array} \right] \left( \text{Pr}_{\mathcal{D}}(X)\mu(Z | X) \right) \\
 &= \inf_{\mu(Z|X)} D\left( \text{Pr}_{\mathcal{D}}(X)\mu(Z | X) \parallel p(X, Z) \right) \\
 &= \inf_{\mu(Z|X)} \mathbb{E}_{x \sim \text{Pr}_{\mathcal{D}}} \log \frac{\mu(z | x) \text{Pr}_{\mathcal{D}}(x)}{p(x, z)} \\
 &= \frac{1}{|\mathcal{D}|} \inf_{\mu(Z|X)} \sum_{x \in \mathcal{D}} \mathbb{E}_{z \sim \mu(Z|x)} \log \frac{\mu(z | x) \text{Pr}_{\mathcal{D}}(x)}{p(x, z)} \frac{p(x)}{p(x)} \\
 &= \frac{1}{|\mathcal{D}|} \inf_{\mu(Z|X)} \sum_{x \in \mathcal{D}} \mathbb{E}_{z \sim \mu} \left[ \log \frac{\mu(z)}{p(z | x)} + \log \frac{1}{p(x)} - \log \frac{1}{\text{Pr}_{\mathcal{D}}(x)} \right] \\
 &= \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \left[ \inf_{\mu(Z)} \left[ D(\mu(Z) \parallel p(Z | x)) \right] + \log \frac{1}{p(x)} \right] - H(\text{Pr}_{\mathcal{D}}) \\
 &= \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \log \frac{1}{p(x)} - H(\text{Pr}_{\mathcal{D}}) = \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} I_p(x) - H(\text{Pr}_{\mathcal{D}}) \\
 &= D(\text{Pr}_{\mathcal{D}} \parallel p(X))
 \end{aligned}$$

□

**Proposition 5.** *The inconsistency of the PDG comprising a probabilistic predictor  $h(Y|X)$ , and a high-confidence empirical distribution  $\text{Pr}_{\mathcal{D}}$  of a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^m$  is equal to the cross-entropy loss (minus the empirical uncertainty in  $Y$  given  $X$ , a constant that depends only on  $\mathcal{D}$ ). That is,*

$$\left\langle\left\langle \begin{array}{c} \text{Pr}_{\mathcal{D}} \text{ }^{(\infty)} \\ \swarrow \quad \searrow \\ \boxed{X} \xrightarrow{h} \boxed{Y} \end{array} \right\rangle\right\rangle = \frac{1}{m} \sum_{i=1}^m \log \frac{1}{h(y_i | x_i)} - H_{\text{Pr}_{\mathcal{D}}}(Y|X).$$

*Proof.*  $\text{Pr}_{\mathcal{D}}$  has high confidence, it is the only joint distribution  $\mu$  with finite score. Since  $f$  is the only other edge, the inconsistency is therefore

$$\begin{aligned} \mathbb{E}_{x \sim \text{Pr}_{\mathcal{D}}} D(\text{Pr}_{\mathcal{D}}(Y | x) \parallel f(Y | x)) &= \mathbb{E}_{x, y \sim \text{Pr}_{\mathcal{D}}} \left[ \log \frac{\text{Pr}_{\mathcal{D}}(y | x)}{f(y | x)} \right] \\ &= \mathbb{E}_{x, y \sim \text{Pr}_{\mathcal{D}}} \left[ \log \frac{1}{f(y | x)} - \log \frac{1}{\text{Pr}_{\mathcal{D}}(y | x)} \right] \\ &= \frac{1}{|\mathcal{D}|} \sum_{(x, y) \in \mathcal{D}} \left[ \log \frac{1}{f(y | x)} \right] - H_{\text{Pr}_{\mathcal{D}}}(Y | X) \end{aligned}$$

□

**Proposition 6.** *Consider functions  $f, h : X \rightarrow Y$  from inputs to labels, where  $h$  is a predictor and  $f$  generates the true labels. The inconsistency of believing  $f$  and  $h$  (with any confidences), and a distribution  $D(X)$  with confidence  $\beta$ , is  $\beta$  times the log accuracy of  $h$ . That is,*

$$\left\langle\left\langle \begin{array}{c} \xrightarrow{D} \\ (\beta) \end{array} \boxed{X} \begin{array}{c} \xrightarrow{h} \text{ }^{(r)} \\ \xrightarrow{f} \text{ }_{(s)} \end{array} \boxed{Y} \right\rangle\right\rangle = -\beta \log \Pr_{x \sim D}(f(x) = h(x)) = \beta I_D[f = h]. \quad (3)$$

