

GP JOURNAL CLUB

On Sparse Variational Methods
and the KL Between Stochastic Processes

Wessel Bruinsma

Machine Learning Group, University of Cambridge

18 Dec 2020

- d. G. Matthews, A. G., Hensman, J., Turner, R. E., & Ghahramani, Z. (2016). On sparse variational methods and the Kullback-Leibler divergence between stochastic processes. In A. Singh & J. Zhu (Eds.), *Proceedings of the 22nd international conference on artificial intelligence and statistics* (Vol. 54). Proceedings of Machine Learning Research, Proceedings of Machine Learning Research. eprint: <https://arxiv.org/abs/1504.07027>

- Titsias (2009) introduced trick to efficiently approximate posteriors of GPs using fewer *inducing points*.
- Very popular, but remained unclear whether Titsias' trick targets *posterior process*.
- d. G. Matthews et al. (2016) established a general condition under which this is the case.
- Puts massive body of work on solid theoretical grounding.

- I. Titsias' Trick and the Posterior Process
- II. Analysis with the ∞ -Dimensional Lebesgue Measure
- III. Formalising the Argument
- IV. Wrap-Up

Titsias' Trick and the Posterior Process

- Consider prior $f \sim \mathcal{GP}(0, k)$ and observations $D = (\mathbf{x}, \mathbf{y})$.
- **Goal:** Efficiently approximate $p(f | D)$.

- ① Augment model with inducing points \mathbf{u} :

$$p(f, \mathbf{u}) = p(f)p(\mathbf{u} | f).$$

↙ augmentation

↑ original prior

- ② Define approximate posterior:

$$q(f, \mathbf{u}) = p(f | \mathbf{u})q(\mathbf{u}).$$

↙ from augmented prior

↑ variational variable

- ③ Optimise approximation:

$$q^*(\mathbf{u}) = \arg \max_{q(\mathbf{u}) \in \mathcal{P}_{\mathbf{G}}} \text{ELBO}(q(\mathbf{u}))$$

✓ Profit!

implicitly depends on $q(\mathbf{u})$: $q(f) = \int p(f | \mathbf{u})q(\mathbf{u}) d\mathbf{u}$

① Is it true that

$$q^*(\mathbf{u}) = \arg \min_{q(\mathbf{u}) \in \mathcal{P}_{\mathbf{G}}} \text{KL}(q(f), p(f | D))?$$

- This is a KL between *processes*!
- ② Can we formally make sense of " $p(f)$ "?

Analysis with the ∞ -Dimensional Lebesgue Measure

- **Thm:** The ∞ -dimensional Lebesgue measure does *not* exist.
- For now, think

$$f = \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix}$$

for some large $n \gg 1$.

- We will later assign meaning to “ $p(f)$ ”.

- ① Is it true that

$$q^*(\mathbf{u}) = \arg \min_{q(\mathbf{u}) \in \mathcal{P}_{\mathbf{G}}} \text{KL}(q(f), p(f | D))?$$

- Strategy: Decompose

$$\text{KL}(q(f, \mathbf{u}), p(f, \mathbf{u} | D))$$

in two ways.

- Note that

$$\frac{p(f, \mathbf{u} | D)}{p(f, \mathbf{u})} = \frac{p(D | f)}{p(D)} \quad \text{and} \quad \frac{q(f, \mathbf{u})}{p(f, \mathbf{u})} = \frac{q(\mathbf{u})}{p(\mathbf{u})}.$$

- Then

$$\begin{aligned} & \text{KL}(q(f, \mathbf{u}), p(f, \mathbf{u} | D)) \\ &= \int \log \frac{q(f, \mathbf{u})/p(f, \mathbf{u})}{p(f, \mathbf{u} | D)/p(f, \mathbf{u})} q(f, \mathbf{u}) \, df \, d\mathbf{u} \\ &= \int \log \frac{q(\mathbf{u})/p(\mathbf{u})}{p(D | f)/p(D)} q(f, \mathbf{u}) \, df \, d\mathbf{u} \\ &= \log p(D) - \underbrace{\left(\mathbb{E}_q[\log p(D | f)] - \text{KL}(q(\mathbf{u}), p(\mathbf{u})) \right)}_{\text{ELBO}(q(\mathbf{u}))}. \end{aligned}$$

