

ML HW4 math

---

r14921A13 鄭皓中

---

---

---

---



# Problem |

$$(a) \min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

Constraint :  $y_i - w^T x_i - b \leq \varepsilon + \xi_i \Rightarrow y_i - w^T x_i - b - \varepsilon - \xi_i \leq 0$

$$w^T x_i + b - y_i \leq \varepsilon + \xi_i \Rightarrow w^T x_i + b - y_i - \varepsilon - \xi_i \leq 0$$

$$\xi_i \geq 0 \Rightarrow -\xi_i \leq 0$$

$$\begin{aligned} L(w, b, \xi, \alpha, \alpha^*) &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i (y_i - w^T x_i - b - \varepsilon - \xi_i) \\ &\quad + \sum_{i=1}^m \alpha_i^* (w^T x_i + b - y_i - \varepsilon - \xi_i) \\ &\quad + \sum_{i=1}^m \beta_i \cdot (-\xi_i) \end{aligned}$$

$$(b) \nabla_w L = w - \sum_{i=1}^m \alpha_i x_i + \sum_{i=1}^m \alpha_i^* x_i = 0 \Rightarrow w = \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^m \alpha_i (-1) + \sum_{i=1}^m \alpha_i^* \cdot 1 = 0 \Rightarrow \sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0$$

$$\frac{\partial L}{\partial \xi_i} = C - \alpha_i - \alpha_i^* - \beta_i = 0 \Rightarrow \beta_i = C - \alpha_i - \alpha_i^*$$

$$\begin{aligned} L &= \frac{1}{2} w^T w - w^T \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i + b \sum_{i=1}^m (\alpha_i^* - \alpha_i) + \sum_{i=1}^m \xi_i (C - \alpha_i - \alpha_i^* - \beta_i) - \varepsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*) + \sum_{i=1}^m (\alpha_i - \alpha_i^*) \gamma_i \\ &= \frac{1}{2} w^T w - w^T w + 0 + 0 + \sum_{i=1}^m (\alpha_i - \alpha_i^*) y_i - \varepsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*) \\ &= -\frac{1}{2} \sum_{i=1}^m \sum_{j \neq i} (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) x_i^T x_j + \sum_{i=1}^m (\alpha_i - \alpha_i^*) y_i - \varepsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*) \end{aligned}$$

Dual problem :  $\max_{\alpha, \alpha^*} -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i^* - \alpha_i) (\alpha_j^* - \alpha_j) x_i^T x_j + \sum_{i=1}^m (\alpha_i - \alpha_i^*) y_i - \varepsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*)$

Constraint :  $\sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0, 0 \leq \alpha_i + \alpha_i^* \leq C$

$$(c) \min_{b, \xi} C \sum_{i=1}^m \xi_i \quad \text{Constraint : } y_i - w^T x_i - b \leq \varepsilon + \xi_i, w^T x_i + b - y_i \leq \varepsilon + \xi_i, \xi_i \geq 0$$

Constraint :  $0 \leq |y_i - (w^T x_i + b)| \leq \varepsilon + \xi_i, \xi_i = 0$

$$\Rightarrow \xi_i \geq |y_i - (w^T x_i + b)| - \varepsilon, \xi_i \geq 0$$

$$\Rightarrow \xi_i \geq \max(0, |y_i - (w^T x_i + b)| - \varepsilon)$$

$$C \sum_{i=1}^m \xi_i = C \sum_{i=1}^m \left( \max(0, |y_i - (w^T x_i + b)| - \varepsilon) \right)$$

$$\min_{b, \xi} C \sum_{i=1}^m \xi_i \Rightarrow \bar{b} = \arg \min_{b \in \mathbb{R}} C \sum_{i=1}^m \max(0, |y_i - (w^T x_i + b)| - \varepsilon)$$

$$(c). \quad e = y_i - (\bar{w}^T x_i + b)$$

1.  $|e| < \varepsilon \Rightarrow$  不在 constraint  $|e| \leq \varepsilon$  的 boundary  $\Rightarrow$  Lagrange multiplier = 0

$$\bar{\alpha}_i = \bar{\alpha}_i^* = 0, \quad \beta_i \xi_i = 0 = (C - \bar{\alpha}_i - \bar{\alpha}_i^*) \xi_i \Rightarrow \xi_i = 0 \quad \forall i = 1, 2, \dots, m$$

2.  $e = \varepsilon \Rightarrow$  滿足  $e = \varepsilon$  的 boundary 但  $-e \leq \varepsilon$  依然沒有 active

$$\Rightarrow \bar{\alpha}_i^* = 0, \quad \xi_i = 0 \quad \forall i = 1, 2, \dots, m$$

