# Data Processing Inequalities

#### TL;DR

Informally, the **Data Processing Inequality (DPI)** states that processing data stochastically can only reduce information. Formally, for distributions  $q(\mathbf{\Theta})$  and  $p(\mathbf{\Theta})$  over a random variable  $\mathbf{\Theta}$  and a stochastic mapping  $Y = f(\mathbf{\Theta})$ , the Kullback-Leibler DPI states:

$$D_{KL}(q(\boldsymbol{\Theta}) \parallel p(\boldsymbol{\Theta})) \geq D_{KL}(q(Y) \parallel p(Y)).$$

Equality holds when  $D_{KL}(q(\boldsymbol{\Theta} \mid Y) \parallel p(\boldsymbol{\Theta} \mid Y)) = 0$ .

The data processing inequality states that if two Example: Image Processing random variables are transformed in this way, they cannot become easier to tell apart.

"Understanding Variational Inference in Function-Space".

Consider an image processing pipeline where X is the original image, Y is a compressed version, and Z is Y after adding blur and pixelation. The DPI tells us that the mutual information Burt et al. (2021)  $I[X;Y] \geq I[X;Z]$ , as each processing step results in information loss.

## Jenson-Shannon Divergence DPI

The Jensen-Shannon divergence (JSD) makes the KL divergence symmetric. For:

$$f(x) = \frac{p(x) + q(x)}{2}$$

$$D_{JSD}(p(x) \parallel q(x)) = \frac{1}{2} D_{KL}(p(x) \parallel f(x)) + \frac{1}{2} D_{KL}(q(x) \parallel f(x)).$$

The square root of the Jensen-Shannon divergence, the *Jensen*-Shannon distance, is symmetric, satisfies the triangle inequality, and is hence a metric.

For p(x) and q(x) and a shared transition function  $f(y \mid x)$  for the model  $X \to Y$ , we apply the KL DPI twice and obtain the JSD DPI:

$$D_{JSD}(p(X) \parallel q(X)) \ge D_{JSD}(p(Y) \parallel q(Y)).$$

## Chain Rule of the P Divergence

An important property of the KL divergence is the chain rule:

$$D_{KL}(q(Y_{n},...) \parallel p(Y_{n},...))$$

$$= D_{KL}(q(Y_{n} \mid Y_{n-1},...) \parallel p(Y_{n} \mid Y_{n-1},...))$$

$$+ D_{KL}(q(Y_{n-1},...) \parallel p(Y_{n-1},...))$$

$$= ... = \sum_{i=1}^{n} D_{KL}(q(Y_{i} \mid Y_{i-1},...) \parallel p(Y_{i} \mid Y_{i-1},...)).$$

This also yields a **chain inequality**:

$$D_{KL}(q(Y_n, ...) \parallel p(Y_n, ...)) \ge D_{KL}(q(Y_{n-1}, ...) \parallel p(Y_{n-1}, ...))$$
  
  $\ge ... \ge D_{KL}(q(Y_1) \parallel p(Y_1)),$ 

by repeatedly applying the chain rule and dropping the conditional term.

# Mutual Information DPI

For any Markov chain  $Z \rightarrow X \rightarrow Y$  with f(z, x, y) = $f(z)f(x \mid z)f(y \mid x)$  for any distribution f(z), we have:

$$I[X; Z] = D_{KL}(f(X \mid Z) \parallel f(X))$$

$$= \mathbb{E}_{f(z)} [D_{KL}(f(X \mid z) \parallel f(X))]$$

$$\stackrel{(1)}{\geq} \mathbb{E}_{f(z)} [D_{KL}(f(Y \mid z) \parallel f(Y))]$$

$$= D_{KL}(f(Y \mid Z) \parallel f(Y))$$

$$= I[Y; Z],$$

where (1) follows from the KL DPI.

#### Proof of the PDPI

Using the chain rule of the KL divergence twice:

$$D_{KL}(p(X) || q(X)) + \underbrace{D_{KL}(p(Y | X) || q(Y | X))}_{=D_{KL}(f(Y | X) || f(Y | X))=0}$$

$$= D_{KL}(p(X, Y) || q(X, Y))$$

$$= D_{KL}(p(Y) || q(Y)) + \underbrace{D_{KL}(p(X | Y) || q(X | Y))}_{\geq 0}$$

We have equality exactly when  $p(x \mid y) = q(x \mid y)$  for (almost) | all x, y. |

# Full Blog Post



Burt, D. R., Ober, S. W., Garriga-Alonso, A., and van der Wilk, M. (2021). Understanding variational inference in function-space. In Third Symposium on Advances in Approximate Bayesian Inference.

