

Why Doubly Stochastic Variational Inference ??

Faezeh Yazdi

October 18, 2019

What Problem?

- ▶ A set of N observations $\mathbf{y} = (y_1, \dots, y_N)^T$ where $y_i \in \mathbb{R}$
- ▶ At N design locations $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ where $\mathbf{x}_i \in \mathbb{R}^d$
- ▶ $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is responsible for generating \mathbf{y} given \mathbf{X}

How to approximate $\hat{f}(\mathbf{x}^*)$ at new location \mathbf{x}^* ??

- ▶ Fitting a GP/ a DGP
- ▶ Computational complexity is $O(N^3)$

A problem for GPs and even more in DGPs (MORE Data needed)

How complex DGP Models?

- ▶ Architecture of a DGP model with L layers corresponds to this graphical model :

$$\mathbf{X} \xrightarrow{\mathbf{f}^1} (\mathbf{F}^1 + \epsilon_1) \xrightarrow{\mathbf{f}^2} \dots \xrightarrow{\mathbf{f}^{L-1}} (\mathbf{F}^{L-1} + \epsilon_{L-1}) \xrightarrow{\mathbf{f}^L} (\mathbf{F}^L + \epsilon_L) = \mathbf{y}$$

- ▶ A nested structure of GPs considering the relationship between the inputs and the final output as a functional composition of GPs :

$$\mathbf{y} = \mathbf{f}^L(\mathbf{f}^{L-1}(\dots \mathbf{f}^2((\mathbf{f}^1(\mathbf{X}) + \epsilon_1)) + \epsilon_2\dots) + \epsilon_{L-1}) + \epsilon_L$$

- ▶ Latent nodes $(\mathbf{F}^l + \epsilon_l) \in \mathbb{R}^{N \times d}$ for $l = 1, 2, \dots, L - 1$
- ▶ $\mathbf{f}^l(\cdot)$ is a multi-output GP for $l = 1, 2, \dots, L - 1$
- ▶ $\epsilon_l \sim N(0, \sigma_{\epsilon_l}^2 \mathbf{I})$

Computationally Difficult for Large N
Collection of lots of Parameters (Latent Variables) to Estimate
Large Parameter Estimation Space

How can do Inference on DGPs?

- ▶ MCMC to get $\pi(\hat{f}(\mathbf{x}^*)|\mathbf{X}, \mathbf{y})$
- ▶ VB approaches to get $\overbrace{\pi(\hat{f}(\mathbf{x}^*)|\mathbf{X}, \mathbf{y})}$

BI into an optimization problem over variational parameters

- ▶ Kulback-Liebler (KL) divergence
- ▶ Inducing Points
- ▶ Minibatches (subsampling)
- ▶ Doubly Stochastic Variational Inference (DSVI)
 - ▶ Doing Inference quickly
 - ▶ Using TensorFlow

Creating a computational graph, which takes multidimensional matrices called “Tensors” and does computations on them
 - ▶ Using GPflow

A package for building GP models in python, using TensorFlow.

Other Features of DSVI

- ▶ A variational algorithm for inference in DGP models that does not force independence or Gaussianity between the layers
- ▶ The variational posterior has the same structure as the full model (it maintains the correlations between layers)
- ▶ Using sparse VI to simplify the correlations within layers
- ▶ This method of inference does not require us to parameterize noisy corruptions separately from the outputs of each GP.
- ▶ They absorb the noise into the kernel

$$k_{\text{noisy}}(x_i, x_j) = k(x_i, x_j) + \sigma^2 \delta_{ij}$$

where δ_{ij} is the Kronecker delta, and σ^2 is the noise variance between layers

- ▶ Complexity of each generative GP mapping is reduced from $O(N^3)$ to $O(NM^2)$

$$X \xrightarrow{f^1} F'_+ + \varepsilon_1 \xrightarrow{f^2} \dots \xrightarrow{f^{L-1}} F^{L-1}_+ + \varepsilon_{L-1} \xrightarrow{f^L} F^L_+ + \varepsilon_L = Y$$

$\xrightarrow{z^t \rightarrow U^t}$

$$f^1 \sim Gp(0, k_{xx})$$

$$f^2 \sim Gp(0, k_{F'_+ + \varepsilon_1, F'_+ + \varepsilon_2})$$

⋮

$$f^L \sim Gp(0, k_{F^{L-1}_+ + \varepsilon_{L-1}, F^L_+ + \varepsilon_L})$$

Full model (In Lawrence paper)

$$P(Y, \{F^t, U^t\}_{t=1}^L, \{F'_+ + \varepsilon_t\}_{t=1}^L) = P(F'_+ + \varepsilon_t | F^t) P(F^t | U^t, X) P(U^t | z^t)$$

first layer

$$\times \prod_{\ell=2}^L P(F^{\ell} | F^{\ell-1}) P(F^{\ell} | U^{\ell}, F^{\ell-1}_+ + \varepsilon_t) P(U^{\ell} | z^{\ell-1})$$