- Result so far:

$$\begin{aligned} \log p(D) - \text{ELBO}(q(\mathbf{u})) \\ = \text{KL}(q(f), p(f | D)) + \mathbb{E}_q[\text{KL}(q(\mathbf{u} | f), p(\mathbf{u} | f))]. \end{aligned}$$

- Therefore, if $q(\mathbf{u} | f) = p(\mathbf{u} | f)$, then

$$\arg \max_{q(\mathbf{u}) \in \mathcal{P}_{\mathbf{G}}} \text{ELBO}(q(\mathbf{u})) = \arg \min_{q(\mathbf{u}) \in \mathcal{P}_{\mathbf{G}}} \text{KL}(q(f), p(f | D)) !$$

- Thm (d. G. Matthews et al., 2016):

$$p(\mathbf{u} | f) = \delta(\mathbf{u} - T(f)) \implies q(\mathbf{u} | f) = \delta(\mathbf{u} - T(f)).$$

\uparrow
 deterministic transform of f

- Example 1 (inducing points):

$$T(f) = (f(z_1), \dots, f(z_m)).$$

- Example 2 (interdomain transform):

$$T(f) = z \mapsto \int h(z, x) f(x) dx.$$

- Proof (ish):

$$\begin{aligned} q(\mathbf{u} | f) &= \frac{q(\mathbf{u})}{q(f)} q(f | \mathbf{u}) = \frac{q(\mathbf{u})}{q(f)} p(f | \mathbf{u}) \\ &= \frac{q(\mathbf{u})}{q(f)} \frac{p(f)}{p(\mathbf{u})} p(\mathbf{u} | f) = \frac{q(\mathbf{u})}{p(\mathbf{u})} \frac{p(f)}{q(f)} \delta(\mathbf{u} - T(f)). \end{aligned}$$

Formalising the Argument

- Consider

$X =$ all functions $f: \mathcal{X} \rightarrow \mathbb{R}$,

$Z =$ all functions $g: \mathcal{Z} \rightarrow \mathbb{R}$,

and assume that these spaces are Polish. ← complete, separable, metrisable

- We will consider probability measures on

$$(X \times Z, \mathcal{B}(X) \otimes \mathcal{B}(Z)).$$



Borel σ -algebra

- For such \mathbb{P} ,

marginals: \mathbb{P}_X and \mathbb{P}_Z ,

conditionals: $\mathbb{P}_{X|Z}$ and $\mathbb{P}_{Z|X}$.

- Generally denote $A \in \mathcal{B}(X)$ and $B \in \mathcal{B}(Z)$.

$$p(\mathbf{x}) = 0 \implies q(\mathbf{x}) = 0$$

- Def. (\mathbb{R}^n): For $q(\mathbf{x})$ and $p(\mathbf{x})$ such that $q(\mathbf{x}) \ll p(\mathbf{x})$, define

$$\text{KL}(q(\mathbf{x}), p(\mathbf{x})) = \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x}.$$

- Def. (general): For \mathbb{Q} and \mathbb{P} such that $\mathbb{Q} \ll \mathbb{P}$, define

$$\text{KL}(\mathbb{Q}, \mathbb{P}) = \int \frac{d\mathbb{Q}}{d\mu} \log \frac{d\mathbb{Q}/d\mu}{d\mathbb{P}/d\mu} d\mu$$

↑ plays role of dx

where μ is any measure such that $\mathbb{P} \ll \mu$.

- Augmented model:

$$\mathbb{P}(A \times B) = \int_A \mathbb{P}_{Z|X}(B | f) d\mathbb{P}_X(f).$$

augmentation
 ↓
 original prior
 ↑

- Approximate posterior:

$$\mathbb{Q}(A \times B) = \int_B \mathbb{P}_{X|Z}(A | g) d\mathbb{Q}_Z(g).$$

from augmented prior
 ↓
 variational variable
 ↑

- Exact posterior:

$$\frac{d\mathbb{P}(\cdot | D)}{d\mathbb{P}}(f, g) = \frac{p(D | f)}{\mathbb{E}_{\mathbb{P}}[p(D | f)]}.$$