由於  $\xi_i = 0 \Rightarrow \beta_i$  不一定為 0  $\Rightarrow C - \alpha_i = \beta_i \geq 0 \Rightarrow 0 \leq \bar{\alpha}_i \leq C$

3.  $e = -\varepsilon \Rightarrow$  滿足  $-e = \varepsilon$  的 boundary.  $e = \varepsilon$  則沒有 active

$$\Rightarrow \alpha_i = 0, \quad \xi_i = 0 \quad \forall i = 1, 2, \dots, m$$

由於  $\xi_i = 0 \Rightarrow \beta_i$  不一定為 0  $\Rightarrow C - \bar{\alpha}_i^* = \beta_i \geq 0 \Rightarrow 0 \leq \bar{\alpha}_i^* \leq C$

4.  $e > \varepsilon \Rightarrow$  引入  $\xi_i$  使得  $e = \varepsilon + \xi_i \Rightarrow -e < \varepsilon + \xi_i$  的 constraint 是 inactive

$$\Rightarrow \bar{\alpha}_i^* = 0, \quad \text{由於 } \xi_i > 0 \Rightarrow \sum_{i=1}^m \beta_i \xi_i = 0 \quad \text{且 } \beta_i \geq 0 \quad \text{得到 } \beta_i = 0$$

$$C - \bar{\alpha}_i = \beta_i = 0 \Rightarrow \bar{\alpha}_i = C$$

$$\Rightarrow \bar{\alpha}_i = C, \quad \xi_i = e - \varepsilon,$$

5.  $e < -\varepsilon \Rightarrow$  引入  $\xi_i > 0$  使得  $-e = \varepsilon - \xi_i \Rightarrow e < \varepsilon - \xi_i$  inactive

$$\Rightarrow \bar{\alpha}_i = 0, \quad \text{由於 } \xi_i > 0 \Rightarrow \sum_{i=1}^m \beta_i \xi_i = 0 \quad \text{且 } \beta_i \geq 0 \quad \text{得到 } \beta_i = 0$$

$$C - \bar{\alpha}_i^* = \beta_i = 0 \Rightarrow \bar{\alpha}_i^* = C$$

$$\Rightarrow \bar{\alpha}_i^* = C, \quad \xi_i = -(e + \varepsilon)$$

(d) <sup>(1)</sup> 由於 dual problem 只有  $x_i^T x_j$  影響  $\Rightarrow$  可替換成 kernel func.  $K(x_i, x_j) = x_i^T x_j$

$$\text{From b. dual problem: } \max_{\alpha_i, \alpha_i^*} -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) K(x_i, x_j) + \sum_{i=1}^m y_i (\alpha_i - \alpha_i^*) - \varepsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*)$$

$$\text{Constraint: } \sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0, \quad 0 \leq \alpha_i + \alpha_i^* \leq C$$

$$\text{決策函式 } f(x) = w^T x + b, \quad \text{Suppose } w = \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i$$

$$f(x) = \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i^T x + b = \sum_{i=1}^m (\alpha_i - \alpha_i^*) K(x_i, x) + b$$

截距  $b$ . 假設一個  $0 < \alpha_k < C$  的 support vector  $x_k$

$$\Rightarrow b = y_k - \varepsilon - \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i^T x_k = y_k - \varepsilon - \sum_{i=1}^m (\alpha_i - \alpha_i^*) K(x_i, x_k)$$

$f, b$  也只受  $x_i^T x$  影響  $\Rightarrow$  SVR algorithm can be kernelized.

## Problem 2.

$$\text{min } f(p, \mu, \xi) = p + \frac{1}{2} \sum_{i=1}^N C_i \xi_i, \quad C_i > 0, \quad 0 < 2 < C_i$$

constraint  $\begin{cases} g_{1,i}(p, \mu, \xi) = \|x_i - \mu\|^2 - p - \xi_i \leq 0 \\ g_{2,i}(p, \mu, \xi) = -\xi_i \leq 0 \end{cases} \quad \forall i = 1, 2, \dots, N.$

$$g_3(p) = -p \leq 0$$

$$\alpha \mapsto g_{1,i}, \quad \beta \mapsto g_{2,i}, \quad \gamma \mapsto g_3$$

$$\begin{aligned} L(p, \xi, \mu, \alpha, \beta, \gamma) &= p + \frac{1}{2} \sum_{i=1}^N C_i \xi_i + \sum_{i=1}^N \alpha_i (\|x_i - \mu\|^2 - p - \xi_i) + \sum_{i=1}^N \beta_i (-\xi_i) + \gamma (-p) \\ &= p \left(1 - \sum_{i=1}^N \alpha_i - \gamma\right) + \frac{1}{2} \sum_{i=1}^N C_i \xi_i - \sum_{i=1}^N \alpha_i \xi_i - \sum_{i=1}^N \beta_i \xi_i + \sum_{i=1}^N \alpha_i \|x_i - \mu\|^2 \\ &= p \left(1 - \sum_{i=1}^N \alpha_i - \gamma\right) + \sum_{i=1}^N \left(\frac{1}{2} C_i - \alpha_i - \beta_i\right) \xi_i + \sum_{i=1}^N \alpha_i \|x_i - \mu\|^2 \end{aligned}$$

