

1(a) Write down the log posterior $\log p(f_x | y, X)$ for the regressors.
 $f_x = [f(x_1), f(x_2), \dots, f(x_n)]$

① prior over latent $f : p(f) = GP(f; m, k)$

② $y_i \in \{-1, 1\}$ input output pairs := $\{(x_i, y_i)\}_{i=1}^N$

$$\textcircled{3} \quad P(y|f_x) = \prod_{i=1}^N \frac{1}{1 + \exp(-y_i f(x_i))}$$

$$\sigma(y_i | f(x_i)) = \frac{1}{1 + \exp(-y_i f(x_i))}$$

$$\begin{aligned} \log p(f_x | y) &= \log p(y | f_x) + \log p(f_x) - \log(p(y)) \\ &= \sum_{i=1}^n \log \sigma(y_i | f(x_i)) - \frac{1}{2} (f_x - m_x)^T K_{xx}^{-1} (f_x - m_x) + \underset{\text{constant term.}}{\text{const.}} \\ &\approx -\sum_{i=1}^n \log (1 + \exp(-y_i f(x_i))) - \frac{1}{2} (f_x - m_x)^T K_{xx}^{-1} (f_x - m_x) \end{aligned}$$

(b) Compute the gradient of this log posterior

$$\nabla_{f_x} \log (P(f_x | y, x))$$

Answer: ① prior over latent $f : p(f) = GP(f; m, k)$

② $y_i \in \{-1, 1\}$ input output pairs := $\{(x_i, y_i)\}_{i=1}^N$

$$③ P(y|f_x) = \prod_{i=1}^n \frac{1}{1 + \exp(-y_i f(x_i))}$$

$$\begin{aligned} \text{log } P(f_x | y) &= \log P(y | f_x) + \log P(f_x) - \log (P(y)) \\ &= \sum_{i=1}^n \log \sigma(y_i f(x_i)) - \frac{1}{2} (f_x - m_x)^T K_{xx}^{-1} (f_x - m_x) + \text{const.} \end{aligned}$$

$$\begin{aligned} \nabla \log (P(f_x | y)) &= \sum_{i=1}^n \nabla_{f_x} \log \sigma(y_i f(x_i)) - K_{xx}^{-1} (f_x - m_x) \\ &= \sum_{i=1}^n - \underbrace{\nabla_{f_x} \log (1 + \exp(-y_i f(x_i)))}_{-\frac{\partial}{\partial f_x} \log (1 + \exp(-y_i f(x_i)))} - K_{xx}^{-1} (f_x - m_x) \\ &= \delta_{ij} \left(\frac{y_i + 1}{2} - \sigma(f_{x_i}) \right) \end{aligned}$$

$$\text{hence as } y_i \in \{-1, 1\} \quad \sigma(f_x) = \frac{1}{1 + \exp(f_x)}$$

$$\frac{y_i + 1}{2} \in \{0, 1\}$$

$$\text{Considerar } -\frac{\partial}{\partial f_{x_i}} \log (1 + \exp(-y_i f_{x_i})) = -\frac{1}{1 + \exp(-y_i f_{x_i})} - y_i \exp(-y_i f_{x_i})$$

$$= \frac{y_i \exp(-y_i f_{x_i})}{1 + \exp(-y_i f_{x_i})}$$

$$\textcircled{1} \quad y_i = -1 \Rightarrow -\frac{\exp(f_{x_i})}{1 + \exp(f_{x_i})} = \frac{-1}{1 + \exp(-f_{x_i})} = -\sigma(f_{x_i})$$

$$\textcircled{2} \quad y_i = 1 \Rightarrow \frac{\exp(-f_{x_i})}{1 + \exp(-f_{x_i})} = 1 - \frac{1}{1 + \exp(f_{x_i})} = 1 - \sigma(f_{x_i})$$

$$\Rightarrow \frac{y_i + 1}{2} - \sigma(f_{x_i})$$

$$\nabla \log(p(f_x | y)) = \begin{bmatrix} \frac{y_1 + 1}{2} - \sigma(f(x_0)) \\ \frac{y_2 + 1}{2} - \sigma(f(x_1)) \\ \vdots \\ \frac{y_n + 1}{2} - \sigma(f(x_n)) \end{bmatrix} - K_x^{-1} (f_n - m_n)$$

1(c). Assume that you have found the mode \hat{f}_x of this posterior (i.e. the MAP-estimate) and used it to construct an appropriate Gaussian posterior $q(f_x) = N(\hat{f}_x; \hat{f}_x, \Sigma)$ on f_x , where Σ is some covariance matrix which does not feature in the following. Show that the MAP estimate for the latent function f which can be written compactly as

$$E_q[f(\circ)] = m_0 + \kappa_{\circ x} \nabla \log p(y | \hat{f}_x)$$

that is using the gradient of the log-likelihood.