*Proof.* Because  $f$  is deterministic, for every  $x$  in the support of a joint distribution  $\mu$  with finite score, we must have  $\mu(Y | x) = \delta_{f(x)}$ , since if  $\mu$  were to place any non-zero mass  $\mu(x, y) = \epsilon > 0$  on a point  $(x, y)$  with  $y \neq f(x)$  results in an infinite contribution to the KL divergence

$$D(\mu(Y | x) \parallel \delta_{f(x)}) = \mathbb{E}_{x, y \sim \mu} \log \frac{\mu(y | x)}{\delta_{f(x)}} \geq \mu(y, x) \log \frac{\mu(x, y)}{\mu(x) \cdot \delta_{f(x)}(y)} = \epsilon \log \frac{\epsilon}{0} = \infty.$$

The same holds for  $h$ . Therefore, for any  $\mu$  with a finite score, and  $x$  with  $\mu(x) > 0$ , we have  $\delta_{f(x)} = \mu(Y | x) = \delta_{h(x)}$ , meaning that we need only consider  $\mu$  whose support is a subset of those points on which  $f$  and  $h$  agree. On all such points, the contribution to the score from the edges associated to  $f$  and  $h$  will be zero, since  $\mu$  matches the conditional marginals exactly, and the total incompatibility of such a distribution  $\mu$  is equal to the relative entropy  $D(\mu \parallel D)$ , scaled by the confidence  $\beta$  of the empirical distribution  $D$ .

So, among those distributions  $\mu(X)$  supported on an event  $E \subset \mathcal{V}(X)$ , which minimizes is the relative entropy of  $D(\mu \parallel D)$ ? It is well known that the conditional distribution  $D \mid E \propto \delta_E(X)D(X) = \frac{1}{D(E)}\delta_E(X)D(X)$  satisfies this property uniquely (see, for instance, **halpernRAU**). Let  $f=h$  denote the event that  $f$  and  $h$  agree. Then

we calculate

$$\begin{aligned}
 \left\langle\!\left\langle \begin{array}{c} \xrightarrow{(\beta) D} X \xrightleftharpoons[f]{h} Y \end{array} \right\rangle\!\right\rangle &= \inf_{\substack{\mu(X) \text{ s.t.} \\ \text{supp}(\mu) \subseteq [f=h]}} \beta \mathbf{D}(\mu(X) \parallel D(X)) \\
 &= \beta \mathbf{D}(D \mid [f=h] \parallel D) \\
 &= \beta \mathbb{E}_{D|f=h} \log \frac{\delta_{f=h}(X) D(X)}{D(f=h) \cdot D(X)} \\
 &= \beta \mathbb{E}_{D|f=h} \log \frac{1}{D(f=h)} \quad \left[ \begin{array}{l} \text{since } \delta_{f=h}(x) = 1 \text{ for all } x \text{ that} \\ \text{contribute to the expectation} \end{array} \right] \\
 &= -\beta \log D(f=h) \quad \left[ \begin{array}{l} \text{since } D(f=h) \text{ is a constant} \end{array} \right] \\
 &= -\beta \log (\text{accuracy}_{f,D}(h)) \\
 &= \beta \mathbf{I}_D[f=h].
 \end{aligned}$$

□

**Proposition 15.** Consider a PDG containing two (distinct) conditional Gaussian distributions on a variable  $Y$ , whose parameters can both depend on a variable  $X$ . Its inconsistency takes the form

$$\begin{aligned}
 \left\langle\!\left\langle \begin{array}{c} \xrightarrow{D!} X \xrightleftharpoons[\sigma_2]{\sigma_1} Y \end{array} \right\rangle\!\right\rangle &= \frac{1}{2} \mathbb{E}_D \left[ \text{HM}(\beta_1, \beta_2) \frac{1}{2} \left( \frac{\mu_1 - \mu_2}{\text{QM}_{\hat{\beta}}(\sigma_1, \sigma_2)} \right)^2 + \text{AM}(\beta_1, \beta_2) \log \frac{\text{QM}_{\hat{\beta}}(\sigma_1, \sigma_2)}{\text{GM}_{\hat{\beta}}(\sigma_1, \sigma_2)} \right] \quad (7) \\
 &= \mathbb{E}_{x \sim D} \left[ \frac{\beta_1 \beta_2}{2} \frac{(f(x) - h(x))^2}{\beta_2 s(x)^2 + \beta_1 t(x)^2} + \frac{\beta_1 + \beta_2}{2} \log \frac{\beta_2 s(x)^2 + \beta_1 t(x)^2}{(\beta_1 + \beta_2)(s(x)^{\beta_2} t(x)^{\beta_1})^{\frac{1}{\beta_1 + \beta_2}}} \right]
 \end{aligned}$$