2. : 確認する Strong duality: ① f, g are convex. ② Slater's condition

$$\textcircled{1} \quad f = p + \frac{1}{2} \sum_{i=1}^N C_i \xi_i \Rightarrow \text{linear} \Rightarrow \text{convex}.$$

$$g_{1,i} = \|x_i - \mu\|^2 - p - \xi_i : \|x_i - \mu\|^2, p, \xi_i \text{ convex} \Rightarrow g_{1,i} \text{ convex}.$$

$$g_{2,i} = -\xi_i : \text{convex}.$$

$$g_3 = -p : \text{convex}.$$

$\Rightarrow$  ① 成立.

② Slater's condition:  $\exists (\hat{p}, \hat{\mu}, \hat{\xi})$  s.t.  $g_{1,i}, g_{2,i}, g_3 < 0$ .

$$\text{Choose } \hat{\mu} = 0 \Rightarrow \|x_i\|^2 - \hat{p} - \hat{\xi}_i < 0 \quad \hat{\xi}_i > 0, \quad \hat{p} > 0$$

$$\text{Choose } \hat{\xi}_i = 1 \quad \forall i \Rightarrow \|x_i\|^2 - \hat{p} - 1 < 0, \quad 1 > 0, \quad \hat{p} > 0$$

$$\text{Choose } \hat{p} = \max_i \|x_i\|^2 \Rightarrow \|x_i\|^2 - \max_i \|x_i\|^2 - 1 < 0, \quad 1 > 0, \quad \max_i \|x_i\|^2 > 0$$

$$\Rightarrow (\max_i \|x_i\|^2, 0, 1) \text{ satisfies } g_{1,i}, g_{2,i}, g_3 < 0$$

$\Rightarrow$  Slater's condition 成立.

①, ② 成立  $\Rightarrow$  Strong duality 成立  $\Rightarrow$  duality gap = 0

$$3. L(\rho, \mu, \xi, \alpha, \beta, \gamma) = \rho(1 - \sum_{i=1}^N \alpha_i - \gamma) + \sum_{i=1}^N \left( \frac{1}{\nu} C_i - \alpha_i - \beta_i \right) \xi_i + \sum_{i=1}^N \alpha_i \|x_i - \mu\|^2$$

$$\frac{\partial L}{\partial \rho} = 1 - \sum_{i=1}^N \alpha_i - \gamma = 0 \Rightarrow \sum_{i=1}^N \alpha_i = 1 - \gamma$$

$$\nabla_\mu L = \sum_{i=1}^N -2(x_i - \mu)\alpha_i = 0 \Rightarrow \mu^* = \frac{\sum_{i=1}^N \alpha_i x_i}{\sum_{i=1}^N \alpha_i}$$

$$\frac{\partial L}{\partial \xi_i} = \frac{C_i}{\nu} - \alpha_i - \beta_i = 0 \Rightarrow \beta_i = \frac{C_i}{\nu} - \alpha_i$$

確保 Lower bound 存在  $\Leftrightarrow \frac{\partial L}{\partial \rho}, \frac{\partial L}{\partial \xi_i} = 0$ , 否則 Lower bound =  $-\infty$

$$\Rightarrow \Theta(\alpha, \beta, \gamma) = \begin{cases} \sum_{i=1}^N \alpha_i \|x_i\|^2 - \frac{\|\sum_{i=1}^N \alpha_i x_i\|^2}{\sum_{i=1}^N \alpha_i} & \text{if } \sum_{i=1}^N \alpha_i = 1 - \gamma, \beta_i = \frac{C_i}{\nu} - \alpha_i \\ -\infty & \text{otherwise.} \end{cases}$$

$$4. \|\alpha\|_1 = \sum_{i=1}^N \alpha_i, \quad \alpha_i = \|\alpha\|_1 \hat{\alpha}_i$$

$$\Theta(\alpha, \beta, \gamma) = \|\alpha\|_1 \sum_{i=1}^N \hat{\alpha}_i \|x_i\|^2 - \|\alpha\|_1 \|\sum_{i=1}^N \hat{\alpha}_i x_i\|^2$$

$$\|\sum_{i=1}^N \hat{\alpha}_i x_i\|^2 = \left( \sum_{i=1}^N \hat{\alpha}_i x_i \right)^T \left( \sum_{i=1}^N \hat{\alpha}_i x_i \right) = \sum_{i=1}^N \sum_{j=1}^N \hat{\alpha}_i \hat{\alpha}_j x_i^T x_j$$