Answer:

① \hat{f}_x is the mode of posterior.

② $q(f_x) = N(f_x; \hat{f}_x, \Sigma)$ is the approximate Gaussian posterior

③ Map Estimate of latent function f can be written as

$$E_q(f(\cdot)) = m(x) + K_{xx} \nabla_{\hat{f}_x} \log p(y | \hat{f}_x)$$

④ We know that at $f_x = \hat{f}_x$ the mode.

$$\nabla_{f_x} \log p(f_x | y, x) = 0$$

$$\nabla_{f_x} \log p(y | f_x) + \nabla_{f_x} \log p(f_x) = 0$$

$$\nabla_{f_x} \log p(y | f_x) + \nabla_{f_x} \left[\frac{1}{2} (f_x - m_x)^T K_{xx}^{-1} (f_x - m_x) \right] = 0$$

$$\nabla_{f_x} \log p(y | f_x) - K_{xx}^{-1} (f_x - m_x) = 0$$

$$\nabla_{f_x} \log (p(y | f_x)) = K_{xx}^{-1} (f_x - m_x)$$

thus $f_x = \hat{f}_x$ we have

$$\nabla_{f_x} \log (p(y | \hat{f})) = K_{xx}^{-1} (\hat{f} - m_x)$$

⑤ As we know \hat{f} the approximate Gaussian posterior at training points

$$q(f_x) = \mathcal{N}(f_x; \hat{f}, \Sigma)$$

⑥ The posterior predictions at f_x (Lecture-14, Slid-17)

$$q(f_x | y) = \int p(f_x | \vec{f}_x) q(\vec{f}_x) d\vec{f}_x$$

$$= \mathcal{N}(f_x; m_x + K_{xx}^{-1} (\hat{f} - m_x), \Sigma_{\text{posterior}})$$

Thus $E_q(f(\cdot)) = m(\cdot) + K_{xx}^{-1} K_{xx}^{-1} (\hat{f} - m_x)$
as expected value is precisely mean of gaussian.

From ④ we have at $f_x = \hat{f}$

$$K_{xx}^{-1} (\hat{f} - m_x) = \nabla_{f_x} \log(p(y|\hat{f}))$$

Thus $E_q(f(\cdot)) = m(\cdot) + K_{xx} \nabla_{f_x} \log(p(y|\hat{f}))$

1(d) By writing down the explicit gradient of $\log p(y|\hat{f}_x)$, make an argument that those training points x_i at which $|\hat{f}(x_i)| \gg 1$, those "far from the decision boundary" do almost not contribute to this estimate $E_y[f(\cdot)]$.

$$\text{Ans: } \nabla_{\hat{f}_x} \log(p(y|\hat{f}_x)) = \sum_{i=1}^n \nabla_{\hat{f}_x} \log(\sigma(y_i + f(x_i))) \\ = \sum_{i=1}^n \nabla_{\hat{f}_x} \left[-\log(1 + \exp(-y_i f(x_i))) \right]$$

$$\frac{\partial}{\partial f_i} -\log(1 + \exp(-y_i f(x_i))) = \frac{+y_i \exp(-y_i f(x_i))}{1 + \exp(-y_i f(x_i))}$$

$$y_i = -1 \Rightarrow \frac{(-1) \exp(f(x_i))}{1 + \exp(f(x_i))} = \frac{-1}{1 + \exp(-f(x_i))} = -\sigma(f(x_i))$$

$$y_i = 1 \Rightarrow \frac{(1) \exp(-f(x_i))}{1 + \exp(-f(x_i))} = 1 - \frac{1}{1 + \exp(-f(x_i))} = 1 - \sigma(f(x_i))$$

$$\begin{aligned} \text{using mapping } & y_i = -1 & ; & y_i = 0 & \} \\ & y_i = 1 & ; & y_i = 1 & \} \end{aligned}$$

$$\frac{\partial}{\partial f_i} -\log(1 + \exp(-y_i f(x_i))) = \frac{y_i + 1}{2} - \sigma(f(x_i))$$

$$\nabla_i \log(p(y|\hat{f}_x)) = \frac{y_i + 1}{2} - \sigma(f(x_i))$$

Now consider $|\hat{f}_x| \gg 1$

$$\sigma(|\hat{f}_x|) = \frac{1}{1 + \exp(-|\hat{f}_x|)} = 1 \quad \text{as } \lim_{a \rightarrow \infty} \exp(-a) = 0$$