 $\geq D_{KL}(p(Y) \parallel q(Y)).$ 

Cover, T. M. (1999). Elements of information theory. John Wiley & Sons.

Rudner, T. G. J., Chen, Z., Teh, Y. W., and Gal, Y. (2022). Tractable function-space variational inference in bayesian neural networks. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K., editors, Advances in Neural Information Processing Systems.

# Function-Space Variational Inference

## TL;DR

Function-space variational inference (FSVI) is a principled approach to Bayesian inference that respects the inherent symmetries and equivalences in overparameterized models. It focuses on approximating the meaningful posterior  $p(|\boldsymbol{\theta}| | \mathcal{D})$  over prediction equivalence classes of the parameters while avoiding the complexities of explicitly constructing and working with equivalence classes. The FSVI-ELBO regularizes towards a data prior using the KL DPI:

$$\mathbb{E}_{q(\boldsymbol{\theta})}\left[-\log p(\mathcal{D} \mid \boldsymbol{\theta})\right] + D_{KL}(q(Y... \mid \boldsymbol{x}...) \parallel p(Y... \mid \boldsymbol{x}...)),$$

(unlike in regular variational inference, where we regularize towards a parameter prior  $D_{KL}(q(\boldsymbol{\Theta}) \parallel p(\boldsymbol{\Theta}))$ ).

## (Regular) Variational Inference & ELBO

The Bayesian posterior  $p(\boldsymbol{\theta} \mid \mathcal{D})$  is approximated with a variational distribution  $q(\boldsymbol{\theta})$  by minimizing  $D_{KL}(q(\boldsymbol{\Theta}) \parallel p(\boldsymbol{\Theta} \mid \mathcal{D}))$ . Dropping constant (intractable) terms yields a simplified tractable objective, which **upper** bounds the information content  $-\log p(\mathcal{D})$  of the data  $\mathcal{D}$ :

$$0 \leq D_{KL}(q(\boldsymbol{\Theta}) \parallel p(\boldsymbol{\Theta} \mid \mathcal{D}))$$

$$= D_{KL}(q(\boldsymbol{\Theta}) \parallel \frac{p(\mathcal{D} \mid \boldsymbol{\Theta}) p(\boldsymbol{\Theta})}{p(\mathcal{D})})$$

$$= \mathbb{E}_{q} \left[ -\log p(\mathcal{D} \mid \boldsymbol{\Theta}) \right] + D_{KL}(q(\boldsymbol{\Theta}) \parallel p(\boldsymbol{\Theta}))$$
Evidence Bound (Simplified Objective)
$$- \left( -\log p(\mathcal{D}) \right).$$
(peg. log) Evidence

This is equivalent to the **evidence lower bound (ELBO)**.

#### Parameter Symmetries

Deep neural networks have many parameter symmetries: for example, in a convolutional neural network, we could swap channels without changing the predictions.  $\implies$  We are not interested in these symmetries but only differing predictions.

#### Equivalence Classes

Equivalence classes group together parameters that lead to the same predictions on a (test) set of data:

$$[\boldsymbol{\theta}] \triangleq \{\boldsymbol{\theta}' : f(x; \boldsymbol{\theta}) = f(x; \boldsymbol{\theta}) \quad \forall x\}.$$

Crucially, different domains for  $\boldsymbol{x}$  will induce different equivalence classes.

#### Consistency of Equivalence Classes with Bayesian Inference

Any distribution over the parameters  $p(\boldsymbol{\theta})$  induces a distribution  $p([\boldsymbol{\theta}])$  over the equivalence classes, which is consistent with Bayesian inference:

$$p([\boldsymbol{\theta}]) \triangleq \sum_{\boldsymbol{\theta}' \in [\boldsymbol{\theta}]} p(\boldsymbol{\theta}'),$$

that is,  $[\boldsymbol{\theta}]$  commutes with Bayesian inference:

$$p([\boldsymbol{\theta}] \mid \mathcal{D}) = \sum_{\boldsymbol{\theta}' \in [\boldsymbol{\theta}]} p(\boldsymbol{\theta}' \mid \mathcal{D}) \Leftrightarrow [\boldsymbol{\Theta} \mid \mathcal{D}] = [\boldsymbol{\Theta}] \mid \mathcal{D}.$$

This commutative property is a general characteristic of applying (stochastic) functions to random variables.

$$\begin{array}{c} \mathbf{\Theta} & \stackrel{\cdot \mid \mathcal{D}}{\longrightarrow} \mathbf{\Theta} \mid \mathcal{D} \\ \downarrow [\cdot] & \downarrow [\cdot] \\ [\mathbf{\Theta}] & \stackrel{\cdot \mid \mathcal{D}}{\longrightarrow} [\mathbf{\Theta}] \mid \mathcal{D} \end{array}$$
Commutative diagram for the equivalence classes.