where  $\hat{\beta} = (\frac{\beta_2}{\beta_1 + \beta_2}, \frac{\beta_1}{\beta_1 + \beta_2})$  represents the normalized and reversed vector of confidences  $\beta = (\beta_1, \beta_2)$  for the two distributions, and  $\mu_1 = f(X)$ ,  $\mu_2 = g(X)$ ,  $\sigma_1 = s(X)$ ,  $\sigma_2 = t(X)$  are random variables over  $X$ .

*Proof.* Let  $\mathcal{M}$  denote the PDG in question. Since  $D$  has high confidence, we know any joint distribution  $\mu$  with a finite score must have  $\mu(X) = D(X)$ . Thus,

$$\begin{aligned}
 \langle \mathcal{M} \rangle_0 &= \inf_{\mu} \mathbb{E}_{x \sim D} \mathbb{E}_{y \sim \mu|x} \left[ \beta_1 \log \frac{\mu(y|x)}{\mathcal{N}(y|f(x), t(x))} + \beta_2 \log \frac{\mu(y|x)}{\mathcal{N}(y|h(x), s(x))} \right] \\
 &= \inf_{\mu} \mathbb{E}_{x \sim D} \mathbb{E}_{y \sim \mu|x} \left[ \beta_1 \log \frac{\mu(y|x)}{\frac{1}{t(x)\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{y-f(x)}{t(x)}\right)^2\right)} + \beta_2 \log \frac{\mu(y|x)}{\frac{1}{s(x)\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{y-h(x)}{s(x)}\right)^2\right)} \right] \\
 &= \inf_{\mu} \mathbb{E}_{x \sim D} \mathbb{E}_{y \sim \mu|x} \left[ \log \mu(y|x)^{\beta_1 + \beta_2} + \beta_1 \log(t(x)\sqrt{2\pi}) + \frac{\beta_1}{2} \left(\frac{y-f(x)}{t(x)}\right)^2 + \beta_2 \log(t(x)\sqrt{2\pi}) + \frac{\beta_2}{2} \left(\frac{y-h(x)}{s(x)}\right)^2 \right]
 \end{aligned}$$

⟨ FINISH PROOF ⟩

□



**Lemma 10.** *The inconsistency of a PDG comprising  $p(X)$  with confidence  $r$  and  $q(X)$  with confidence  $s$  is given in closed form by*

$$\left\langle \left\langle \frac{p}{(r)} \rightarrow X \leftarrow \frac{q}{(s)} \right\rangle \right\rangle = -(r+s) \log \sum_x \left( p(x)^r q(x)^s \right)^{\frac{1}{r+s}}.$$

*Proof.*

$$\begin{aligned} \left\langle \left\langle \frac{p}{(\beta:r)} \rightarrow X \leftarrow \frac{q}{(\beta:s)} \right\rangle \right\rangle &= \inf_{\mu} \mathbb{E}_{\mu} \log \frac{\mu(x)^{r+s}}{p(x)^r q(x)^s} \\ &= (r+s) \inf_{\mu} \mathbb{E}_{\mu} \left[ \log \frac{\mu(x)}{(p(x)^r q(x)^s)^{\frac{1}{r+s}}} \cdot \frac{Z}{Z} \right] \\ &= \inf_{\mu} (r+s) D \left( \mu \parallel \frac{1}{Z} p^{\frac{r}{r+s}} q^{\frac{s}{r+s}} \right) - (r+s) \log Z \end{aligned}$$

where  $Z := (\sum_x p(x)^r q(x)^s)^{\frac{1}{r+s}}$  is the constant required to normalize the denominator as a distribution. Since this is now a relative entropy, it achieves its minimum when  $\mu$  is the other distribution, at which point it contributes zero, so our formula becomes

$$\begin{aligned} &= -(r+s) \log Z \\ &= -(r+s) \log \sum_x \left( p(x)^r q(x)^s \right)^{\frac{1}{r+s}} \quad \text{as promised.} \end{aligned}$$

