

Machine Learning Final Exam Date: Jan 7, 2021 Time: 9:00-12:10
 Name: Leung Ko Tsun (梁高浚) SID: 71092303
 電考: 陳芳

$$Q6a). L(\alpha, b, \xi, w, \beta, \gamma) = \frac{1}{2} \sum_{i=1}^N \alpha_i + \sum_{i=1}^N C_i \xi_i + \sum_{i=1}^N w_i \cdot (1 - \xi_i - y_i) \left(\sum_{j=1}^N \alpha_j y_j x_{ij} \cdot x_{ij} + b \right) \\ + \sum_{i=1}^N p_i (-\alpha_i) + \sum_{i=1}^N \gamma_i (-\xi_i)$$

b). By strong Duality theorem, [if X is a nonempty convex set in \mathbb{R}^N .
 If f, g_1, g_2, g_3 are convex, and $\exists x \in X$ st. $g_1(x) < 0, g_2(x) < 0$ and
 $g_3(x) < 0$, then $\inf f(x) : x \in X, g_1(x) \leq 0, g_2(x) \leq 0, g_3(x) \leq 0 \} = \sup \{ Q(w, \beta, \gamma) : w \geq 0, \beta \geq 0, \gamma \geq 0 \}$,
 $w, \beta, \gamma \in \mathbb{R}^N$]

We want to ensure these conditions are satisfied in order to apply Strong Duality theorem.

Obviously, X is nonempty. X is also convex set as $tx + (1-t)x$ is contained in X .
 for function f , $\sum_{i=1}^N \alpha_i$ is quadratic hence convex, $\sum_{i=1}^N C_i \xi_i$ is affine

hence convex.

for function g_{11} , it is also convex as second derivative ≥ 0 .
 $\cup g_{21}, g_{31}$ they are also convex as they are just affine functions.

and $\exists x \in X$ st. $g_{11}(x) < 0, g_{11}(x) < 0$ and $g_{11}(x) < 0$ can be also

satisfied if α_j, b picked to be 0 and $\xi_j \geq 1$.

$g_{21}(x), g_{31}(x)$ can be < 0 by picking positive α_i and ξ_i .

\Rightarrow therefore, by strong duality theorem, the duality gap ≥ 0 .

(a6c).

$$\nabla_{\mathbf{x}} L = \mathbf{z} - \sum_{i=1}^N \sum_{j=1}^N w_i y_j y_j x_i x_j - \sum_{i=1}^N \beta_i$$
$$\frac{\partial L}{\partial b} = -\sum_{i=1}^N w_i y_i$$

$$\frac{\partial L}{\partial z_i} = c_i - w_i - \gamma_i$$

(*) $\sum_{i=1}^N w_i y_i = 0, w_i, y_i = c_i, i=1, \dots, N$

If (*) is satisfied, then $\Omega(w, \beta, \gamma) = L(\alpha, b, \mathbf{z}, w, \beta, \gamma)$ iff $\vec{\alpha} = \sum_{i=1}^N \sum_{j=1}^N w_i y_j y_j x_i x_j + \sum_{i=1}^N \beta_i$

otherwise, one can make $\Omega(w, \beta, \gamma)$ to be $-\infty$

$$\begin{aligned} \text{so } \Omega(w, \beta, \gamma) &= \frac{1}{2} \sum_{i=1}^N \alpha_i^2 + \sum_{i=1}^N c_i \alpha_i + \sum_{i=1}^N w_i (1 - \alpha_i - y_i (\sum_{j=1}^N w_j y_j x_i x_j + \sum_{j=1}^N \beta_j)) + \sum_{i=1}^N \beta_i (\alpha_i) + \sum_{i=1}^N \gamma_i (\alpha_i) \\ &= \frac{1}{2} \sum_{i=1}^N \alpha_i^2 + \sum_{i=1}^N w_i - \cancel{\sum_{i=1}^N \sum_{j=1}^N w_i y_j y_j x_i x_j + \sum_{j=1}^N \beta_j} \\ &= \sum_{i=1}^N w_i - \frac{1}{2} \|\mathbf{x}\|^2 \\ &= \sum_{i=1}^N w_i - \frac{1}{2} \left(\sum_{i=1}^N \sum_{j=1}^N w_i y_j y_j x_i x_j + \sum_{j=1}^N \beta_j \right)^T \left(\sum_{i=1}^N \sum_{j=1}^N w_k y_k y_k x_k x_k + \sum_{k=1}^N \beta_k \right) \end{aligned}$$

d). Simplifying the above, we have:

