

### Question 1

(a) From Bayes' theorem,  $p(t|x) \propto p(x|t) p(t)$

$$p(x) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h} k\left(\frac{x-x_n}{h}\right) \quad p(x|t) = \frac{1}{N_t} \sum_{n=1}^{N_t} \frac{1}{2\kappa} k(x, x_n) \delta(t, t_n)$$

$N_t$  is the number of input vectors with label  $t$  ( $+1$  or  $-1$ ) and  $N = N_1 + N_{-1}$ .  $\delta(t, t_n)$  equals 1 if  $t = t_n$  and 0 otherwise.

$2\kappa$  is the normalization constant for the kernel. The minimum misclassification rate is achieved if, for each new input vector,  $\hat{x}$ , we chose  $\hat{t}$  to maximize  $p(t|\hat{x})$ . With equal class priors, this is equivalent to maximizing  $p(t|\hat{x})$  and thus

$$\hat{t} = \begin{cases} +1 & \text{if } \frac{1}{N_t} \sum_{n=1}^{N_t} k(\hat{x}, x_n) \geq \frac{1}{N_{-1}} \sum_{j=1}^{N_{-1}} k(\hat{x}, x_j) \\ -1 & \text{otherwise} \end{cases}$$

Dropped the factor  $1/2\kappa$ , because it's a scaling factor. Use the encode scheme.

$$\hat{t} = \text{sign}\left(\sum_{n=1}^{N_t} \frac{t_n}{N_t} k(\hat{x}, x_n)\right)$$

$$(b) \text{ let } K(x, x_n) = x^T x_n \quad p(x|t=+1) = \frac{1}{N_t} \sum_{n: t_n=+1} x^T x_n = x^T x^+$$

The sum in the middle expression runs over all vectors  $x_n$  for which  $t_n = +1$  and  $x^+$  indicates the mean of these vectors, and similar definition for the negative class.

$$\hat{t} = \begin{cases} +1 & \text{if } x^T x^+ \geq x^T x^- \\ -1 & \text{otherwise} \end{cases}$$

(c) The same argument also applies for feature space  $\phi(x)$

$$\hat{t} = \begin{cases} +1 & \text{if } \phi(\hat{x})^T \phi(x^+) \geq \phi(\hat{x})^T \phi(x^-) \\ -1 & \text{otherwise} \end{cases}$$

### Question 2

$$p(t=y) = \sigma(y) \Rightarrow p(t=-1|y) = 1 - p(t=+1|y) = 1 - \sigma(y) = \sigma(-y)$$

data  $D = \{(t_1, x_1), \dots, (t_N, x_N)\}$

$$p(D) = \prod_{n=1}^N \sigma(y_n) \prod_{n=1}^N \sigma(-y_n) = \prod_{n=1}^N \sigma(t_n y_n)$$

where  $y_n = y(x_n)$ . Taking the negative logarithm of this, we get

$$-\ln p(D) = -\ln \prod_{n=1}^N \sigma(t_n y_n) = \sum_{n=1}^N \ln \sigma(t_n y_n) = \sum_{n=1}^N \ln (1 + e^{-t_n y_n})$$

Combining with regularization term  $\lambda \|w\|^2$ , we get

$$\sum_{n=1}^N \text{Elo}(y_n) + \lambda \|w\|^2$$

### Question 3

$$p(t|x, w, \beta) = \prod_{n=1}^N p(t_n|x_n, w, \beta^{-1})$$

$$p(w|d) = \prod_{i=1}^I N(w_i | 0, \alpha_i^{-1})$$

$$p(t|x, w, \beta) = \int p(t|x, w, \beta) p(w|d) dw$$

$$p(t|x, w, \beta) = \left(\frac{\beta}{2\pi}\right)^{M/2} \frac{1}{\det(\Sigma)} \prod_{i=1}^I \int \exp\{-\text{E}(w_i)\} dw$$

where  $\text{E}(w) = \frac{1}{2} (w - m)^T \Sigma^{-1} (w - m) + \text{E}(c) = A \cdot \text{diag}(a)$

$$\text{E}(w) = \frac{1}{2} (w - m)^T \Sigma^{-1} (w - m) + \text{E}(c) \quad m = \beta \beta^T c \quad \Sigma = (A + \beta \beta^T)^{-1}$$