□

**Proposition 8.** *Suppose you have a parameterized model  $p(Y|\Theta)$ , a prior  $q(\Theta)$ , and a trusted distribution  $D(Y)$ . The inconsistency of also believing  $\Theta = \theta$  is the cross entropy loss, plus the regularizer:  $\log \frac{1}{q(\theta)}$  times your confidence in  $q$ . That is,*

$$\left\langle \left\langle \begin{array}{c} q \\ (\beta) \end{array} \rightarrow \Theta \xrightarrow{p} Y \right\rangle \right\rangle_{D \uparrow_{(\infty)}} = \mathbb{E}_{y \sim D} \log \frac{1}{p(y|\theta)} + \beta \log \frac{1}{q(\theta)} - H(D) \quad (4)$$

*Proof.* **FINISH PROOF**

□

**Proposition 11.** *The negative ELBO of  $x$  is the inconsistency of the PDG containing  $p, q$ , and  $X=x$ , with high confidence in  $q$ . That is,*

$$-\text{ELBO}_{p,q}(x) = \left\langle \left\langle \frac{q}{(\infty)} \rightarrow Z \xrightarrow{p} X \leftarrow x \right\rangle \right\rangle.$$

*Proof.* Every distribution that does marginalize to  $q(Z)$  or places any mass on  $x' \neq x$  will have infinite score. Thus the only distribution that could have a finite score is  $\mu(X, Z)$ . Thus,

$$\begin{aligned}
 \left\langle \left\langle \begin{array}{c} q! \\ \rightarrow \\ Z \end{array} \begin{array}{c} \swarrow^p \\ X \end{array} \leftarrow x \right\rangle \right\rangle &= \inf_{\mu} \left[ \left\langle \begin{array}{c} q! \\ \rightarrow \\ Z \end{array} \begin{array}{c} \swarrow^p \\ X \end{array} \leftarrow x \right\rangle \right] (\mu) \\
 &= \left[ \left\langle \begin{array}{c} q! \\ \rightarrow \\ Z \end{array} \begin{array}{c} \swarrow^p \\ X \end{array} \leftarrow x \right\rangle \right] (\delta_x(X)q(Z)) \\
 &= \mathbb{E}_{\substack{x' \sim \delta_x \\ z \sim q}} \log \frac{\delta_x(x')q(z)}{p(x', z)} = - \mathbb{E}_{z \sim q} \frac{p(x, z)}{q(z)} = -\text{ELBO}_{p,q}(x).
 \end{aligned}$$

□

We prove both [Proposition 12](#) and [Proposition 16](#) at the same time.

**Proposition 12.** *The VAE loss of a sample  $x$  is the inconsistency of the PDG comprising the encoder (with high confidence, as it defines the encoding), decoder  $d$ , prior  $p$ , and  $x$ . That is,*

$$-\text{ELBO}_{p,e,d}(x) = \left\langle \begin{array}{c} p \\ \rightarrow \\ Z \end{array} \begin{array}{c} \xrightarrow{d} \\ X \end{array} \begin{array}{c} \xleftarrow{e} \\ x \end{array} \right\rangle.$$

**Proposition 16.** *The following analog of [Proposition 12](#) for a whole dataset  $\mathcal{D}$  holds:*

$$-\mathbb{E}_{\text{Pr}_{\mathcal{D}}} \text{ELBO}_{p,e,d}(X) = \left\langle \begin{array}{c} p \\ \rightarrow \\ Z \end{array} \begin{array}{c} \xrightarrow{d} \\ X \end{array} \begin{array}{c} \xleftarrow{e!} \\ \text{Pr}_{\mathcal{D}}! \end{array} \right\rangle + H(\text{Pr}_{\mathcal{D}}).$$