(continued)

$$\text{maximize } \Omega(w, \beta, \gamma) = \sum_{i=1}^N w_i - \frac{1}{2} \left(\sum_{i=1}^N \sum_{j=1}^N w_i y_j y_j x_i x_j + \sum_{j=1}^N \beta_j \right)^T \left(\sum_{i=1}^N \sum_{j=1}^N w_i y_j y_j x_i x_j + \sum_{j=1}^N \beta_j \right)$$

$$\Leftrightarrow \text{maximize } \Omega(w, \beta, \gamma) = \sum_{i=1}^N w_i - \frac{1}{2} \cdot \sum_{i=1}^N \left(\sum_{j=1}^N y_i y_j x_j \cdot x_i \right)^2,$$

equivalent to

(as $\beta_i \geq 0$ is one of the condition in Lagrangian dual problem)

And also, the conditions $w_i + \beta_i = c_i, w_i \geq 0, y_i \geq 0$

can be simplified to

$$0 \leq w_i \leq c_i, i=1, \dots, N.$$

Therefore, the dual problem can be simplified as:

$$\text{maximize } \sum_{i=1}^N w_i - \frac{1}{2} \sum_{i=1}^N \left(\sum_{j=1}^N w_j y_j x_j \cdot x_i \right)^2$$

$$\text{subject to } \sum_{i=1}^N w_i y_i = 0, \quad 0 \leq w_i \leq c_i, i=1, \dots, N$$

Q66) Explicit form:

$$\Theta(w, \beta, r) = \sum_{i=1}^N w_i - \frac{1}{2} \sum_{i,j,k,l} w_i w_j y_i y_k y_l x_i \cdot x_j x_k x_l - \sum_{i=1}^N \sum_{j=1}^N w_i y_i x_i \cdot x_j \beta_j - \frac{1}{2} \sum_{j=1}^N \beta_j^2$$

Q66, i.e., if $y_i(\bar{w} \cdot x_i + b) < 1$, then $w_i > 0$ by QP, $\beta_i = 0$ by C₂, $\bar{w} = g_i$ by QP.

Q6 (c), (i).

Note that the dual problem can be rewritten as: (from reordering term in grad
g_i for b_i)

$$\text{minimize}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N C_i \max(1 - y_i (\sum_{j=1}^N w_j y_j x_i \cdot x_j + b_i), 0)$$

variables $\mathbf{w} \in \mathbb{R}^N, b \in \mathbb{R}$.

Given the optimal weighting vector \mathbf{w} , the optimal bias is given by

$$b = \arg \min_{b \in \mathbb{R}} \sum_{i=1}^N C_i \max(1 - y_i (\mathbf{w} \cdot \mathbf{x}_i + b), 0)$$

$(\bar{\mathbf{w}}, \bar{b}, \bar{\mathbf{z}})$, $(\bar{\mathbf{w}}, \bar{b}, \bar{\mathbf{y}})$ are primal/dual optimal solutions iff KKT conditions hold.

