

## [ Paper review 36 ]

---

# Gaussian Processes for Big Data

---

( Hensman, et al., 2013 )

## [ Contents ]

---

1. Abstract
2. Introduction
3. Sparse GPs Revisited
4. SVI for GPs
  1. Global variables
  2. Natural gradients
  3. Latent variables
  4. Non-Gaussian likelihood
5. Discussion

## 1. Abstract

---

introduce **SVI** (Stochastic Variational Inference) for **GP** (Gaussian Process) models

- enable GP models to be **scalable**
- show that GPs can be variationally decomposed, to be dependent on a set of **globally relevant inducing variables**

## 2. Introduction

---

GP, used for regression, classification, unsupervised learning..

drawback : complexity of  $O(n^3)$

To deal with this problem, various approximate techniques have been proposed

- 1) partition data set into separate groups
- 2) low rank approximation to the covariance matrix ( complexity of  $O(nm^2)$  )
- 3) (by this paper)

**"recent advances in VI can be combined with the idea of INDUCING VARIABLES to develop a practical algorithm for fitting GPs using SVI"**

# 3. Sparse GPs Revisited

## inducing variables of Titsias (2009)

notation

- $\mathbf{y}$  : data vector  
( consists of  $y_i$ , which are noisy observation of the function  $f(\mathbf{x}_i)$  )  
( independent Gaussian, with precision  $\beta$  )
- $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$  : all the datapoints
- **INDUCING VARIABLES** :  $\mathbf{u}$  values of the function  $f$  at the points  $\mathbf{Z} = \{\mathbf{z}_i\}_{i=1}^m$

$$p(\mathbf{y} | \mathbf{f}) = \mathcal{N}(\mathbf{y} | \mathbf{f}, \beta^{-1}\mathbf{I})$$

$$p(\mathbf{f} | \mathbf{u}) = \mathcal{N}(\mathbf{f} | \mathbf{K}_{nm}\mathbf{K}_{mm}^{-1}\mathbf{u}, \tilde{\mathbf{K}}).$$

$$p(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \mathbf{0}, \mathbf{K}_{mm})$$

- $\mathbf{K}_{mm}$  : covariance function evaluated between **all the inducing points**
- $\mathbf{K}_{nm}$  : covariance function between **all inducing points** and **training points**
- $\tilde{\mathbf{K}} = \mathbf{K}_{nn} - \mathbf{K}_{nm}\mathbf{K}_{mm}^{-1}\mathbf{K}_{mn}$

Apply Jensen's inequality on the conditional probability  $p(\mathbf{y} | \mathbf{u})$

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{u}) &= \log \langle p(\mathbf{y} | \mathbf{f}) \rangle_{p(\mathbf{f}|\mathbf{u})} \\ &\geq \langle \log p(\mathbf{y} | \mathbf{f}) \rangle_{p(\mathbf{f}|\mathbf{u})} \triangleq \mathcal{L}_1 \end{aligned}$$

- $\langle \cdot \rangle_{p(x)}$  : expectation under  $p(x)$ .
- $\log \langle p(\mathbf{y} | \mathbf{f}) \rangle_{p(\mathbf{f}|\mathbf{u})}$  : computed by  $\mathcal{O}(n^3)$
- $\langle \log p(\mathbf{y} | \mathbf{f}) \rangle_{p(\mathbf{f}|\mathbf{u})} \triangleq \mathcal{L}_1$  : computed by  $\mathcal{O}(m^3)$ .

Interpretation : belows are equivalent

- $\mathbf{u} = \mathbf{f}$
- $\mathbf{K}_{mm} = \mathbf{K}_{nn} = \mathbf{K}_{mm}$
- $m = n$  inducing variables, and they are placed at training data locations
- no computational / storage advantage

when  $p(\mathbf{y} | \mathbf{f})$  factorizes across the data,  $p(\mathbf{y} | \mathbf{f}) = \prod_{i=1}^n p(y_i | f_i)$ .

→ then  $\exp(\mathcal{L}_1) = \prod_{i=1}^n \mathcal{N}(y_i | \mu_i, \beta^{-1}) \exp(-\frac{1}{2}\beta\tilde{k}_{i,i})$ .

- $\mu = \mathbf{K}_{nm}\mathbf{K}_{mm}^{-1}\mathbf{u}$ .
- $\tilde{k}_{i,i}$  :  $i$  th diagonal element of  $\tilde{\mathbf{K}}$ .

