# Theory Question

a)

$$p(f_x \mid y, X) = \frac{p(y \mid f_x, X) \, p(f_x \mid X)}{p(y \mid X)}, \text{ using Bayes' Rule}$$

$$p(f \mid y) = \frac{p(y \mid f) \, p(f)}{p(y)} \quad \longrightarrow \quad \text{Without loss of generality for continuous } X$$

$$\log(p(f_x \mid y, X)) = \log p(y \mid f_x) + \log p(f_x) - \log p(y)$$

$$= \log \prod_{i=1}^{N} \sigma(y_i \, f(x_i)) + \log p(f_x) + \text{constant term}$$

$$p(f_x) = \mathcal{N}(\underbrace{f_x}_{m_x} ; m(x), \underbrace{k(X,X)}_{k_{xx}})$$

$$\log(p(f_x \mid y, X)) \cong \left( -\sum_{i=1}^{N} \log(1 + \exp(-y_i f(x_i))) - \frac{1}{2}(f_x - m_x)^T k_{xx}^{-1} (f_x - m_x) \right)$$

b) $\nabla_{fx} \log(p(f_x | y, x)) = \sum_{i=1}^{N} \nabla_{fx} \log(\sigma(y_i, f(x_i))) - k_{xx}^{-1}(f_x - m_x)$

$$= \sum_{i=1}^{N} -\nabla_{fx} \underline{\log(1 + \exp(-y_i f_{x_i}))} - k_{xx}^{-1}(f_x - m_x)$$

where $-\dfrac{\partial(\log(1+\exp(-y_i f_{x_i})))}{\partial f_{x_j}} = \delta_{ij}\left(\dfrac{y_i+1}{2} - \sigma(f_{x_i})\right)$

where $y_i \in \{-1, 1\}$

$\Rightarrow \dfrac{y_i+1}{2} \in \{0, 1\}$

Note that $-\dfrac{\partial}{\partial f_{x_i}} \log(1 + \exp(-y_i f_{x_i})) = \dfrac{y_i \exp(-y_i f_{x_i})}{1 + \exp(-y_i f_{x_i})}$

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

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

$\begin{matrix}①\\②\end{matrix} \Rightarrow \dfrac{y_i+1}{2} - \sigma(f_{x_i})$

$\nabla_{fx} \log(p(f_x | y)) = \begin{bmatrix} \frac{y_1+1}{2} - \sigma(f(x_1)) \\ \frac{y_2+1}{2} - \sigma(f(x_2)) \\ \vdots \\ \frac{y_n+1}{2} - \sigma(f(x_n)) \end{bmatrix} - k_{xx}^{-1}(f_x - m_x)$

## c)

We know that at the mode, $f_x = \hat{f}_x$ :

$$\nabla_{f_x} \log p(f_x \mid y, x)\big|_{f_x = \hat{f}_x} = 0$$

Using part b $\rightarrow$ $\nabla_{f_x} \log(p(y \mid f_x)) = k_{xx}^{-1}(f_x - m_x) = 0$ , at the mode

$$\circledast \quad \nabla \log(p(y \mid \hat{f}_x)) = k_{xx}^{-1}(\hat{f}_x - m_x) \quad , \text{ at the mode}$$

$$k_{xx} \nabla \log(p(y \mid f_x)) = \underbrace{k_{xx} k_{xx}^{-1}}_{=1, \text{ spd}}(f_x - m_x) \quad , \text{ at the mode}$$

**Way 1**

$$\Rightarrow \quad m_x + k_{xx} \nabla \log(p(y \mid f_x)) = f_x \quad , \text{ at the mode}$$

From the last result, it can be seen that

$$E_q[f(\cdot)] = m_\cdot + k_{\cdot x} \nabla \log(p(y \mid \hat{f}_x))$$

**Way 2**

More rigorously, as depicted also in Lecture 14, slide 17:

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

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

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

Thus, $E_q[f(\cdot)] = m(\cdot) + k_{\cdot x} k_{xx}^{-1}(\hat{f}_x - m_x)$ as expected value
is precisely the mean of gaussian.

Using $\circledast$ to replace $k_{xx}^{-1}(\hat{f}_x - m_x)$ with $\nabla \log(p(y \mid \hat{f}_x))$, the MAP estimate can be written as:

$$\boxed{E_q[f(\cdot)] = m(\cdot) + k_{\cdot x} \nabla \log(p(y \mid \hat{f}_x))}$$

Q.E.D.

**d)** Utilizing the resulte of part b, the explicit form can be written as:

$$\nabla \log(p(y)\,\hat{f}x) = \sum_{i=1}^{N} \nabla \hat{f}x \,\log(\sigma(y_i \hat{f}(x_i)))$$

$$= \sum_{i=1}^{N} \nabla \hat{f}x\,[-\log(1+\exp(-y_i \hat{f}(x_i)))]$$

with $\quad \dfrac{\partial}{\partial \hat{f}x_i}(-\log(1+\exp(-y_i \hat{f}(x_i)))) = \dfrac{y_i+1}{2} - \sigma(\hat{f}(x_i))$

and $\quad \dfrac{\partial \log \sigma(y_i \hat{f}(x_i))}{\partial \hat{f}x_j} = \delta_{ij}\left(\dfrac{y_i+1}{2} - \sigma(\hat{f}(x_i))\right)$

Observe that when $|\hat{f}(x_i)| \gg 1$, there are two cases: either the $\hat{f}(x_i)$ is too small or too large:

① when $\hat{f}(x_i) \gg 1$, $\sigma(\hat{f}(x_i)) = \dfrac{1}{1+\exp(-\hat{f}(x_i))} \approx 1$ as $\lim\limits_{x\to\infty} \exp(-x) = 0$

and the gradient term becomes, $\dfrac{y_i+1}{2} - \sigma(\hat{f}(x_i)) = \dfrac{1+1}{2} - 1 = 0$

② when $\hat{f}(x_i) \ll 1$, $\sigma(\hat{f}(x_i)) = \dfrac{1}{1+\exp(-\hat{f}(x_i))} \approx 0$ as $\lim\limits_{x\to-\infty} \exp(-x) = \infty$

and the gradient term becomes, $\dfrac{y_i+1}{2} - \sigma(\hat{f}(x_i)) = \dfrac{-1+1}{2} - 0 = 0$

therefore, training points far from the decision boundary do almost not contribute to this estimate $E_q(f(\cdot))$.

e)

We consider the normalization of the suggested likelihood

$$\sum_{y_i \in \{-1,1\}} c \, r(y_i; f(x_i)) \quad \text{with} \quad r(y_i; f(x_i)) = \exp[-\max(0, 1 - y_i f(x_i))]$$

Separate negatives and positives

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

without loss of generality, consider for individual training points in the positive and negative classes. Then for two such points, we have

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

Without loss of generality, consider the case $f(x_i) > 1$, then:

$$c \exp[-(1 + f(x_i))] + c \exp[0] = c[1 + e^{-(1 + f(x_i))}]$$

The normalization constant, $c$ should include the term:

$$c = c_1 \frac{1}{1 + e^{-(1 + f(x_i))}} \qquad \text{yet this is not independent of } f(x_i)$$

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