KKT conditions:

S1: Stationary condition:

$$\sum_{j=1}^N w_j y_j = 0$$

$$S3: \mathbf{z} = \sum_{i=1}^N (1 - y_i (\mathbf{w} \cdot \mathbf{x}_i + b_i)) \geq 0$$

Primal feasibility:

$$P1: y_i (\sum_{j=1}^N w_j y_j x_i \cdot x_j + b_i) \geq 1 - z_i$$

$$P2: \epsilon_i \geq 0$$

$$P3: \alpha_i \geq 0$$

$$D1: w_i \geq 0$$

$$D2: b \geq 0$$

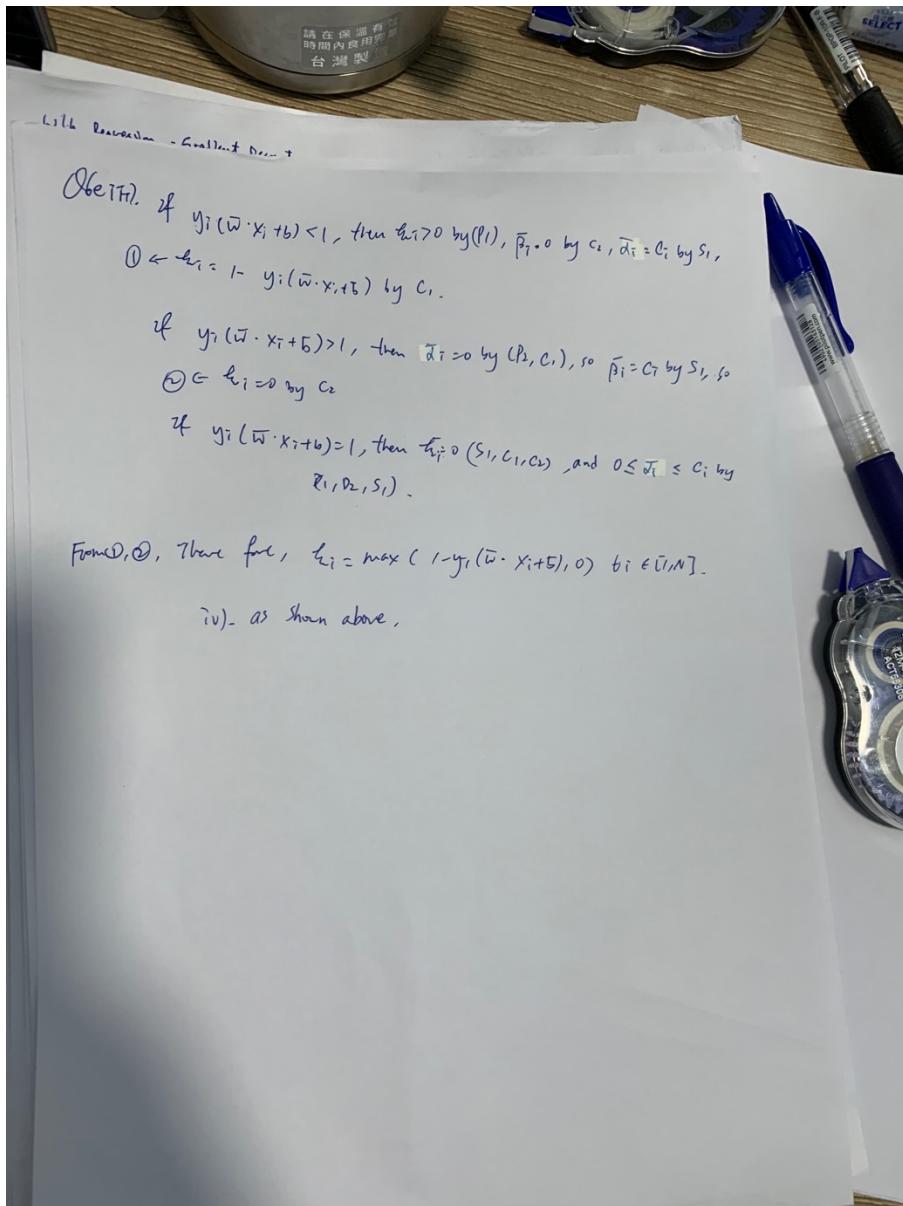
Complementary slackness

$$C_1: w_i (1 - z_i - y_i (\sum_{j=1}^N w_j y_j x_i \cdot x_j + b_i)) = 0$$

$$C_2: \epsilon_i \alpha_i = 0$$

i). From S3, we can see that $\alpha_i = \max(\sum_{j=1}^N w_j y_j y_i x_i \cdot x_j + b_i, 0)$
and $b_i = \max(\sum_{j=1}^N w_j y_j y_i x_i \cdot x_j, 0)$ since $\beta_1 > \beta_2$.

