

# MLHW4

Q1. Given LSTM model and

$z_i, z_f, z_o \rightarrow \text{gate control input}$   
 $c \rightarrow \text{cell memory}$   
 $f, g, h \rightarrow \text{activation function } (f: \text{sigmoid} / g-h \text{ are linear})$   
 $w, w_i, w_f, w_o \rightarrow \text{weight} / b, b_i, b_f, b_o \rightarrow \text{bias}$   
 updated cell memory  $\rightarrow c' = f(z_i)g(z) + cf(z_f)$

Input Sequence  $\begin{cases} x^1 = [0, 1, 0, 3] \\ x^2 = [1, 0, 1, -2] \\ x^3 = [1, 1, 1, 4] \\ x^4 = [0, 1, 1, 0] \end{cases}$

Weight  $\begin{cases} w = [0, 0, 0, 1] \\ w_i = [100, 100, 0, 0] \\ w_f = [-100, -100, 0, 0] \\ w_o = [0, 0, 100, 0] \end{cases}$

Bias  $\begin{cases} b = 0 \\ b_i = -10 \\ b_f = 110 \\ b_o = -10 \end{cases}$

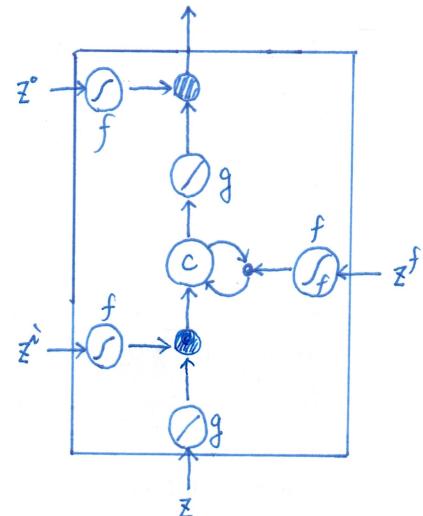
- $x^1 = [0, 1, 0, 3], c = 0$

$$z = w x^1 + b = [0, 0, 0, 1] \begin{bmatrix} 0 \\ 1 \\ 0 \\ 3 \end{bmatrix} + 0 = [3] \Rightarrow g(z) = 3$$

$$z^i = w^i x^1 + b^i = [100, 100, 0, 0] \begin{bmatrix} 0 \\ 1 \\ 0 \\ 3 \end{bmatrix} + (-10) = [90] \Rightarrow f(z^i) = 1$$

$$z^f = w^f x^1 + b^f = [-100, -100, 0, 0] \begin{bmatrix} 0 \\ 1 \\ 0 \\ 3 \end{bmatrix} + 110 = [10] \Rightarrow f(z^f) = 1$$

$$z^o = w^o x^1 + b^o = [0, 0, 100, 0] \begin{bmatrix} 0 \\ 1 \\ 0 \\ 3 \end{bmatrix} - 10 = [-10] \Rightarrow f(z^o) = 0$$



$$c' = f(z^i)g(z) + cf(z^f) = 1 \cdot 3 + 0 \cdot 1 = [3] \Rightarrow h(c') = 3$$

