

BRIDGING THE DATA PROCESSING INEQUALITY AND FUNCTION-SPACE VARIATIONAL INFERENCE

Andreas Kirsch
University of Oxford²⁰²³

Data Processing Inequalities

TL;DR

Informally, the **Data Processing Inequality (DPI)** states that processing data stochastically can only reduce information. Formally, for distributions $q(\Theta)$ and $p(\Theta)$ over a random variable Θ and a stochastic mapping $Y = f(\Theta)$, the DPI is expressed as:

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

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

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

“Understanding Variational Inference in Function-Space”,
Burt et al. (2021)

Jensen-Shannon Divergence DPI

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

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

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

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

For $p(x)$ and $q(x)$ and shared transition function $f(y \mid x)$ for the model $X \rightarrow Y$:

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

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 $f(z)$:

$$\begin{aligned} I[X; Z] &= D_{\text{KL}}(f(X \mid Z) \parallel f(X)) \\ &= \mathbb{E}_{f(z)} [D_{\text{KL}}(f(X \mid z) \parallel f(X))] \\ &\stackrel{(1)}{\geq} \mathbb{E}_{f(z)} [D_{\text{KL}}(f(Y \mid z) \parallel f(Y))] \\ &= D_{\text{KL}}(f(Y \mid Z) \parallel f(Y)) \\ &= I[Y; Z], \end{aligned}$$

where (1) follows from the KL DPI.

Example: Image Processing

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 $I[X; Y] \geq I[X; Z]$, as each processing step results in information loss.

Chain Rule of the Divergence

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

$$\begin{aligned} D_{\text{KL}}(q(Y_n, \dots) \parallel p(Y_n, \dots)) \\ &= \sum_{i=1}^n D_{\text{KL}}(q(Y_i \mid Y_{i-1}, \dots) \parallel p(Y_i \mid Y_{i-1}, \dots)). \end{aligned}$$

This chain rule also yields a **chain inequality**:

$$\begin{aligned} D_{\text{KL}}(q(Y_n, \dots) \parallel p(Y_n, \dots)) &\geq D_{\text{KL}}(q(Y_{n-1}, \dots) \parallel p(Y_{n-1}, \dots)) \\ &\dots \\ &\geq D_{\text{KL}}(q(Y_1) \parallel p(Y_1)), \end{aligned}$$

where we start from the KL DPI and then use the chain rule.

Proof of the DPI

Using the chain rule of the KL divergence twice:

$$\begin{aligned} D_{\text{KL}}(p(X) \parallel q(X)) &+ \underbrace{D_{\text{KL}}(p(Y \mid X) \parallel q(Y \mid X))}_{=D_{\text{KL}}(f(Y \mid X) \parallel f(Y \mid X))=0} \\ &= D_{\text{KL}}(p(X, Y) \parallel q(X, Y)) \\ &= D_{\text{KL}}(p(Y) \parallel q(Y)) + \underbrace{D_{\text{KL}}(p(X \mid Y) \parallel q(X \mid Y))}_{\geq 0} \\ &\geq D_{\text{KL}}(p(Y) \parallel q(Y)). \end{aligned}$$

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

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 $p([\Theta] \mid \mathcal{D})$ while avoiding the complexities of explicitly constructing and working with equivalence classes. The FSVI-ELBO regularizes towards a data prior:

$$\mathbb{E}_{q(\Theta)} [-\log p(\mathcal{D} \mid \Theta)] + D_{\text{KL}}(q(Y_{\dots} \mid \mathbf{x}_{\dots}) \parallel p(Y_{\dots} \mid \mathbf{x}_{\dots})),$$