ii). From P1 and rewritten primal problem, we can get T as the desired form



(6.2.1). The joint probability function

$$= \prod_{i=1}^n p_\theta(\varepsilon | x_i)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\varepsilon - f_\theta(x_i))^2}{2\sigma^2}\right)$$

$$\approx (2/\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (\varepsilon - f_\theta(x_i))^2\right)$$

The log-likelihood function

$$\ell = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 = -\frac{1}{2\sigma^2} \sum_{i=1}^n (\varepsilon - f_\theta(x_i))^2$$

Maximizing the log-likelihood is equivalent to minimizing

$$\sum_{i=1}^n (\varepsilon - f_\theta(x_i))^2$$

$$\text{or simply } \sum_{i=1}^n (\varepsilon - f_\theta(x_i))^2$$

$$\therefore L_{ML}(\theta) = \sum_{i=1}^n (\varepsilon - f_\theta(x_i))^2$$

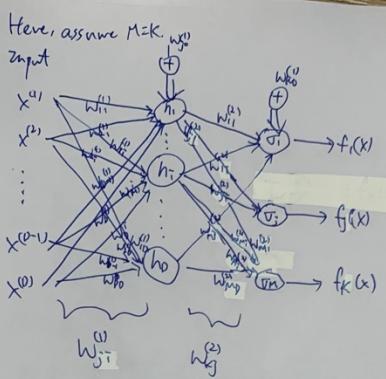
$$\begin{aligned} b). P(\theta | D_x) P(D_y | \theta, D_x) &= \prod_{i=1}^k N(\theta^{(i)}; \theta, \Sigma) \prod_{i=1}^N \text{Inv}(y_i; f_\theta(x_i), \theta^2) \\ &= \prod_{i=1}^k \frac{1}{\sqrt{2\pi\Sigma}} \exp\left(-\frac{(\theta^{(i)} - \theta)^2}{2\Sigma}\right) \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - f_\theta(x_i))^2}{2\sigma^2}\right) \\ &\propto \left(\frac{1}{\sqrt{2\pi\Sigma}}\right)^k \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^N \exp\left(-\frac{1}{2\Sigma} \sum_{i=1}^k (\theta^{(i)})^2 - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_\theta(x_i))^2\right) \end{aligned}$$

$$P(\theta | D) = \frac{P(\theta, D_y | D_x)}{P(D_y | D_x)} \propto P(\theta, D_y | D_x) = P(\theta | D_x) P(D_y | \theta, D_x)$$

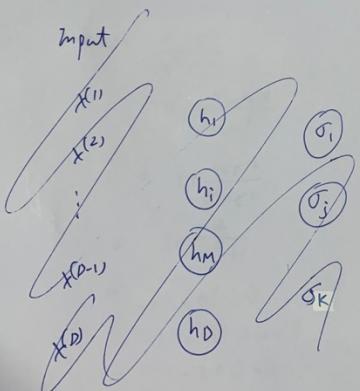
\exp(-L_{Bayes}(\theta))

Hence proved.

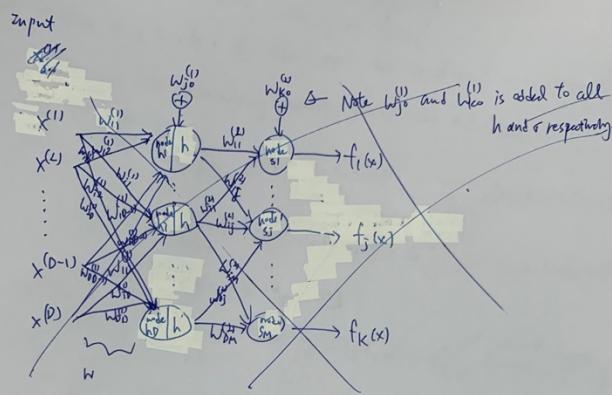
Q(1a).



j denote index in x k denote index in $f(x)$
 i denote index in h .