$$y = f(z^o)h(c') = 0$$

- $x^2 = [1, 0, 1, -2], c = 3$

$$z = [-2] \Rightarrow g(z) = -2$$

$$z^i = [90] \Rightarrow f(z^i) = 1$$

$$z^f = [10] \Rightarrow f(z^f) = 1$$

$$z^o = [90] \Rightarrow f(z^o) = 1$$

- $x^3 = [1, 1, 1, 4], c = 1$

$$z = [4] \Rightarrow g(z) = 4$$

$$z^i = [90] \Rightarrow f(z^i) = 1$$

$$z^f = [90] \Rightarrow f(z^f) = 0$$

$$z^o = [90] \Rightarrow f(z^o) = 1$$

- $x^4 = [0, 1, 1, 0], c = 4$

$$z = [0] \Rightarrow g(z) = 0$$

$$z^i = [90] \Rightarrow f(z^i) = 1$$

$$z^f = [10] \Rightarrow f(z^f) = 1$$

$$z^o = [90] \Rightarrow f(z^o) = 1$$

$$c' = 1 \times -2 + 3 \cdot 1 = [1]$$

$$c' = [4], h(c') = 4, y = 4$$

$$c' = [4], h(c') = 4, y = 4$$

$$h(c') = 1$$

$$y = 1$$

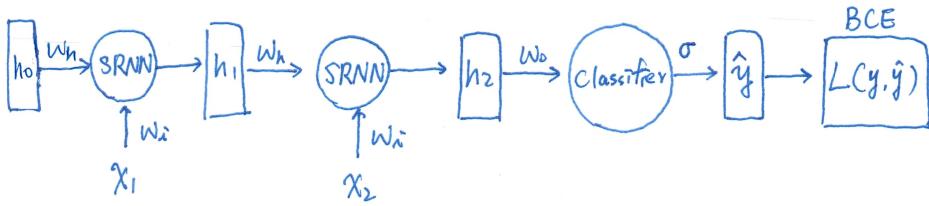
t	$z^i$	$f(z^i)$	$z^f$	$f(z^f)$	$z^o$	$f(z^o)$	$z$	$c'$	y
1	90	1	10	1	-10	0	3	3	0
2	90	1	10	1	90	1	-2	1	1
3	190	1	90	0	90	1	4	4	4
4	90	1	10	1	90	1	0	4	4

\*

Q2. Given Objective function is Binary Cross Entropy  $L(y, \hat{y})$  where  $y$  is label

$$\hat{y} = \sigma(w_0 h_2) = \frac{1}{1 + \exp(-w_0 h_2)}$$

$$h_t = \tanh(w_i x_t + w_h h_{t-1})$$



Goal: derive  $\frac{\partial L(y, \hat{y})}{\partial w_0}$ ,  $\frac{\partial L(y, \hat{y})}{\partial w_h}$ ,  $\frac{\partial L(y, \hat{y})}{\partial w_i}$

$$L(y, \hat{y}), \hat{y}(w_0, h_2), h_2(w_i, x_2, w_h, h_1), h_1(w_i, x_1, w_h, h_0), h_0=0$$

$$(1) \frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_0}$$

$$(2) \begin{aligned} \frac{\partial L}{\partial w_h} &= \frac{\partial L}{\partial h_2} \frac{\partial h_2}{\partial w_h} + \frac{\partial L}{\partial h_1} \frac{\partial h_1}{\partial w_h} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_2} \frac{\partial h_2}{\partial w_h} + \frac{\partial L}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_h} \\ &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_2} \frac{\partial h_2}{\partial w_h} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_h} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_2} \left[ \frac{\partial h_2}{\partial w_h} + \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_h} \right] \end{aligned}$$

$$(3) \begin{aligned} \frac{\partial L}{\partial w_i} &= \frac{\partial L}{\partial h_2} \frac{\partial h_2}{\partial w_i} + \frac{\partial L}{\partial h_1} \frac{\partial h_1}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_2} \frac{\partial h_2}{\partial w_i} + \frac{\partial L}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_i} \\ &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_2} \frac{\partial h_2}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_2} \left[ \frac{\partial h_2}{\partial w_i} + \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_i} \right] \end{aligned}$$

$$(1) \frac{\partial \hat{y}}{\partial w_0} = \frac{h_2 \exp(-w_0 h_2)}{\left[ 1 + \exp(-w_0 h_2) \right]^2}$$

$\because \begin{cases} h_2 = \tanh(w_i x_2 + w_h h_1) \\ h_1 = \tanh(w_i x_1 + w_h h_0) \end{cases}$

$$(2) \frac{\partial \hat{y}}{\partial w_h} = \frac{w_0 \exp(-w_0 h_2)}{\left[ 1 + \exp(-w_0 h_2) \right]^2}, \quad \frac{\partial h_2}{\partial w_h} = h_1 \operatorname{sech}^2(w_i x_2 + w_h h_1)$$

$$\frac{\partial h_2}{\partial h_1} = w_h \operatorname{sech}^2(w_i x_2 + w_h h_1) \Rightarrow \frac{\partial h_1}{\partial w_h} = h_0 \operatorname{sech}^2(w_i x_1 + w_h h_0) = 0$$