Thus when $y_i = 1$ and $|\hat{f}(x_i)| \gg 1$ and $\hat{f}(x_i) > 0$

$$\nabla_i \log(p(y|\hat{f}_x)) = \frac{1+1}{2} - 1 = 0$$

Thus x_i will not contribute to $E_g[f(\cdot)]$

Similarly when $y_i = -1$ and $|\hat{f}(x_i)| \gg 1$ and $f(x_i) < 0$

$$\begin{aligned} \nabla_i \log(p(y|\hat{f}_x)) &= -\frac{1+1}{2} - \frac{1}{1 + \exp(-f(x_i))} \\ &= 0 - \underbrace{\frac{1}{1 + \exp(-f(x_i))}}_{\text{lim as } -f(x_i) \rightarrow \infty} = 0 \end{aligned}$$

When the signs of $y_i \& f(x_i)$ agree that is

$y_i f(x_i) > 0$ and $|\hat{f}(x_i)| \gg 1$ $\sigma(y_i f(x_i)) \rightarrow 1$ these are highly likely training points as given by logistic likelihood at the same time the gradient of their log-logistic likelihood

$$\nabla_i \log(p(y|\hat{f}_x)) \approx 0$$

thus those training point will

not contribute to update equation. as the gradient of log likelihood

$$E_g[f(\cdot)] = m_0 + \kappa_x \nabla \log p(y|\hat{f}_x) \text{ these component is zero.}$$

1(e) This observation is reminiscent of the notion of support vectors in SVMs, where only a few training points contribute to the solution. In SVMs this is achieved using the hinge loss

$$l(y_i; f(x_i)) = \max(0, 1 - y_i f(x_i)) \text{ instead of the logistic loss.}$$

Because this loss has zero gradient whenever $y_i f(x_i) > 1$

the associated training pairs do not contribute at all to the optimization problem, which can be leveraged for efficient numerical optimization. It would thus be tempting to use:

$$\begin{aligned} r(y_i; f(x_i)) &= \exp(-l(y_i; f(x_i))) \text{ as a likelihood.} \\ &= \exp(-\max(0, 1 - y_i f(x_i))) \end{aligned}$$

as a likelihood in our GP classification model. Unfortunately, this is not a valid likelihood.

To make this clear show that ; there is no constant c such that.

$$\sum_{y_i \in \{-1, 1\}} c r(y_i; f(x_i)) = 1 \quad \forall f(x_i) \in \mathbb{R}$$

thus $r(y_i; f(x_i))$ is not a family of probability distributions, and can't be a likelihood.

Answer: We consider the normalization of suggested likelihood.

$$\sum_{y_i \in \{-1, 1\}} c \cdot r(y_i; f(x_i))$$

$$\text{where } r(y_i; f(x_i)) = \exp[-\max(0, 1 - y_i f(x_i))]$$

we separate out +ve and -ve classifications.

s.t.

$$\sum_{\substack{\text{where} \\ y_i = 1}} c \exp[-\max(0, 1 + f(x_i))] + \sum_{\substack{\text{where} \\ y_i = -1}} c \exp[-\max(0, 1 - f(x_i))]$$

Without loss of generality consider training points in +ve and negative classes with same $f(x_i) = f'$ value f' . Then either $f' \geq 1$ or $f' \leq 1$

Let $f' > 1$

without loss of generality consider training points in the +ve and -ve classes with $f(x_i) = f'$

then for two such points we have.

$$S = \underbrace{c \exp[-\max(0, 1 + f')]}_{y_i = -1} + \underbrace{c \exp[-\max(0, 1 - f')]}_{y_i = 1}$$

without loss of generality let $f' > 1$
then $\max(0, 1 + f') = 1 + f'$
 $\max(0, 1 - f') = 0$

$$S = c \left[\exp[-(1+f')] + P^o \right]$$

$$= c \left[1 + e^{-\frac{1}{1+f'}} \right]$$

The only normalizer for S is when

$$c = \frac{1}{1 + e^{-(1+f')}} \quad \text{but this is not independent of } f(x_i)$$

thus $r(y_i, f(x_i))$ is not a family of probability distributions, and can't be a likelihood.