

Variational Learning for Sparse GP

Background

Given observation vectors \mathbf{x} and \mathbf{y} , we want to find the latent function value \mathbf{f} . Our goal is the posterior distribution $p(\mathbf{f}|\mathbf{y}) \propto (\mathbf{y}|\mathbf{f})p(\mathbf{f})$ through maximizing the log marginal likelihood $\log p(\mathbf{y}) = \log[N(\mathbf{y}|\mathbf{0}, \sigma^2 I + K_{nn})]$, which uses $O(n^3)$ time.

Principles behind Variational Inference in Sparse GP

Besides latent variable \mathbf{f} , let us introduce two new latent variables \mathbf{f}_m and \mathbf{z} . \mathbf{f}_m , named inducing variables, are function values evaluated at pseudo-inputs X_m , which live in the same space as training inputs \mathbf{x} but are independent from them. For latent variable \mathbf{z} , it is any set of true function values.

Remember that we want posterior $p(\mathbf{f}|\mathbf{y})$, we can describe it by predictive distribution $p(\mathbf{z}|\mathbf{y})$

$$\begin{aligned} p(\mathbf{z}|\mathbf{y}) &= \iint p(\mathbf{z}, \mathbf{f}, \mathbf{f}_m|\mathbf{y}) d\mathbf{f} d\mathbf{f}_m \\ &= \iint p(\mathbf{z}|\mathbf{f}, \mathbf{f}_m, \mathbf{y}) p(\mathbf{f}, \mathbf{f}_m|\mathbf{y}) d\mathbf{f} d\mathbf{f}_m \\ &= \iint p(\mathbf{z}|\mathbf{f}, \mathbf{f}_m) p(\mathbf{f}, \mathbf{f}_m|\mathbf{y}) d\mathbf{f} d\mathbf{f}_m \end{aligned}$$

Now we want to use $q(\mathbf{z})$ to approximate $p(\mathbf{z}|\mathbf{y})$, i.e., $q(\mathbf{z}) \approx p(\mathbf{z}|\mathbf{y})$:

If we can approximate the posterior distribution $p(\mathbf{f}, \mathbf{f}_m|\mathbf{y})$ by variational distribution $q(\mathbf{f}, \mathbf{f}_m)$ and the inducing variables \mathbf{f}_m are good enough to be sufficient statistics for \mathbf{f} , i.e., $q(\mathbf{f}, \mathbf{f}_m) \approx p(\mathbf{f}, \mathbf{f}_m|\mathbf{y})$, $p(\mathbf{z}|\mathbf{f}, \mathbf{f}_m) = p(\mathbf{z}|\mathbf{f}_m)$

Thus

$$\begin{aligned} p(\mathbf{z}|\mathbf{y}) &\approx q(\mathbf{z}) = \iint p(\mathbf{z}|\mathbf{f}_m) q(\mathbf{f}, \mathbf{f}_m) d\mathbf{f} d\mathbf{f}_m \\ &= \iint p(\mathbf{z}|\mathbf{f}_m) p(\mathbf{f}|\mathbf{f}_m) \phi(\mathbf{f}_m) d\mathbf{f} d\mathbf{f}_m \\ &= \int p(\mathbf{z}|\mathbf{f}_m) \left\{ \int p(\mathbf{f}|\mathbf{f}_m) d\mathbf{f} \right\} \phi(\mathbf{f}_m) d\mathbf{f}_m \\ &= \int p(\mathbf{z}|\mathbf{f}_m) \phi(\mathbf{f}_m) d\mathbf{f}_m \\ &= \int q(\mathbf{z}, \mathbf{f}_m) d\mathbf{f}_m \end{aligned}$$

where $\phi(\mathbf{f}_m)$ is free variation Gaussian distribution for \mathbf{f}_m with mean μ and covariance A . We can obtain the mean and covariance of this approximate posterior $q(\mathbf{z})$:

$$m(\mathbf{x}) = K_{xm} K_{mm}^{-1} \mu \quad \text{equation(1)}$$

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}') - K_{xm} K_{mm}^{-1} K_{mx'} + K_{xm} K_{mm}^{-1} A K_{mm}^{-1} K_{mx'}$$

The above defines the general form of sparse posterior GP, which is computed in $O(nm^2)$

Now our task is to find the optimal variational distribution $q(\mathbf{f}, \mathbf{f}_m) \approx p(\mathbf{f}, \mathbf{f}_m|\mathbf{y})$

In the theory of Variational Inference, the optimal variational distribution is always derived from a restricted distribution family, which has a factorization property. Here, the variational distribution $q(\mathbf{f}, \mathbf{f}_m)$ must satisfy the factorization: $q(\mathbf{f}, \mathbf{f}_m) = p(\mathbf{f}|\mathbf{f}_m) \phi(\mathbf{f}_m)$

I omit detailed derivations for the following equation:

$$\log p(\mathbf{y}) = \underbrace{\int q(\mathbf{f}, \mathbf{f}_m) \log \frac{p(\mathbf{f}, \mathbf{f}_m, \mathbf{y})}{q(\mathbf{f}, \mathbf{f}_m)} d\mathbf{f} d\mathbf{f}_m}_{ELBO} + KL(q(\mathbf{f}, \mathbf{f}_m) || p(\mathbf{f}, \mathbf{f}_m | \mathbf{y}))$$