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

(Regular) Variational Inference & ELBO

In standard VI, we approximate a Bayesian posterior $p(\Theta \mid \mathcal{D})$ with a variational distribution $q(\Theta)$ by minimizing $D_{\text{KL}}(q(\Theta) \parallel p(\Theta \mid \mathcal{D}))$. This yields an information-theoretic evidence (**upper**) bound on the information content $-\log p(\mathcal{D})$ of the data \mathcal{D} under the variational distribution $q(\Theta)$:

$$\begin{aligned} 0 &\leq H(p(\mathcal{D})) + D_{\text{KL}}(q(\Theta) \parallel p(\Theta \mid \mathcal{D})) \\ &= H(p(\mathcal{D})) + D_{\text{KL}}(q(\Theta) \parallel \frac{p(\mathcal{D} \mid \Theta)p(\Theta)}{p(\mathcal{D})}) \\ &= \underbrace{\mathbb{E}_q [-\log p(\mathcal{D} \mid \Theta)] + D_{\text{KL}}(q(\Theta) \parallel p(\Theta))}_{\text{Evidence Bound}} \\ &\quad - \underbrace{(-\log p(\mathcal{D}))}_{\text{Evidence } H(p(\mathcal{D}))} \end{aligned}$$

In the literature, the negative of this bound is called the **evidence lower bound (ELBO)** and is maximized.

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 in the predictions.*

Equivalence Classes

We can use **equivalence classes** to group together parameters that lead to the same predictions on a (test) set of data:

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

Crucially, *different domains for \mathbf{x} will induce different equivalence classes*.

Consistency of Equivalence Classes with Bayesian Inference

Any distribution over the parameters $p(\Theta)$ induces a distribution $\hat{p}([\Theta])$ over the equivalence classes:

$$\hat{p}([\Theta]) \triangleq \sum_{\Theta' \in [\Theta]} p(\Theta').$$

$[\Theta]$ commutes with Bayesian inference:

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

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

$$\begin{array}{ccc} \Theta & \xrightarrow{\cdot \mid \mathcal{D}} & \Theta \mid \mathcal{D} \\ \downarrow [\cdot] & & \downarrow [\cdot] \\ [\Theta] & \xrightarrow{\cdot \mid \mathcal{D}} & [\Theta] \mid \mathcal{D} \end{array}$$

Equality in the Infinite Data Limit

Using the DPI:

$$\begin{aligned} D_{\text{KL}}(q(\Theta) \parallel p(\Theta)) &\geq D_{\text{KL}}(q([\Theta]) \parallel p([\Theta])) \\ &\geq D_{\text{KL}}(q(Y_{\dots} \mid \mathbf{x}_{\dots}) \parallel p(Y_{\dots} \mid \mathbf{x}_{\dots})). \end{aligned}$$

Unless there are no parameter symmetries, the **first inequality will not be tight**. For the second inequality to be tight, we need $D_{\text{KL}}(q([\Theta] \mid Y_n, \mathbf{x}_n, \dots) \parallel p([\Theta] \mid Y_n, \mathbf{x}_n, \dots)) \rightarrow 0$ for $n \rightarrow \infty$, which *converges* as it is monotonically increasing and bounded by $D_{\text{KL}}(q([\Theta]) \parallel p([\Theta]))$ from above, and thanks of Bernstein von Mises' theorem we have:

$$\begin{aligned} D_{\text{KL}}(q([\Theta]) \parallel p([\Theta])) &= \\ &= \sup_{n \in \mathbb{N}} D_{\text{KL}}(q(Y_n, \dots \mid \mathbf{x}_n, \dots) \parallel p(Y_n, \dots \mid \mathbf{x}_n, \dots)). \end{aligned}$$

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 that they have support*, which is true for $[\Theta]$.

Function-Space Variational Inference & ELBO

FSVI's ELBO is just the regular ELBO but for $[\Theta]$ and approximations via chain rule of the DPI:

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

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

$$\begin{aligned} &= \mathbb{E}_{q(\Theta)} [-\log p(\mathcal{D} \mid \Theta)] \\ &\quad + \sup_n D_{\text{KL}}(q(Y_{n\dots} \mid \mathbf{x}_{n\dots}) \parallel p(Y_{n\dots} \mid \mathbf{x}_{n\dots})) \\ &\geq \mathbb{E}_{q(\Theta)} [-\log p(\mathcal{D} \mid \Theta)] \\ &\quad + D_{\text{KL}}(q(Y_{n\dots} \mid \mathbf{x}_{n\dots}) \parallel p(Y_{n\dots} \mid \mathbf{x}_{n\dots})) \quad \forall n. \end{aligned}$$

More Info



More References

- [1] Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- [2] Tim G. J. Rudner, Zonghao Chen, Yee Whye Teh, and Yarin Gal. Tractable function-space variational inference in bayesian neural networks. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.