*Proof.* The two proofs are similar. For [Proposition 12](#), the optimal distribution must be  $\delta_x(X)e(Z | X)$ , and for [Proposition 16](#), it must be  $\text{Pr}_{\mathcal{D}}(X)e(Z | X)$ , because  $e$  and the data both have infinite confidence, so any other distribution gets an infinite score. At the same time,  $d$  and  $p$  define a joint distribution, so the inconsistency in question becomes

$$D\left(\delta_x(X)e(Z | X) \parallel p(Z)d(X | Z)\right) = \mathbb{E}_{z \sim e|x} \left[ \log \frac{p(z)d(x | z)}{e(z | x)} \right] = \text{ELBO}_{p,e,d}(x)$$

in the first, case, and

$$\begin{aligned}
 D\left(\text{Pr}_{\mathcal{D}}(X)e(Z | X) \parallel p(Z)d(X | Z)\right) &= \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \mathbb{E}_{z \sim e|x} \left[ \log \frac{p(z)d(x | z)}{e(z | x)} + \log \frac{1}{\text{Pr}_{\mathcal{D}}(x)} \right] \\
 &= \text{ELBO}_{p,e,d}(x) - H(\text{Pr}_{\mathcal{D}})
 \end{aligned}$$

in the second.

□

Now, we formally state and prove the more general result for  $\beta$ -VAEs.

**Proposition 17.** *The negative  $\beta$ -ELBO objective for a prior  $p(X)$ , encoder  $e(Z | X)$ , decoder  $d(X | Z)$ , at a sample  $x$ , is equal to the inconsistency of the corresponding PDG, where the prior has confidence equal to  $\beta$ . That is,*

$$-\beta\text{-ELBO}_{p,e,d}(x) = \left\langle \begin{array}{c} \begin{array}{c} (\beta) \\ p \end{array} \\ \rightarrow \\ Z \end{array} \begin{array}{c} \xrightarrow{d} \\ X \end{array} \begin{array}{c} \xleftarrow{e!} \\ x \end{array} \right\rangle$$

*Proof.*

$$\begin{aligned}
 \left\langle \left( \begin{array}{c} \xrightarrow{(\beta; \beta')} \\ \text{ } \\ \xleftarrow{e!} \end{array} \right) \begin{array}{c} Z \\ \xrightarrow{d} \\ X \end{array} \xleftarrow{x} \right\rangle &= \inf_{\mu} \left[ \left( \begin{array}{c} \xrightarrow{(\beta; \beta')} \\ \text{ } \\ \xleftarrow{e!} \end{array} \right) \begin{array}{c} Z \\ \xrightarrow{d} \\ X \end{array} \xleftarrow{x} \right] (\mu) \\
 &= \inf_{\mu} \mathbb{E}_{\mu(X, Z)} \left[ \beta \log \frac{\mu(Z)}{p(Z)} + \log \frac{\mu(X, Z)}{\mu(Z)d(X | Z)} \right]
 \end{aligned}$$

As before, the only candidate for a joint distribution with finite score is  $\delta_x(X)e(Z | X)$ . Note that the marginal on  $Z$  for this distribution is itself, since  $\int_x \delta_x(X)e(Z | X) dx = e(Z | x)$ . Thus, our equation becomes

$$\begin{aligned}
 &= \mathbb{E}_{\delta_x(X)e(Z | X)} \left[ \beta \log \frac{e(Z | x)}{p(z)} + \log \frac{\delta_x(X)e(Z | X)}{e(Z | x)d(x | Z)} \right] \\
 &= \mathbb{E}_{e(Z | x)} \left[ \beta \log \frac{e(Z | x)}{p(Z)} + \log \frac{1}{d(x | Z)} \right] = -\beta \text{ELBO}_{p, e, d}(x).
 \end{aligned}$$

□

**Proposition 13.** *For any weighted factor graph  $\Psi$  we have  $\langle \mathbf{m}_{\Psi} \rangle_1 = -\log Z_{\Psi}$ .*

*Proof.* Let the  $(\{x\}) := x$  be a function that extracts the unique element singleton set. We showed in the original paper (Corollary 4.4.1) that

$$\text{the} \langle (\mathbf{n}_{\Phi}, \theta, \theta) \rangle_1^* = \Pr_{\Phi, \theta}(\mathbf{w}) = \frac{1}{Z_{\Psi}} \prod_j \phi_j(\mathbf{w}_j)^{\theta_j}.$$