$$\text{E}(c) = \frac{1}{2} (\beta \beta^T + m^T \Sigma m)$$

$$\int \exp\{-\text{E}(w)\} dw = \exp\{-\text{E}(c)\} (2\pi)^{\frac{M}{2}} |I|^{-\frac{1}{2}}$$

$$C = \beta \beta^T + \frac{1}{2} A^{-1} \Sigma A$$

$$\text{Using } (A + BD^T C)^{-1} = A^{-1} - A^{-1} B C D^{-1} C^T B^{-1}$$

$$\text{We can rewrite } \text{E}(c) = \frac{1}{2} (\beta \beta^T + m^T \Sigma m) = \frac{1}{2} (\beta \beta^T + \beta \beta^T \Sigma \Sigma^{-1} \Sigma \beta^T + \beta \beta^T)$$

$$= \frac{1}{2} \beta^T (\beta^T - \beta \beta^T \Sigma \Sigma^{-1} \beta) \beta = \frac{1}{2} \beta^T (\beta^T - \beta \beta^T \Sigma \Sigma^{-1} \beta)^T \beta = \frac{1}{2} \beta^T C^{-1} \beta$$

This is the last term of the right hand side. The first and second terms are given implicitly, as they form the normalization constant for the posterior Gaussian distribution  $p(t|x, w, \beta)$ .

### Question 4

Using results from Question 3

$$\text{we can rewrite } \ln(p(t|x, w, \beta)) = \ln N(t|0, C) = -\frac{1}{2} \{N \ln(2\pi) + \ln|C| + t^T C^{-1} t\} = \frac{1}{2} \ln \beta + \frac{1}{2} \sum_{i=1}^I \ln(a_i) - \frac{1}{2} \ln|I| - \frac{1}{2} \ln(2\pi)$$

By taking the derivative w.r.t.  $\alpha$ :

$$\frac{\partial}{\partial \alpha_i} \ln(p(t|x, w, \beta)) = \frac{1}{2\alpha_i} - \frac{1}{2} \sum_{i=1}^I -\frac{1}{2} \alpha_i^{-2}$$

Setting the derivative to zero

$$\alpha_i = \frac{|\phi(x_i)|}{\sigma_i^2} = \frac{\gamma_i}{\sigma_i^2} \quad \gamma_i = |\phi(\Sigma_{ii})|$$

$$\frac{\partial}{\partial \beta} \ln p(t|x, \alpha, \beta) = \frac{1}{2} \left( \frac{N}{\beta} - \|t - \phi w\|^2 - \text{Tr}[\Sigma \phi^T \phi] \right)$$

$$\begin{aligned} \Sigma \phi^T \phi &= \Sigma \phi^T \phi + A^{-1} \Sigma A - \beta^{-1} \Sigma A \\ &= \Sigma (\phi^T \phi + A \beta^{-1} - \beta^{-1} \Sigma A) \\ &= (I + \beta \phi^T \phi + A \beta^{-1} - \beta^{-1} \Sigma A) \\ &= (I - A \Sigma) \beta^{-1} \end{aligned}$$

$(I - A \Sigma)$  equals  $\gamma_i^{-1} \alpha_i \Sigma_{ii}$  in matrix form

Set  $\frac{\partial}{\partial \beta} \ln p(t|x, \alpha, \beta) = 0$   
 By rearranging  $(\beta^{\text{new}})^{-1} = \frac{\|t - \phi w\|^2}{N - \sum \gamma_i \alpha_i}$

### Question 5

(a)  $K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = \langle \phi(x_j), \phi(x_i) \rangle = K(x_j, x_i)$

(b)  $\| \phi(x_i) - \phi(x_j) \|^2$   
 $= \langle \phi(x_i), \phi(x_i) \rangle + \langle \phi(x_j), \phi(x_j) \rangle - 2 \cdot \langle \phi(x_i), \phi(x_j) \rangle = K(x_i, x_i) + K(x_j, x_j) - 2 \cdot K(x_i, x_j) = (1 - 2 \exp(-\frac{1}{2} \|x_i - x_j\|^2)) \leq 2$

### Question 6

$\|x_{\text{far}} - x_i\| \gg 0 \quad \forall i \in SV$

$\Rightarrow K(x_{\text{far}}, x_i) \approx 0 \quad \forall i \in SV$

$\Rightarrow \sum_{i \in SV} y_i \alpha_i K(x_i, x) \approx 0$

$\Rightarrow f(x, \alpha, \hat{w}_0) \approx \hat{w}_0$