$$\left\{ \begin{array}{l} \Theta(\alpha, \beta, \gamma) = \|\alpha\|_1 \left( \sum_{i=1}^N \hat{\alpha}_i \|x_i\|^2 - \sum_{i=1}^N \sum_{j=1}^N \hat{\alpha}_i \hat{\alpha}_j x_i^T x_j \right) \\ \text{Constraint: } \sum_{i=1}^N \alpha_i = 1 - \gamma \Rightarrow \sum_{i=1}^N \alpha_i \leq 1 \quad \text{since } \gamma \geq 0 \\ \frac{C_i}{\nu} - \alpha_i = \beta_i \geq 0 \Rightarrow 0 \leq \alpha_i \leq \frac{C_i}{\nu} \end{array} \right.$$

$$5. (a) \bar{\mu} = \mu^* \quad (\text{in 3.}) \Rightarrow \sum_{i=1}^N -2(x_i - \bar{\mu}) \alpha_i = 0 \Rightarrow \|\bar{x}\|_1 \bar{\mu} = \sum_{i=1}^N \bar{\alpha}_i x_i$$

$$(b) \xi_i \geq 0, \quad \xi_i \geq \|x_i - \bar{\mu}\|^2 - \rho \Rightarrow \xi_i \geq \max(\|x_i - \bar{\mu}\|^2 - \rho, 0)$$

$$f = \rho + \frac{1}{\nu} \sum_{i=1}^N C_i \xi_i = \rho + \frac{1}{\nu} \sum_{i=1}^N C_i \left( \max(\|x_i - \bar{\mu}\|^2 - \rho, 0) \right)$$

$$\text{Hence, } \bar{\rho} = \arg \min_{\rho \geq 0} \left( \rho + \frac{1}{\nu} \sum_{i=1}^N C_i \left( \max(\|x_i - \bar{\mu}\|^2 - \rho, 0) \right) \right)$$

$$(c) J(\rho) = \rho + \frac{1}{\nu} \sum_{i=1}^N C_i \left( \max(\|x_i - \bar{\mu}\|^2 - \rho, 0) \right), \quad \|x_i - \bar{\mu}\|^2 = d_i^2$$

$$\frac{\partial}{\partial \rho} \left( \max(\|x_i - \bar{\mu}\|^2 - \rho, 0) \right) = \begin{cases} -1 & d_i^2 > \rho \\ 0 & d_i^2 \leq \rho \\ [1, 0] & d_i^2 = \rho \end{cases}$$

$$\Rightarrow \partial J(\rho) = 1 - \frac{1}{\nu} \sum_{d_i^2 > \rho} C_i - \frac{1}{\nu} \sum_{d_i^2 \leq \rho} C_i [0, 1]$$

$$\bar{\rho} \text{ 是 min } \Rightarrow 0 \in \partial J(\rho) \Rightarrow 1 - \frac{1}{\nu} \sum_{d_i^2 \leq \bar{\rho}} C_i \geq 0, \quad 1 - \frac{1}{\nu} \sum_{d_i^2 > \bar{\rho}} C_i \leq 0$$

$$\Rightarrow \sum_{d_i^2 > \bar{\rho}} C_i \leq \nu$$

$$\sum_{d_i^2 \leq \bar{\rho}} C_i \geq \nu$$

Consider set  $A = \{p \geq 0 \mid \sum_{i: d_i^2 > p} C_i \leq \nu\}$

$\bar{p}$  satisfies  $\sum_{i: d_i^2 > \bar{p}} C_i \leq \nu \Rightarrow \bar{p} \in A \Rightarrow \bar{p} \geq \min(\{p \geq 0 \mid \sum_{i: \|x_i - \mu\|^2 > p} C_i \leq \nu\})$

Consider set  $B = \{p \geq 0 \mid \sum_{i: d_i^2 > p} C_i < \nu\}$

$$\forall \hat{p} \in B, \quad \partial J(\hat{p}) = 1 - \frac{1}{\nu} \sum_{i: d_i^2 > \hat{p}} C_i > 0$$

$\Rightarrow J$  is increasing if  $\hat{p} \in B \Rightarrow \bar{p} \leq \min(\{p \geq 0 \mid \sum_{i: \|x_i - \mu\|^2 > p} C_i < \nu\})$

$\Rightarrow \min(\{p \geq 0 \mid \sum_{i: \|x_i - \mu\|^2 > p} C_i \leq \nu\}) \leq \bar{p} \leq \min(\{p \geq 0 \mid \sum_{i: \|x_i - \mu\|^2 > p} C_i < \nu\})$