Recall the statement of Prop 4.6 from the original paper,

$$\langle \mathbf{m} \rangle_{\gamma}(\mu) = \mathbb{E}_{\mathbf{w} \sim \mu} \left\{ \sum_{X \xrightarrow{L} Y} \left[ \beta_L \log \frac{1}{\mathbf{p}_L(y^{\mathbf{w}} | x^{\mathbf{w}})} + (\gamma \alpha_L - \beta_L) \log \frac{1}{\mu(y^{\mathbf{w}} | x^{\mathbf{w}})} \right] - \gamma \log \frac{1}{\mu(\mathbf{w})} \right\}, \quad (8)$$

but note that since  $\gamma = 1$ , and  $\alpha, \beta$  are both equal to  $\theta$  for our PDG (since  $\mathbf{m}_{\Psi} = \mathbf{m}_{(\Phi, \theta)} = (\mathbf{n}_{\Phi}, \theta, \theta)$ ), the middle term disappears, yielding the standard variational free energy  $VFE(\mu)$ . Recall also that  $\langle \mathbf{m} \rangle_{\gamma} = \inf_{\mu} \langle \mathbf{m} \rangle_{\gamma}(\mu)$  and  $\langle \mathbf{m} \rangle_{\gamma}^* = \arg \min \langle \mathbf{m} \rangle_{\gamma}(\mu)$ , so (with a minor abuse of notation),  $\langle \mathbf{m} \rangle_{\gamma} = \langle \mathbf{m} \rangle_{\gamma}(\langle \mathbf{m} \rangle_{\gamma}^*)$ . We now compute the value of the inconsistency  $\langle (\mathbf{n}_{\Phi}, \theta, \theta) \rangle_1$ .

$$\begin{aligned}
 \langle (\mathbf{n}_{\Phi}, \theta, \theta) \rangle_1 &= \langle (\mathbf{n}_{\Phi}, \theta, \theta) \rangle_1 \left( \Pr_{\Phi, \theta}(\mathbf{w}) \right) \\
 &= \mathbb{E}_{\mathbf{w} \sim \mu} \left\{ \sum_{X \xrightarrow{L} Y} \left[ \beta_L \log \frac{1}{\mathbf{p}_L(y^{\mathbf{w}} | x^{\mathbf{w}})} \right] - \log \frac{1}{\Pr_{\Phi, \theta}(\mathbf{w})} \right\} \quad \left[ \text{by (8)} \right] \\
 &= \mathbb{E}_{\mathbf{w} \sim \mu} \left\{ \sum_j \left[ \theta_j \log \frac{1}{\phi_j(\mathbf{w}_j)} \right] - \log \frac{Z_{\Psi}}{\prod_j \phi_j(\mathbf{w}_j)^{\theta_j}} \right\} \quad \left[ \begin{array}{l} \text{cpds } \mathbf{p}_L \text{ correspond} \\ \text{to factors } \phi_j \end{array} \right] \\
 &= \mathbb{E}_{\mathbf{w} \sim \mu} \left\{ \sum_j \left[ \theta_j \log \frac{1}{\phi_j(\mathbf{w}_j)} \right] - \sum_j \left[ \theta_j \log \frac{1}{\phi_j(\mathbf{w}_j)} \right] - \log Z_{\Psi} \right\} \\
 &= \mathbb{E}_{\mathbf{w} \sim \mu} [-\log Z_{\Psi}] \\
 &= -\log Z_{\Psi} \quad \left[ Z_{\Psi} \text{ is constant in } \mathbf{w} \right]
 \end{aligned}$$

□

## D More Notes

### D.1 Surprise

A common justification for using  $I_p(x)$  as a cost for updating a probabilistic model  $p(x)$  based on an observed sample  $x$ , is that by minimizing it, you “maximize the probability of seeing your data”.<sup>8</sup> But this explanation applies just as well to  $-p(x)$ . Why include the logarithm? There are plenty of answers to this question; among them:  $I_p$  is convex in  $p$ , it decomposes products into arguably simpler sums, is more numerically stable, has a well-defended physical analogue in thermodynamics, and is a primitive of information theory.

For those after a quick and rigorous justification (as opposed to handwaving or a thermodynamics textbook), none of these answers are entirely satisfying. They suggest that  $I_p$  has certain nice properties, but not that it enjoys them uniquely, or that no other loss function satisfies nicer ones. Pedagogically speaking, the situation is more straightforward for us. Although PDG semantics themselves require non-trivial justification, they give us in return uniform answers to many questions, starting with: Why use the surprise  $I_p(x)$ , to measure the loss of a model  $p(X)$  on sample  $x$ ? Because it is the inconsistency of simultaneously believing  $X = x$  and  $X \sim p$ .

---

<sup>8</sup>this justification should not be taken too seriously without constraints on  $p$ , because the optimal value of  $p$  is  $\delta_x$ , which does not generalize.