Variational posterior (In Lawrence paper)

$$Q = P(F^t | U^t, X) q(U^t) \times \prod_{\ell=2}^L P(F^{\ell} | U^{\ell}, F^{\ell-1}_+ + \varepsilon_t) q(U^{\ell}) \boxed{q(F'_+ + \varepsilon_t)}$$

\downarrow
 $\prod_{\ell=1}^L N(\mu, S)$

$$X \xrightarrow{f^1} F^1_+ + \varepsilon_1 \xrightarrow{f^2} \dots \xrightarrow{f^{L-1}} F^{L-1}_+ + \varepsilon_{L-1} \xrightarrow{f^L} F^L_+ + \varepsilon_L = Y$$

$\xrightarrow{z^t \rightarrow U^t}$

$$f^1 \sim Gp(0, k_{xx}^{noisy})$$

$$f^2 \sim Gp(0, k_{F^1_+ + \varepsilon_1, F^1_+ + \varepsilon_2}^{noisy})$$

⋮

$$f^L \sim Gp(0, k_{F^{L-1}_+ + \varepsilon_{L-1}, F^L_+ + \varepsilon_L}^{noisy})$$

Full model (In Salimbeni Inference)

$$P(Y, \{F^t, U^t\}_{t=1}^L) = P(Y | F^L) \prod_{t=1}^L P(F^t | U^t, F^{\ell-1}_+ + \varepsilon_t) P(U^t | z^{\ell-1})$$

Variational posterior (In salimbeni Inference)

$$Q(\{F^t, U^t\}_{t=1}^L) = \prod_{\ell=1}^L P(F^{\ell} | U^{\ell}, F^{\ell-1}_+ + \varepsilon_t) q(U^{\ell})$$

Why Doubly?

The evidence lower bound of the DGP in this inference method is

$$\mathcal{L}_{DGP} = \sum_{i=1}^N \mathbb{E}_{q(\mathbf{f}_i^L)} [\log p(\mathbf{y}_n | \mathbf{f}_n^L)] - \sum_{l=1}^L \text{KL}[q(\mathbf{U}^l) || p(\mathbf{U}^l; \mathbf{Z}^{l-1})]$$

Approximate ELBO using two sources of stochasticity

- ▶ Approximate the expectation with a Monte Carlo sample from the variational posterior
- ▶ Samples from the marginals at the top layer are obtained from marginals of all the other layers
- ▶ The bound factorizes over the data we achieve scalability through sub-sampling the data.
- ▶ This bound has complexity $O(NM^2Ld)$ to evaluate

Easy to Set up the Problem !

Let's start with a simple example (Fitting a single layer DGP) :

- ▶ Randomly sample 500 points $x \in [0, 1]^2$
- ▶ Generate Y from $f(x) = \text{Sin}(2\pi x_1) + \text{Cos}(2\pi x_2)$
- ▶ Randomly sample 100 inducing points $z \in [0, 1]^2$
- ▶ kernel=RBF(d, lengthscales=0.2, variance=1.)
- ▶ mdgp = DGP(X, Y, Z, Kernel, Gaussianlik(), minibatchsize=None)
- ▶ AdamOptimizer(0.001).minimize(mdgp, maxiter=2000)
- ▶ samples, m,v = mgp.predict.y(Xp)

- ▶ Adam ("Adaptive moments") optimizer is an adaptive learning rate optimization algorithm that's been designed specifically for training deep neural networks.
- ▶ adaptive learning rate optimization algorithm which means, it computes individual learning rate for different parameters
- ▶ Adam uses estimations of first and second moments of gradients to adapt the learning rate for each parameter

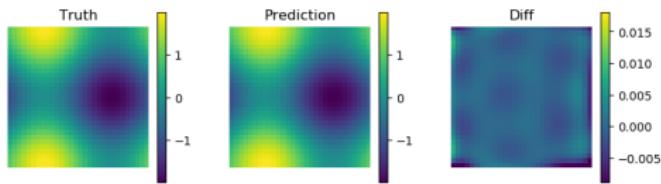
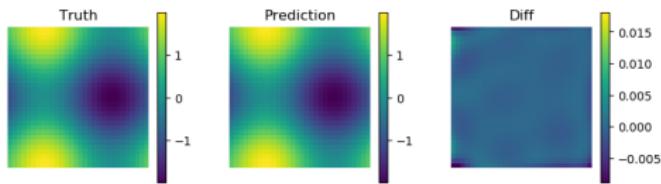
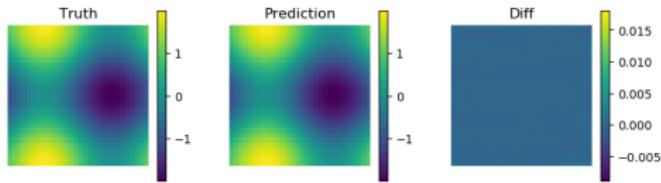


Figure: From top to bottom: Prediction using GP, One single layer DGP and SVGP

Table: Results

Model	RMSPE	Opt Time
GP	0.00015139	1.79 s
DGP($L=1$)	0.00111246	21.74 s
Sparse GP	0.00208732	17.00 s