$$(3) \frac{\partial h_2}{\partial w_i} = x_2 \operatorname{sech}^2(w_i x_2 + w_h h_1) \Rightarrow \frac{\partial h_1}{\partial w_i} = x_1 \operatorname{sech}^2(w_i x_1 + w_h h_0)$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{\partial -[y \log \hat{y} + (1-y) \log(1-\hat{y})]}{\partial \hat{y}} = -\left[ \frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}} \right] = \frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}}$$

$$(1) \quad \frac{\partial L}{\partial w_0} = \left( \frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}} \right) \frac{h_2 \exp(-w_0 h_2)}{\left[ 1 + \exp(-w_0 h_2) \right]^2}$$

$$(2) \quad \frac{\partial L}{\partial w_h} = \left( \frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}} \right) \frac{w_0 \exp(-w_0 h_2)}{\left[ 1 + \exp(-w_0 h_2) \right]^2} \left[ h_1 \operatorname{sech}^2(w_0 x_2 + w_h h_1) + o \right]$$

$$(3) \quad \frac{\partial L}{\partial w_x} = \left( \frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}} \right) \frac{w_0 \exp(-w_0 h_2)}{\left[ 1 + \exp(-w_0 h_2) \right]^2} \left[ x_2 \operatorname{sech}^2(w_0 x_2 + w_h h_1) + w_h \operatorname{sech}^2(w_0 x_2 + w_h h_1) \cdot x_1 \operatorname{sech}^2(w_0 x_1 + w_h h_0) \right]$$

Q.E.D.

Q3. Given input space  $X$ , # of classes  $K$ , training data set  $\{(x_i, \hat{y}_i)\}_{i=1}^m$   
 $\hookrightarrow \in X$   
collection of multi-class classifiers  $\mathcal{F}$

We want to find function  $g_{T+1}^k(x) = \sum_{t=1}^T \alpha_t f_t^k(x)$ ,  $k \in [1, K]$  where  $f \in \mathcal{F}$ ,  $\alpha_t \in \mathbb{R}$   
 $t \in [1, T]$

Objective: to minimize  $L(g_{T+1}^1, \dots, g_{T+1}^K) = \sum_{i=1}^m \exp\left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{T+1}^k(x_i) - g_{T+1}^{\hat{y}_i}(x_i)\right)$

and show how use  $f_t$  and  $\alpha_t$  to achieve this goal.

<sol> In Binary classifier, that is  $y = \{\pm 1\}$ , we aim to minimize  $\hat{R}_S(g)$  just like this problem.

\* determine direction  $f_t = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \frac{\partial}{\partial \alpha} \hat{R}_S(g_t + \alpha f) \Big|_{\alpha=0}$

\* determine step-size  $\alpha_t = \underset{\alpha \in \mathbb{R}}{\operatorname{argmin}} \hat{R}_S(g_t + f_t)$

$f_t = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \frac{\partial}{\partial \alpha} \sum_{i=1}^m \exp\left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{T+1}^k(x_i) - g_{T+1}^{\hat{y}_i}(x_i)\right) \Big|_{\alpha=0}$  take  $g_t + \alpha f$  in Loss function

$\Rightarrow f_t = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \frac{\partial}{\partial \alpha} \sum_{i=1}^m \exp\left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} (g_t^k(x_i) + \alpha f^k(x_i)) - (g_t^{\hat{y}_i}(x_i) + \alpha f^{\hat{y}_i}(x_i))\right) \Big|_{\alpha=0}$

$$= \underset{f \in \mathcal{H}}{\operatorname{argmin}} \frac{\partial}{\partial \alpha} \sum_{i=1}^m \exp\left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} (g_t^k(x_i) + \alpha f^k(x_i))\right] \exp(-g_t^{\hat{y}_i}(x_i) - \alpha f^{\hat{y}_i}(x_i)) \Big|_{\alpha=0}$$