(d)  $\|x_i - \bar{\mu}\|^2 - \bar{p} > 0 \Rightarrow \xi_i \geq \|x_i - \mu\|^2 - \bar{p} > 0 \Rightarrow \xi_i \text{ 不滿足 } -\xi_i \leq 0 \text{ 的 boundary.}$

$$\Rightarrow \bar{\beta}_i = 0 \Rightarrow \alpha_i = \frac{C_i}{\nu} > 0 \Rightarrow \alpha_i \text{ 存在} \Rightarrow \xi_i = \|x_i - \bar{\mu}\|^2 - \bar{p}$$

(e)  $\|x_i - \bar{\mu}\|^2 - \bar{p} < 0, \quad \text{Suppose } \xi_i > 0 \Rightarrow \xi_i \text{ 不滿足 } -\xi_i \leq 0 \text{ 的 boundary.} \Rightarrow \beta_i = 0$

$$\Rightarrow \alpha_i = \frac{C_i}{\nu} > 0 \Rightarrow \alpha_i \text{ 存在} \Rightarrow \xi_i = \|x_i - \bar{\mu}\|^2 - \bar{p} < 0 \rightarrow *$$

Hence  $\bar{\xi}_i = 0$

$$\Rightarrow \bar{\xi}_i = \max(\|x_i - \bar{\mu}\|^2 - \bar{p}, 0)$$

(e) From (d), we have  $\alpha_i = \frac{C_i}{\nu}$  if  $\|x_i - \bar{\mu}\|^2 > \bar{p}$

if  $\|x_i - \bar{\mu}\|^2 < \bar{p} \Rightarrow \bar{\xi}_i = 0 > \|x_i - \bar{\mu}\|^2 - \bar{p} \Rightarrow \xi_i \text{ 不滿足 } g_{1,i} \text{ 的 boundary.}$

$$\Rightarrow \bar{\alpha}_i = 0 \quad \forall i$$

if  $\|x_i - \bar{\mu}\|^2 = \bar{p} \Rightarrow \xi_i = 0 \Rightarrow \text{同時在 } g_{1,i}, g_{2,i} \text{ 的 boundary.}$

$$\Rightarrow \bar{\beta}_i = \frac{C_i}{\nu} - \bar{\alpha}_i \geq 0 \Rightarrow 0 \leq \bar{\alpha}_i \leq \frac{C_i}{\nu}$$

6. From 5.(e). For  $\|x_i - \bar{\mu}\|^2 > \bar{p} \Rightarrow \bar{\alpha}_i = \frac{1}{n\nu}, \quad \sum_{i=1}^r \bar{\alpha}_i = 1$

$$\Rightarrow 1 = \sum_{i=1}^n \bar{\alpha}_i \geq \sum_{i: \|x_i - \bar{\mu}\|^2 > \bar{p}} \bar{\alpha}_i = n_{out} \cdot \frac{1}{n\nu} \Rightarrow \nu \geq \frac{n_{out}}{n}$$

For  $\|x_i - \bar{\mu}\|^2 = \bar{p} \Rightarrow 0 \leq \bar{\alpha}_i \leq \frac{1}{n\nu}, \quad \|x_i - \bar{\mu}\|^2 < \bar{p} \Rightarrow \bar{\alpha}_i = 0$

$$\Rightarrow \sum_{i=1}^n \bar{\alpha}_i = 1 = \sum_{i: \|x_i - \bar{\mu}\|^2 > \bar{p}} \bar{\alpha}_i + \sum_{i: \|x_i - \bar{\mu}\|^2 = \bar{p}} \bar{\alpha}_i \leq \frac{n_{out}}{n\nu} + \frac{n_{sv}}{n\nu}$$

$$\Rightarrow \nu \leq \frac{n_{out} + n_{sv}}{n} \Rightarrow \frac{n_{out}}{n} \leq \nu \leq \frac{n_{out} + n_{sv}}{n}$$

至少有  $\nu$  的 data 在 hypersphere 的 boundary 上或外面

至多有  $\nu$  的 data 在 hypersphere 的外部

$\Rightarrow$  這個 hypersphere 至少包含了  $1 - \nu$  的 data

### Problem 3.

$$(a) C_t = f_t \odot C_{t-1} + (1-f_t) \odot \tilde{C}_t$$

$$\frac{\partial C_t}{\partial f_t} = C_{t-1} - \tilde{C}_t$$

$$\frac{\partial C_t}{\partial \tilde{C}_t} = 1 - f_t$$

$$\frac{\partial C_t}{\partial C_{t-1}} = f_t$$

$$\frac{\partial L}{\partial f_t} = \frac{\partial L}{\partial C_t} \odot \frac{\partial C_t}{\partial f_t} = \frac{\partial L}{\partial C_t} \odot (C_{t-1} - \tilde{C}_t)$$