WHAT WE WROTE

Observations:

$$\frac{p(f, \mathbf{u} | D)}{p(f, \mathbf{u})} = \frac{p(D | f)}{p(D)},$$

$$\frac{q(f, \mathbf{u})}{p(f, \mathbf{u})} = \frac{q(\mathbf{u})}{p(\mathbf{u})}.$$

WHAT WE MEANT

$$\frac{d\mathbb{P}(\cdot | D)}{d\mathbb{P}}(f, g) = \frac{p(D | f)}{\mathbb{E}_{\mathbb{P}}[p(D | f)]},$$

$$\frac{d\mathbb{Q}}{d\mathbb{P}}(f, g) = \frac{d\mathbb{Q}_Z}{d\mathbb{P}_Z}(g).$$

ELBO:

$$\text{KL}(q(f, \mathbf{u}), p(f, \mathbf{u} | D)) \int q(f, \mathbf{u}) df d\mathbf{u}$$

$$= \int \log \frac{q(f, \mathbf{u})/p(f, \mathbf{u})}{p(f, \mathbf{u} | D)/p(f, \mathbf{u})} dq$$

$$= \int \log \frac{q(\mathbf{u})/p(\mathbf{u})}{p(D | f)/p(D)} dq$$

$$= \log p(D) - \text{ELBO}(q(\mathbf{u})).$$

$$\text{KL}(\mathbb{Q}, \mathbb{P}(\cdot | D))$$

$$= \int \log \frac{d\mathbb{Q}/d\mathbb{P}}{d\mathbb{P}(\cdot | D)/d\mathbb{P}} d\mathbb{Q}$$

$$= \int \log \frac{d\mathbb{Q}_Z/d\mathbb{P}_Z}{p(D | f)/p(D)} d\mathbb{Q}$$

$$= \log p(D) - \text{ELBO}(\mathbb{Q}_Z).$$

- Nothing to do! Chain rule works for general KL.

- Result so far:

$$\begin{aligned} \log \mathbb{E}_{\mathbb{P}}[p(D | f)] - \text{ELBO}(\mathbb{Q}_Z) \\ = \text{KL}(\mathbb{Q}_X, \mathbb{P}_X(\cdot | D)) + \mathbb{E}_{\mathbb{Q}}[\text{KL}(\mathbb{Q}_{Z|X}, \mathbb{P}_{Z|X})]. \end{aligned}$$

- Therefore, if $\mathbb{Q}_{Z|X} = \mathbb{P}_{Z|X}$, then

$$\arg \max_{\mathbb{Q}_Z \in \mathcal{P}_{\mathcal{G}}} \text{ELBO}(\mathbb{Q}_Z) = \arg \min_{\mathbb{Q}_Z \in \mathcal{P}_{\mathcal{G}}} \text{KL}(\mathbb{Q}_X, \mathbb{P}_X(\cdot | D)) !$$

- Thm (d. G. Matthews et al., 2016):

$$\mathbb{P}_{Z|X}(\cdot | f) = \delta_{T(f)} \implies \mathbb{Q}_{Z|X}(\cdot | f) = \delta_{T(f)}.$$

- Proof: More involved. See paper.

Wrap-Up

- Titsias' trick targets the **posterior process** if *inducing function* is a deterministic transform of f (d. G. Matthews et al., 2016).
- ⇒ Formally justifies inducing points and interdomain transforms.
- Although ∞ -dimensional Lebesgue measure does not exist, manipulations can directly translate to formal manipulations with $p(f) df = d\mathbb{P}$.

Appendix

References

- d. G. Matthews, A. G., Hensman, J., Turner, R. E., & Ghahramani, Z. (2016). On sparse variational methods and the Kullback-Leibler divergence between stochastic processes. In A. Singh & J. Zhu (Eds.), *Proceedings of the 22nd international conference on artificial intelligence and statistics* (Vol. 54). Proceedings of Machine Learning Research, Proceedings of Machine Learning Research. eprint: <https://arxiv.org/abs/1504.07027>
- Titsias, M. K. (2009). Variational learning of inducing variables in sparse Gaussian processes. In D. van Dyk & M. Welling (Eds.), *Proceedings of the 12th international conference on artificial intelligence and statistics* (Vol. 12, pp. 567–574). Proceedings of Machine Learning Research. Proceedings of Machine Learning Research. Retrieved from <http://proceedings.mlr.press/v5/titsias09a/titsias09a.pdf>