$$= \underset{f \in \mathcal{H}}{\operatorname{argmin}} \sum_{i=1}^m \exp\left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} (g_t^k(x_i) + \alpha f^k(x_i))\right] \cdot \left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} f^k(x_i)\right] \exp(-g_t^{\hat{y}_i}(x_i) - \alpha f^{\hat{y}_i}(x_i)) \\ + \exp\left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} (g_t^k(x_i) + \alpha f^k(x_i))\right] \cdot \exp(-g_t^{\hat{y}_i}(x_i) - \alpha f^{\hat{y}_i}(x_i)) \cdot (-f^{\hat{y}_i}(x_i)) \Big|_{\alpha=0}$$

$$= \underset{f \in \mathcal{H}}{\operatorname{argmin}} \sum_{i=1}^m \exp\left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} (g_t^k(x_i) + \alpha f^k(x_i))\right] \exp(-g_t^{\hat{y}_i}(x_i) - \alpha f^{\hat{y}_i}(x_i)) \cdot \left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} f^k(x_i) - f^{\hat{y}_i}(x_i)\right] \Big|_{\alpha=0}$$

$$= \underset{f \in \mathcal{H}}{\operatorname{argmin}} \sum_{i=1}^m \exp\left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} (g_t^k(x_i) + \alpha f^k(x_i)) - (g_t^{\hat{y}_i}(x_i) + \alpha f^{\hat{y}_i}(x_i))\right] \left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} f^k(x_i) - f^{\hat{y}_i}(x_i)\right] \Big|_{\alpha=0}$$

$$= \underset{f \in \mathcal{H}}{\operatorname{argmin}} \sum_{i=1}^m \exp\left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_t^k(x_i) - g_t^{\hat{y}_i}(x_i)\right] \left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} f^k(x_i) - f^{\hat{y}_i}(x_i)\right]$$

$$\text{Let } Z_t = \sum_{i=1}^m \exp\left(\frac{1}{K-1} \sum_{k \neq \hat{y}_t} g_t^k(x_i) - g_t^{\hat{y}_t}(x_i)\right), D_t(i) = \frac{\exp\left[\frac{1}{K-1} \sum_{k \neq \hat{y}_t} g_t^k(x_i) - g_t^{\hat{y}_t}(x_i)\right]}{Z_t}$$

$$\Rightarrow \underset{f \in H}{\operatorname{argmin}} \sum_{i=1}^m Z_t D_t(i) \left[ \frac{1}{K-1} \sum_{k \neq \hat{y}_t} f_t^k(x_i) - f_t^{\hat{y}_t}(x_i) \right]$$

$$= \underset{f \in H}{\operatorname{argmin}} \underbrace{Z_t}_{\text{fixed}} \mathbb{E} \left[ \frac{1}{K-1} \sum_{k \neq \hat{y}_t} f_t^k(x_i) - f_t^{\hat{y}_t}(x_i) \right] \quad \begin{array}{l} \text{this expectation just has 2 values including} \\ \text{correct} \\ \text{wrong} \end{array}$$

$$\text{相當於} \Rightarrow \underset{f \in H}{\operatorname{argmin}} \underset{k \neq \hat{y}_t}{\mathbb{P}} \left[ \frac{1}{K-1} \sum_{k \neq \hat{y}_t} f_t^k(x_i) - f_t^{\hat{y}_t}(x_i) = \text{判錯} \right] \quad \text{Q.E.D.}$$

$$\Rightarrow \mathbb{P} \left[ \underset{i \sim D_t}{\frac{1}{K-1} \sum_{k \neq \hat{y}_t} f_t^k(x_i) - f_t^{\hat{y}_t}(x_i)} = \text{判對} \right] \times \begin{array}{l} \text{value of} \\ \text{correct} \end{array} + \mathbb{P} \left[ \underset{i \sim D_t}{\frac{1}{K-1} \sum_{k \neq \hat{y}_t} f_t^k(x_i) - f_t^{\hat{y}_t}(x_i)} = \text{判錯} \right] \times \begin{array}{l} \text{value of} \\ \text{wrong} \end{array}$$