$$\frac{\partial L}{\partial \tilde{C}_t} = \frac{\partial L}{\partial C_t} \odot \frac{\partial C_t}{\partial \tilde{C}_t} = \frac{\partial L}{\partial C_t} \odot (1 - f_t)$$

$$\frac{\partial L}{\partial C_{t-1}} = \frac{\partial L}{\partial C_t} \odot \frac{\partial C_t}{\partial C_{t-1}} = \frac{\partial L}{\partial C_t} \odot f_t$$

$$(b) \quad \begin{array}{c} \textcircled{1} \\ C_t \end{array} \longrightarrow h_t \longrightarrow L \quad \begin{array}{c} \textcircled{2} \\ C_t \end{array} \longrightarrow C_{t+1} \longrightarrow L$$

$$\textcircled{1} \quad h_t = o_t \odot \tanh(c_t)$$

$$\textcircled{2} \quad \frac{\partial C_{t+1}}{\partial C_t} = f_{t+1}$$

$$\frac{\partial h_t}{\partial C_t} = o_t \odot (1 - \tanh^2(c_t))$$

$$\Rightarrow \frac{\partial L}{\partial C_t} = \frac{\partial L}{\partial h_t} \odot \frac{\partial h_t}{\partial C_t} + \frac{\partial L}{\partial C_{t+1}} \odot \frac{\partial C_{t+1}}{\partial C_t}$$

$$= \delta h_t \odot o_t \odot (1 - \tanh^2(c_t)) + \delta C_{t+1} \odot f_{t+1}$$

$$\frac{\partial L}{\partial f_t} = \delta C_t \odot (C_{t-1} - \tilde{C}_t) \quad \text{From (a)}$$

$$z_{f,t} = W_f x_t + U_f h_{t-1} + b_f, \quad f_t = \sigma(z_{f,t}), \quad \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

$$\frac{\partial L}{\partial z_{f,t}} = \frac{\partial L}{\partial f_t} \odot \frac{\partial f_t}{\partial z_{f,t}} = \delta f_t \odot f_t (1 - f_t)$$

$$z_{o,t} = W_o x_t + U_o h_{t-1} + b_o, \quad o_t = \sigma(z_{o,t})$$

$$\frac{\partial L}{\partial z_{o,t}} = \frac{\partial L}{\partial o_t} \odot \frac{\partial o_t}{\partial z_{o,t}} = \delta o_t \odot o_t (1 - o_t)$$

$$z_{c,t} = W_c x_t + U_c h_{t-1} + b_c, \quad \tilde{C}_t = \tanh(z_{c,t})$$

$$\frac{\partial L}{\partial z_{c,t}} = \frac{\partial L}{\partial \tilde{C}_t} \odot \frac{\partial \tilde{C}_t}{\partial z_{c,t}} = \delta \tilde{C}_t \odot (1 - \tilde{C}_t^2)$$

$$(c) \quad \frac{\partial z_{f,t}}{\partial w_f} = x_t^T, \quad \frac{\partial L}{\partial w_f} = \sum_t \left( \frac{\partial L}{\partial z_{f,t}} \cdot \frac{\partial z_{f,t}}{\partial w_f} \right) = \sum_t \delta z_{f,t} x_t^T$$

$$\frac{\partial z_{f,t}}{\partial u_f} = h_{t-1}^T, \quad \frac{\partial L}{\partial u_f} = \sum_t \left( \frac{\partial L}{\partial z_{f,t}} \cdot \frac{\partial z_{f,t}}{\partial u_f} \right) = \sum_t \delta z_{f,t} h_{t-1}^T$$

$$\frac{\partial z_{f,t}}{\partial b_f} = 1, \quad \frac{\partial L}{\partial b_f} = \sum_t \left( \frac{\partial L}{\partial z_{f,t}} \cdot \frac{\partial z_{f,t}}{\partial b_f} \right) = \sum_t \delta z_{f,t}$$

## Problem 4.

$$(a) \quad \alpha_t = 1 - \beta_t, \quad \bar{\alpha}_t = \prod_{s=1}^t \alpha_s = \prod_{s=1}^t (1 - \beta_s)$$

$$g(x_t | x_{t-1}) = N(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

$$x_t = \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_{t-1}$$

$$= \sqrt{\alpha_t} x_{t-1} + \sqrt{1 - \alpha_t} \varepsilon_{t-1}$$

$$x_{t-1} = \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-1}} \varepsilon_{t-2}$$

$$\Rightarrow x_t = \sqrt{\alpha_t} (\sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-1}} \varepsilon_{t-2}) + \sqrt{1 - \alpha_t} \varepsilon_{t-1}$$

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{\alpha_t (1 - \alpha_{t-1})} \varepsilon_{t-2} + \sqrt{1 - \alpha_t} \varepsilon_{t-1}$$