Bound of Titsias (2009)

- by marginalizing the inducing variables  $\mathbf{u}$

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{X}) &= \log \int p(\mathbf{y} | \mathbf{u}) p(\mathbf{u}) d\mathbf{u} \\ &\geq \log \int \exp\{\mathcal{L}_1\} p(\mathbf{u}) d\mathbf{u} \triangleq \mathcal{L}_2 \end{aligned}$$

- $\mathcal{L}_2 = \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} + \beta^{-1} \mathbf{I}) - \frac{1}{2} \beta \text{tr}(\tilde{\mathbf{K}})$ .
  - $\mathbf{\Lambda} = \beta \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} + \mathbf{K}_{mm}^{-1}$ .
  - $\hat{\mathbf{u}} = \beta \mathbf{\Lambda}^{-1} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} \mathbf{y}$ .

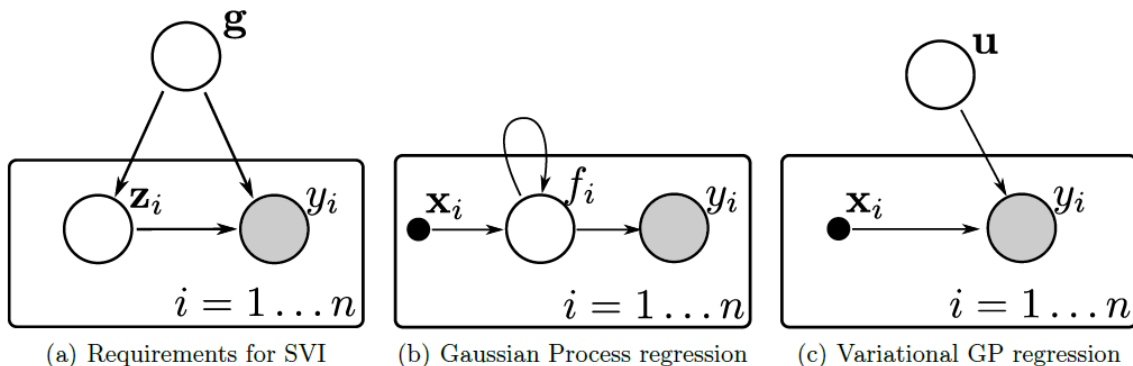
## 4. SVI for GPs

novelties of Titsias bound : rather than explicitly representing variational distribution for  $q(\mathbf{u})$ , these are collapsed!

but for SVI to work on GP, we need to "maintain an explicit representation of these inducing variables"

### SVI (Stochastic Variational Inference)

- works on large dataset
- but can be only applied to (probabilistic) models
  - which have **global variables**
  - which factorizes in the observations and latent variables
- by introducing  $\mathbf{u}$ , we satisfies the condition!



- But in above we have found lower bound as below  
 $(\log \int \exp\{\mathcal{L}_1\} p(\mathbf{u}) d\mathbf{u} \triangleq \mathcal{L}_2)$

### 4-1. Global Variables

New lower bound : ( $\mathcal{L}_2 \geq \mathcal{L}_3$  ,)

$$\log p(\mathbf{y} | \mathbf{X}) \geq \langle \mathcal{L}_1 + \log p(\mathbf{u}) - \log q(\mathbf{u}) \rangle_{q(\mathbf{u})} \triangleq \mathcal{L}_3.$$

Now, parameterize our variational distribution as  $q(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S})$ .

$$\mathcal{L}_3 = \sum_{i=1}^n \left\{ \log \mathcal{N}(y_i | \mathbf{k}_i^\top \mathbf{K}_{mm}^{-1} \mathbf{m}, \beta^{-1}) - \frac{1}{2} \beta \tilde{k}_{i,i} - \frac{1}{2} \text{tr}(\mathbf{S} \boldsymbol{\Lambda}_i) \right\} - \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})).$$

- $\mathbf{k}_i$  : vector of the  $i^{\text{th}}$  column of  $\mathbf{K}_{mn}$
- $\boldsymbol{\Lambda}_i = \beta \mathbf{K}_{mm}^{-1} \mathbf{k}_i \mathbf{k}_i^\top \mathbf{K}_{mm}^{-1}$ .
- $\mathcal{L}_3$  can be written as sum of  $n$  terms!