(Q1a).



$$\begin{aligned}
 b). \quad \tanh(z) &= \frac{e^z - e^{-z}}{e^z + e^{-z}} \\
 &= \frac{e^z + e^{-z} - 2e^{-z}}{e^z + e^{-z}} \\
 &= 1 + \frac{-2e^{-z}}{e^z + e^{-z}} \\
 &= 1 - \frac{2}{e^{2z} + 1} \\
 &= 1 - 2h(-2z) \\
 &= 1 - 2(1 - h(2z)) \\
 &= 1 - 2 + 2h(2z) \\
 &= 2h(2z) - 1
 \end{aligned}$$

So, they differ by linear transformation. We can transform the ~~tanh~~ by the above formula $2h(2z) - 1$, therefore there exist a ~~real~~ NN s.t. they perform the same function.

(Q3): Our goal is to minimize L wrt. both α_t and $f_t(x)$.
 We suppose f_1, \dots, f_{t-1} are fixed, and their coefficients $\alpha_1, \dots, \alpha_{t-1}$,
 so we are minimizing L only wrt. α_t and $f_t(x)$.

$$L = C_1 \sum_{i, y_i=1} e^{y_i g_{t-1}(x_i) - \alpha_t f_t(x_i)} + C_2 \sum_{i, y_i=-1} e^{g_{t-1}(x_i) - \alpha_t f_t(x_i)}$$

$$= C_1 \sum_{i, y_i=1} w_i^{(t)} \exp\{-\alpha_t f_t(x_i)\} + C_2 \sum_{i, y_i=-1} w_i^{(t)} \exp\{-\alpha_t f_t(x_i)\}$$

where $w_i^{(t)} = \exp\{-g_{t-1}(x_i)\}$

$$\min_{\alpha_t} L \quad \frac{\partial L}{\partial \alpha_t} = \cancel{2e^{-\alpha_t}} + \cancel{\sum_{i=1}^n w_i^{(t)} I(f_t(x_i) \neq y_i)} - \cancel{2e^{-\alpha_t + \sum_{i=1}^n w_i^{(t)}}} \cancel{w_i^{(t)}}$$

Solving the above equation, we have

$$\alpha_t = \log \left(\frac{\sum_{i, y_i=1} w_i^{(t)} f_t(x_i)}{\sum_{i, y_i=-1} w_i^{(t)} f_t(x_i)} \right)$$

$$\frac{\partial L}{\partial \alpha_t} = -C_1 \sum_{i, y_i=1} w_i^{(t)} f_t(x_i) \exp\{-\alpha_t f_t(x_i)\} - C_2 \sum_{i, y_i=-1} w_i^{(t)} f_t(x_i) \exp\{-\alpha_t f_t(x_i)\} > 0$$

Solving above equation,

$$-C_1 \sum_{i, y_i=1} w_i^{(t)} f_t(x_i) = C_2 \sum_{i, y_i=-1} w_i^{(t)} f_t(x_i)$$

$$C_1 = \frac{-C_2 \sum_{i, y_i=-1} w_i^{(t)} f_t(x_i)}{\sum_{i, y_i=1} w_i^{(t)} f_t(x_i)}$$

(Q4.

$$\begin{aligned} \text{Given } |f'(x) - f'(y)| &\leq \gamma|x-y|, \\ |f(x) - f(a) - f'(a)(x-a)| &= \left| \int_a^x f'(t) dt - \int_a^x f'(a) dt \right| \\ &= \left| \int_a^x (f'(t) - f'(a)) dt \right| \\ &\leq \int_a^x |f'(t) - f'(a)| dt \\ &\leq \int_a^x \gamma(t-a) dt \\ &= \gamma(x-a)^2/2 \end{aligned}$$