$\varepsilon_{t-1}, \varepsilon_{t-2}$  : Gaussian Noise,  $\varepsilon_{t-2}, \varepsilon_{t-1} \in N(0, 1)$

$$\sqrt{\alpha_t (1 - \alpha_{t-1})} \varepsilon_{t-2} \in N(0, \alpha_t - \alpha_{t-1} \alpha_{t-1}) \quad . \quad \sqrt{1 - \alpha_t} \varepsilon_{t-1} \in N(0, 1 - \alpha_t)$$

$$\Rightarrow \sqrt{\alpha_t (1 - \alpha_{t-1})} \varepsilon_{t-2} + \sqrt{1 - \alpha_t} \varepsilon_{t-1} \in N(0, 1 - \alpha_t \alpha_{t-1})$$

$$\Rightarrow \text{Let } \sqrt{1 - \alpha_t \alpha_{t-1}} \varepsilon_{t-2} = \sqrt{\alpha_t (1 - \alpha_{t-1})} \varepsilon_{t-2} + \sqrt{1 - \alpha_t} \varepsilon_{t-1}$$

$$\text{Then } x_t = \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\varepsilon}_{t-2}$$

根據遞迴，我們可以把  $x_t$  看作  $x_0$ ,  $\varepsilon$  表示

$$x_t = \sqrt{\alpha_t \alpha_{t-1} \dots \alpha_1} x_0 + \sqrt{1 - \alpha_t \alpha_{t-1} \dots \alpha_1} \varepsilon, \quad \varepsilon \in N(0, 1)$$

$$= \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \varepsilon$$

$$\mathbb{E}[x_t | x_0] = \sqrt{\alpha_t} x_0, \quad \text{Var}(x_t | x_0) = (1 - \alpha_t) I$$

$$g(x_t | x_0) = N(x_t; \sqrt{\alpha_t} x_0, (1 - \alpha_t) I)$$

$$(b) \quad g(x_{t-1} | x_t, x_0) = \frac{g(x_t | x_{t-1}, x_0) g(x_{t-1} | x_0)}{g(x_t | x_0)}$$

$$\text{Markov property} : \quad g(x_t | x_{t-1}, x_0) = g(x_t | x_{t-1})$$

$$\Rightarrow g(x_{t-1} | x_t, x_0) \propto g(x_t | x_{t-1}) g(x_{t-1} | x_0)$$

$$g(x_t | x_{t-1}) = N(x_t; \sqrt{\alpha_t} x_{t-1}, \beta_t I)$$

$$g(x_{t-1} | x_0) = N(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}} x_0, (1 - \bar{\alpha}_{t-1}) I)$$

考慮 exponential term (ignore  $-\frac{1}{2}$ )

$$\begin{aligned}
 L &= \frac{(\bar{x}_t - \sqrt{\alpha_t} \bar{x}_{t-1})^2}{\beta_t} + \frac{(\bar{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}} x_0)^2}{1 - \bar{\alpha}_{t-1}} \\
 &= \frac{\bar{x}_t^2 - 2\sqrt{\alpha_t} \bar{x}_t \bar{x}_{t-1} + \alpha_t \bar{x}_{t-1}^2}{\beta_t} + \frac{\bar{x}_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}} x_0 \bar{x}_{t-1} + \bar{\alpha}_{t-1} x_0^2}{1 - \bar{\alpha}_{t-1}} \\
 &= \bar{x}_{t-1}^2 \left( \frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) - 2 \left( \frac{\sqrt{\alpha_t}}{\beta_t} \bar{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right) \bar{x}_{t-1} + \frac{1}{\beta_t} \bar{x}_t^2 + \frac{\bar{\alpha}_{t-1}}{1 - \bar{\alpha}_{t-1}} x_0^2
 \end{aligned}$$

$$q_f(x_{t-1} | x_t, x_0) = N(x_{t-1}; \tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t I) \Rightarrow \text{exponential term} = \frac{(x_{t-1} - \mu)^2}{\tilde{\beta}_t}$$

$$\Rightarrow \frac{1}{\tilde{\beta}_t} = \frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} = \frac{\alpha_t(1 - \bar{\alpha}_{t-1}) + \beta_t}{\beta_t(1 - \bar{\alpha}_{t-1})} = \frac{\alpha_t - \alpha_t \bar{\alpha}_{t-1} + \beta_t}{\beta_t(1 - \bar{\alpha}_{t-1})} = \frac{1 - \bar{\alpha}_t}{(1 - \bar{\alpha}_{t-1})\beta_t}$$

$$\Rightarrow \tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

$$\frac{2\mu_t}{\tilde{\beta}_t} = 2 \left( \frac{\sqrt{\alpha_t}}{\beta_t} \bar{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right)$$

$$\mu_t = \left( \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \right) \left( \frac{\sqrt{\alpha_t}}{\beta_t} \bar{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right)$$