Gradients of lower bound  $\mathcal{L}_3$

- $\frac{\partial \mathcal{L}_3}{\partial \mathbf{m}} = \beta \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} \mathbf{y} - \boldsymbol{\Lambda} \mathbf{m}$
- $\frac{\partial \mathcal{L}_3}{\partial \mathbf{S}} = \frac{1}{2} \mathbf{S}^{-1} - \frac{1}{2} \boldsymbol{\Lambda}$
- setting the above to zero...
  - $\mathbf{S} = \boldsymbol{\Lambda}^{-1}, \mathbf{m} = \hat{\mathbf{u}}$ .
  - This is when  $\mathcal{L}_2 = \mathcal{L}_3$

## 4-2. Natural Gradients

SVI works by taking steps in the direction of **approximate natural gradient** ( $= \tilde{\mathbf{g}}(\boldsymbol{\theta}) = G(\boldsymbol{\theta})^{-1} \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}}$ )

canonical and expectation parameters

- $\boldsymbol{\theta}_1 = \mathbf{S}^{-1} \mathbf{m}$   
 $\boldsymbol{\theta}_2 = -\frac{1}{2} \mathbf{S}^{-1}$
- $\boldsymbol{\eta}_1 = \mathbf{m}$   
 $\boldsymbol{\eta}_2 = \mathbf{m} \mathbf{m}^\top + \mathbf{S}$

simplification of natural gradient :  $\tilde{\mathbf{g}}(\boldsymbol{\theta}) = G(\boldsymbol{\theta})^{-1} \frac{\partial \mathcal{L}_3}{\partial \boldsymbol{\theta}} = \frac{\partial \mathcal{L}_3}{\partial \boldsymbol{\eta}}$ .

Therefore , by  $\boldsymbol{\theta}_{(t+1)} = \boldsymbol{\theta}_{(t)} + \ell \frac{d\mathcal{L}_3}{d\boldsymbol{\eta}}$ ,

$$\begin{aligned} \boldsymbol{\theta}_{2(t+1)} &= -\frac{1}{2} \mathbf{S}_{(t+1)}^{-1} \\ &= -\frac{1}{2} \mathbf{S}_{(t)}^{-1} + \ell \left( -\frac{1}{2} \boldsymbol{\Lambda} + \frac{1}{2} \mathbf{S}_{(t)}^{-1} \right) \\ \boldsymbol{\theta}_{1(t+1)} &= \mathbf{S}_{(t+1)}^{-1} \mathbf{m}_{(t+1)} \\ &= \mathbf{S}_{(t)}^{-1} \mathbf{m}_{(t)} + \ell \left( \beta \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} \mathbf{y} - \mathbf{S}_{(t)}^{-1} \mathbf{m}_{(t)} \right) \end{aligned}$$

## 4-3. Latent Variables

- enable **online learning** for GPR using SVI
- to perform SVI with latent variable, need **factorization** ( like figure 1-(a) )

$$\begin{aligned} \log p(\mathbf{y}) &= \log \int p(\mathbf{y} | \mathbf{X}) p(\mathbf{X}) d\mathbf{X} \\ &\geq \int q(\mathbf{X}) \{ \mathcal{L}_3 + \log p(\mathbf{X}) - \log q(\mathbf{X}) \} d\mathbf{X} \end{aligned}$$

- $q(\mathbf{X}) = \prod_{i=1}^n q(x_i)$ .

To perform SVI in this model, alternate between...

- (1) selecting a mini-batch of data
- (2) optimizing relevant variables of  $q(\mathbf{X})$  ( with  $q(\mathbf{u})$  fixed ) and updating  $q(\mathbf{u})$  using the approximate natural gradient

## 4-4. Non-Gaussian likelihood

---

Advantage of using

$$\mathcal{L}_3 = \sum_{i=1}^n \left\{ \log \mathcal{N}(y_i | \mathbf{k}_i^\top \mathbf{K}_{mm}^{-1} \mathbf{m}, \beta^{-1}) - \frac{1}{2} \beta \tilde{k}_{i,i} - \frac{1}{2} \text{tr}(\mathbf{S} \mathbf{\Lambda}_i) \right\} - \text{KL}(q(\mathbf{u}) || p(\mathbf{u}))$$

→ enable inference with non-Gaussian likelihoods

ex) binary, probit likelihood ....

## 5. Discussion

---

method for inference in GP using SVI ( enable scalability )

discuss the bound on  $p(\mathbf{y} | \mathbf{u})$  in detail

( becomes tight when  $\mathbf{Z} = \mathbf{X}$  )

complexity becomes  $O(m^3)$