b). Apply g to (a) with $x = x_k$ and $a = x_{k-1}$,

$$\begin{aligned} |f(k) - f(x_{k-1}) - \nabla f(x_{k-1}) \cdot (x_k - x_{k-1})| &\leq \frac{\gamma(x_k - x_{k-1})^2}{2} \\ f(k) &\leq f(x_{k-1}) + \nabla f(x_{k-1}) \cdot (x_k - x_{k-1}) + \frac{\gamma(x_k - x_{k-1})^2}{2} \\ &\leq f(x_{k-1}) + \nabla f(x_{k-1}) \geq (-\eta \nabla f(x_{k-1}) + \frac{\gamma(-\eta \nabla f(x_{k-1}))^2}{2}) \\ &\geq f(t_{k-1}) - \eta(1 - \gamma\eta/2) \| \nabla f(x_{k-1}) \|^2_2 \end{aligned}$$

(Q5) Posterior prob. dist. of latent variables z_i based. on current parameters $\theta^{(t)}$:

$$P[z_i=k|x_i; \theta^{(t)}] = \frac{p(x_i, z_i=k; \theta^{(t)})}{\sum_{j=1}^k P(x_i, z_i=j; \theta^{(t)})} = s_{ik}^{(t)}$$

log-likelihood function:

$$\log p(x_i|\theta) = \sum_{j=1}^N \log \left(\pi_k \frac{\partial}{\partial z_j} \mu_{kj}^{(t)} (1-\mu_{kj}^{(t)})^{1-x_{ij}} \right)$$

log-likelihood of parameter θ given data $x_{i,j}$ and latent variable z_i

$$\begin{aligned} \log p(x_i, z_i=k; \theta) &= \log \pi_k \frac{\partial}{\partial z_j} \mu_{kj}^{(t)} (1-\mu_{kj}^{(t)})^{1-x_{ij}} \\ &= \pi_k \sum_{j=1}^k \log (\mu_{kj}^{(t)} (1-\mu_{kj}^{(t)})^{1-x_{ij}}) \\ &= \pi_k \sum_{j=1}^k [x_{ij} \log (\mu_{kj}) + (1-x_{ij}) \log (1-\mu_{kj})] \end{aligned}$$

Expectation of log-likelihood:

$$\begin{aligned} Q(\theta|\theta^{(t)}) &= \sum_{i=1}^N E_{z_i|x_i, \theta^{(t)}} [\log p(x_i, z_i; \theta)] \\ &= \sum_{i=1}^N \sum_{k=1}^K \frac{\partial}{\partial z_k} s_{ik}^{(t)} [x_{ij} \log (\mu_{kj}) + (1-x_{ij}) \log (1-\mu_{kj})] \end{aligned}$$

M-step: choose $\theta^{(t+1)} = \underset{\theta}{\operatorname{argmax}} Q(\theta|\theta^{(t)})$

partial derivative over μ_{kj} ,

$$\frac{\partial}{\partial \mu_{kj}} Q(\theta|\theta^{(t)}) = \sum_{i=1}^N s_{ik}^{(t)} \left[\frac{x_{ij}}{\mu_{kj}} - \frac{1-x_{ij}}{1-\mu_{kj}} \right]$$

Setting derivatives to be 0,

$$\mu_{kj}^{(t+1)} = \frac{\sum_{i=1}^N s_{ik}^{(t)} x_{ij}}{\sum_{i=1}^N s_{ik}^{(t)}}$$

Partial derivatives over π_k with Lagrange multiplier constrained in $\sum_{k=1}^K \pi_k = 1$

$$\nabla_{\pi_k} (\mathcal{L}(0|0^t) - \frac{\lambda}{\sum_{k=1}^K \pi_k - 1)) = \sum_{i=1}^N \frac{s_{ik}^{(t)}}{\pi_k} - \lambda$$

Setting derivatives = 0,

$$\pi_k^{(t+1)} = \lambda^{-1} \sum_{i=1}^N s_{ik}^{(t)}$$

the constraint $\sum_{k=1}^K \pi_k = 1$ implies $\lambda = \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^N s_{ik}^{(t)} = 1$

$$\pi_k^{(t+1)} = \frac{1}{N} \sum_{i=1}^N s_{ik}^{(t)}$$