$$= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \bar{x}_t + \frac{\beta_t \sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} x_0$$

$$= \frac{\sqrt{1 - \beta_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \bar{x}_t + \frac{\beta_t \sqrt{\alpha_t}}{1 - \bar{\alpha}_t} x_0$$

$$(C) \quad \hat{x}_0 = \frac{1}{\sqrt{\alpha_t}} \left( \bar{x}_t - \sqrt{1 - \bar{\alpha}_t} \varepsilon_0(x_t, t) \right)$$

$$\begin{aligned}
 \mu_\theta(x_t, \hat{x}_0) &= \frac{\sqrt{1 - \beta_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \bar{x}_t + \frac{\beta_t \sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \cdot \frac{1}{\sqrt{\alpha_t}} \left( \bar{x}_t - \sqrt{1 - \bar{\alpha}_t} \varepsilon_0 \right) \\
 &= \left( \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} + \frac{\beta_t}{(1 - \bar{\alpha}_t)\sqrt{\alpha_t}} \right) \bar{x}_t - \frac{\beta_t}{(1 - \bar{\alpha}_t)\sqrt{\alpha_t}} \sqrt{1 - \bar{\alpha}_t} \varepsilon_0 \\
 &= \frac{\alpha_t(1 - \bar{\alpha}_{t-1}) + \beta_t}{(1 - \bar{\alpha}_t)\sqrt{\alpha_t}} \bar{x}_t - \frac{\beta_t}{(1 - \bar{\alpha}_t)\sqrt{\alpha_t}} \sqrt{1 - \bar{\alpha}_{t-1}} \varepsilon_0 \\
 &= \frac{\alpha_t + \beta_t - \bar{\alpha}_t}{(1 - \bar{\alpha}_t)\sqrt{\alpha_t}} \bar{x}_t - \frac{\beta_t}{(1 - \bar{\alpha}_t)\sqrt{\alpha_t}} \sqrt{1 - \bar{\alpha}_t} \varepsilon_0 \\
 &= \frac{1}{\sqrt{\alpha_t}} \bar{x}_t - \frac{\beta_t}{\sqrt{\alpha_t} \sqrt{1 - \bar{\alpha}_t}} \varepsilon_0 \\
 &= \frac{1}{\sqrt{1 - \beta_t}} \left( \bar{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_0 \right)
 \end{aligned}$$

$$(d) \text{ Show: } D_{KL} \left( q(x_{t+1} | x_t, x_0) \| p_\theta(x_{t+1} | x_t) \right) = \frac{1}{2\beta_t} \| \tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t) \|^2 + \text{const}$$

$$q(x_{t+1} | x_t, x_0) \sim N(\tilde{\mu}_t, \tilde{\beta}_t I)$$

$$p_\theta(x_{t+1} | x_t) \sim N(\mu_\theta, \tilde{\beta}_t I)$$

$$D_{KL} \left( N(\mu_1, \Sigma) \| N(\mu_2, \Theta) \right) = \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2)$$

$$\begin{aligned} D_{KL} \left( q(x_{t+1} | x_t, x_0) \| p_\theta(x_{t+1} | x_t) \right) &= \frac{1}{2\beta_t} (\mu_t - \mu_\theta)^T (\mu_t - \mu_\theta) \\ &= \frac{\| \mu_t - \mu_\theta \|^2}{2\tilde{\beta}_t} \end{aligned}$$

$$\mu_t = \frac{\sqrt{1-\alpha_t} (1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_t + \frac{\beta_t \sqrt{\bar{\alpha}_{t-1}}}{1-\bar{\alpha}_t} x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \varepsilon_0 \right)$$

$$\mu_\theta = \frac{1}{\sqrt{1-\beta_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \varepsilon_\theta \right)$$

$$\mu_t - \mu_\theta = - \frac{\beta_t}{\sqrt{\alpha_t(1-\bar{\alpha}_t)}} (\varepsilon - \varepsilon_\theta)$$

$$\| \mu_t - \mu_\theta \|^2 = \frac{\beta_t^2}{\alpha_t(1-\bar{\alpha}_t)} \| \varepsilon - \varepsilon_\theta \|^2$$

$$D_{KL}(q \| p_\theta) = \frac{\beta_t^2}{2\tilde{\beta}_t \alpha_t(1-\bar{\alpha}_t)} \| \varepsilon - \varepsilon_\theta \|^2$$

$$\Rightarrow -ELBO(\theta) = \text{const} + \sum_{t=2}^T E \left[ \frac{\beta_t^2}{2\tilde{\beta}_t \alpha_t(1-\bar{\alpha}_t)} \| \varepsilon - \varepsilon_\theta(x_t, t) \|^2 \right] + E_g[-\log p_\theta(x_0 | x_1)]$$