To get the optimal $q(\mathbf{f}, \mathbf{f}_m)$, we will minimize the KL divergence, $KL(q(\mathbf{f}, \mathbf{f}_m) || p(\mathbf{f}, \mathbf{f}_m | \mathbf{y}))$, which is equivalent to maximize the evidence lower bound (ELBO), $\int q(\mathbf{f}, \mathbf{f}_m) \log \frac{p(\mathbf{f}, \mathbf{f}_m, \mathbf{y})}{q(\mathbf{f}, \mathbf{f}_m)} d\mathbf{f} d\mathbf{f}_m$. Now our objective function is ELBO, and our optimized variables are $\underbrace{X_m, \sigma^2, \theta}_{\text{variational parameters}}$ and $\underbrace{\phi(\mathbf{f}_m)}_{\text{free variational distribution}}$.

Therefore, I denote ELBO with $F_v(X_m, \sigma^2, \theta, \phi)$

We optimize $F_v(X_m, \sigma^2, \theta, \phi)$ with respect to $\phi(\mathbf{f}_m)$, and we get $\phi^*(\mathbf{f}_m) = N(\mathbf{f}_m | \mu, A)$, where $\mu = \sigma^{-2} K_{mm} \Sigma K_{mn} \mathbf{y}$, $A = K_{mm} \Sigma K_{mm}$, $\Sigma = (K_{mm} + \sigma^{-2} K_{mn} K_{nm})^{-1}$. Substitute μ and Σ into equation (1) and (2), then we can get the approximate posterior mean and variance.

Now with fixed optimal $\phi^*(\mathbf{f}_m)$, our objective function becomes

$$F_v = F_v(X_m, \sigma^2, \theta) = \log[N(\mathbf{y} | 0, \sigma^2 I + Q_{nn})] - \frac{1}{2\sigma^2} \underbrace{Tr(K_{nn} - Q_{nn})}_{\text{total variance of } p(\mathbf{f} | \mathbf{f}_m)}$$

where $Q_{nn} = K_{nm} K_{mm}^{-1} K_{mn}$, and $K_{nn} - Q_{nn} = Var[\mathbf{f} | \mathbf{f}_m]$

So far, our goal is to maximize the $F_v = F_v(X_m, \sigma^2, \theta)$ with respect to variational parameters (X_m, σ^2, θ) . In standard variational inference, we implement the optimization with gradient descent method as the following algorithm shows:

Algorithm 1 standard variational inference

Initialize inducing inputs X_m

repeat

$\theta = \argmax F_v$, fixed σ^2, X_m

$\sigma^2 = \frac{1}{n} \int \phi^*(\mathbf{f}_m) || \mathbf{y} - K_{nm} K_{mm}^{-1} \mathbf{f}_m ||^2 d\mathbf{f}_m + \frac{1}{n} Tr(K_{nn} - Q_{nn})$, fixed θ, X_m

$X_m = \argmax F_v$, fixed θ, σ^2

until Convergence

However, the standard gradient descent method will be difficult to implement; we can instead use greedy selection method. Greedy selection method results in a suboptimal solution, and it endures an easier algorithm. We have n training inputs, and basically, m inducing inputs will be selected among them. We start with an empty inducing set $m = \emptyset$ and a remaining set $n - m = \{1, \dots, n\}$. At each iteration, we add a training point $j \in J \subset n - m$, where J is a randomly chosen working set with the size of what we choose, into the inducing set that maximizes the selection criterion Δ_j . For me, I personally would like to choose the trace $Tr(K_{nn} - Q_{nn})$ to be the selection criterion, because it represents the total variance of the conditional prior $p(\mathbf{f} | \mathbf{f}_m)$. Thus smaller $Tr(K_{nn} - Q_{nn})$ means that the inducing variables \mathbf{f}_m are more likely to contain more information of \mathbf{f} .

Here is the algorithm using greedy selection method.

Algorithm 2 variational inference using greedy selection method

Start with with an empty inducing set $m = \emptyset$ and a remaining set $n - m = \{1, \dots, n\}$

repeat

- (1). Add a training point $j \in J \subset n - m$, where J is a randomly chosen working set with the size of what we choose, into the inducing set that minimize $Tr(K_{nn} - Q_{nn})$, fixed θ, σ^2
- (2). $\theta = \operatorname{argmax} F_v$, fixed σ^2, X_m
- (3). $\sigma^2 = \frac{1}{n} \int \phi^*(\mathbf{f}_m) \| \mathbf{y} - K_{nm} K_{mm}^{-1} \mathbf{f}_m \|^2 d\mathbf{f}_m + \frac{1}{n} Tr(K_{nn} - Q_{nn})$, fixed θ, X_m

until Convergence

This summary is my latest understanding after I read the Variational Inference Chapter in Pattern Recognition and Machine Learning, and the 2013 paper Stochastic Variational Inference. Thus the algorithm part of this summary is a little different from what I did in my previous code for the GSM kernel with variational inference. I will try to adapt my code to this summary, and to see whether I will obtain new improvements in my work.