#### Equality in the Infinite Data Limit

$$D_{KL}(q(\boldsymbol{\Theta}) \parallel p(\boldsymbol{\Theta})) \geq D_{KL}(q([\boldsymbol{\Theta}]) \parallel p([\boldsymbol{\Theta}]))$$
$$\geq D_{KL}(q(Y_n... \mid \boldsymbol{x}_n...) \parallel p(Y_n... \mid \boldsymbol{x}_n...)).$$

Unless there are no parameter symmetries, the **first inequal**ity will not be tight  $(D_{KL}(q(\boldsymbol{\Theta} \mid [\boldsymbol{\Theta}]) \parallel p(\boldsymbol{\Theta} \mid [\boldsymbol{\Theta}])) > 0).$ The **second inequality will be tight**: to see this, we need  $D_{\mathrm{KL}}(q([\boldsymbol{\Theta}]|Y_n,\boldsymbol{x}_n,...)\|p([\boldsymbol{\Theta}]|Y_n,\boldsymbol{x}_n,...)) \to 0 \text{ for } n \to \infty.$  The term converges as it is monotonically increasing and bounded by  $D_{KL}(q([\boldsymbol{\Theta}]) \parallel p([\boldsymbol{\Theta}]))$  from above. Further, thanks to Bernstein von Mises' theorem, we have  $\rightarrow 0$  and thus:

$$D_{\mathrm{KL}}(\mathbf{q}([\boldsymbol{\Theta}]) \parallel \mathbf{p}([\boldsymbol{\Theta}])) =$$

$$= \sup_{n \in \mathbb{N}} D_{\mathrm{KL}}(\mathbf{q}(Y_n, \dots \mid \boldsymbol{x}_n, \dots) \parallel \mathbf{p}(Y_n, \dots \mid \boldsymbol{x}_n, \dots)).$$

#### Bernstein von Mises' Theorem

BvM states that a posterior distribution converges to the maximum likelihood estimate (MLE) as the number of data points tends to infinity as long as the model parameters are identifiable, that is the true parameters we want to learn are unique, and they have support. This is true for  $[\Theta]$  as  $[\Theta]$  is a unique parameterization of the same distribution as  $\Theta$ , and we can ensure support.

## Function-Space Variational Inference & ELBO

FSVI's ELBO is just the regular ELBO but for  $[\Theta]$  and uses further approximations via chain inequalities:

$$\begin{split} H[\mathcal{D}] &\leq H[\mathcal{D}] + D_{KL}(q([\boldsymbol{\Theta}]) \parallel p([\boldsymbol{\Theta}] \mid \mathcal{D})) \\ &= H[\mathcal{D}] + D_{KL}(q([\boldsymbol{\Theta}]) \parallel \frac{p(\mathcal{D} \mid [\boldsymbol{\Theta}]) p([\boldsymbol{\Theta}])}{p(\mathcal{D})}) \\ &= \mathbb{E}_{q([\boldsymbol{\theta}])} \left[ -\log p(\mathcal{D} \mid [\boldsymbol{\theta}]) \right] + D_{KL}(q([\boldsymbol{\Theta}]) \parallel p([\boldsymbol{\Theta}])). \end{split}$$

Then, we can apply the chain rule together with BvM:

$$= \mathbb{E}_{\mathbf{q}(\boldsymbol{\theta})} \left[ -\log \mathbf{p}(\boldsymbol{\mathcal{D}} \mid \boldsymbol{\theta}) \right] + \sup_{n} D_{\mathrm{KL}}(\mathbf{q}(\cdot) \parallel \mathbf{p}(\cdot)) [Y_n \dots \mid \boldsymbol{x}_n \dots]$$

$$\geq \mathbb{E}_{q(\boldsymbol{\theta})} \left[ -\log p(\mathcal{D} \mid \boldsymbol{\theta}) \right] + D_{KL}(q(\cdot) \parallel p(\cdot)) [Y_n ... \mid \boldsymbol{x}_n ...] \quad \forall n.$$