$$\Rightarrow 1 - \text{判對的機率} = \text{判錯的機率.}$$

$$\begin{aligned} & \left[ 1 - \mathbb{P} \left[ \underset{i \sim D_t}{\frac{1}{K-1} \sum_{k \neq \hat{y}_t} f_t^k(x_i) - f_t^{\hat{y}_t}(x_i)} = \text{判錯} \right] \right] \times \begin{array}{l} \text{value of} \\ \text{correct} \end{array} \\ & + \mathbb{P} \left[ \underset{i \sim D_t}{\frac{1}{K-1} \sum_{k \neq \hat{y}_t} f_t^k(x_i) - f_t^{\hat{y}_t}(x_i)} = \text{判錯} \right] \times \begin{array}{l} \text{value of} \\ \text{wrong} \end{array} \\ & = \underbrace{\left[ \begin{array}{l} \text{value of} \\ \text{correct} \end{array} \right] + \left[ \begin{array}{l} \text{value of} \\ \text{wrong} \end{array} \right]}_{\text{fixed value}} \times \mathbb{P} \left[ \underset{i \sim D_t}{\frac{1}{K-1} \sum_{k \neq \hat{y}_t} f_t^k(x_i) - f_t^{\hat{y}_t}(x_i)} = \text{判錯} \right] \end{aligned}$$

+ value of  
correct  
fixed value

$$\text{In Binary classifier, } d_t = \underset{\alpha \in \mathbb{R}}{\operatorname{argmin}} \hat{R}_S(g_t + \alpha f_t)$$

$$\Rightarrow d_t = \underset{\alpha \in \mathbb{R}}{\operatorname{argmin}} \sum_{i=1}^m \exp \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_t} [g_t^k(x_i) - \alpha f_t^k(x_i)] - (g_t^{\hat{y}_t}(x_i) + \alpha f_t^{\hat{y}_t}(x_i)) \right)$$

$$= \sum_{i: f_t^{\hat{y}_t}(x_i) \text{判錯}} \exp \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_t} [g_t^k(x_i) - \alpha f_t^k(x_i)] - (g_t^{\hat{y}_t}(x_i) + \alpha f_t^{\hat{y}_t}(x_i)) \right) +$$

$$\sum_{i: f_t^{\hat{y}_t}(x_i) \text{判對}} \exp \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_t} [g_t^k(x_i) - \alpha f_t^k(x_i)] - (g_t^{\hat{y}_t}(x_i) + \alpha f_t^{\hat{y}_t}(x_i)) \right), \quad \begin{array}{l} \text{if value of correct is } C \\ \text{and value of wrong is } W \end{array}$$

$$= \sum_{i: f_t^{\hat{y}_t}(x_i) = W} Z_t D_t(x_i) \exp(-g_t^{\hat{y}_t}(x_i) + \alpha W) + \sum_{i: f_t^{\hat{y}_t}(x_i) = C} Z_t D_t(x_i) \exp(-g_t^{\hat{y}_t}(x_i) + \alpha C)$$

$$\text{Let } \varepsilon_t = \underset{i \sim D_t}{\mathbb{P}} [f_t^{\hat{y}_t}(x_i) = \text{判錯}] = Z_t \varepsilon_t \exp(-g_t^{\hat{y}_t}(x_i) + \alpha W) + Z_t (1 - \varepsilon_t) \exp(-g_t^{\hat{y}_t}(x_i) + \alpha C)$$

$$\text{Let } \frac{\partial}{\partial \alpha} [\varepsilon_t \exp(-g_t^{\hat{y}_t}(x_i) + \alpha W) + (1 - \varepsilon_t) \exp(-g_t^{\hat{y}_t}(x_i) + \alpha C)]$$

$$\begin{aligned} & = W \varepsilon_t \exp(-g_t^{\hat{y}_t}(x_i) + \alpha W) + C \exp(-g_t^{\hat{y}_t}(x_i) + \alpha C) - C \varepsilon_t \exp(-g_t^{\hat{y}_t}(x_i) + \alpha C) \\ & = G \left[ W \exp(\alpha W) + \frac{C}{\varepsilon_t} \exp(\alpha C) - C \exp(\alpha C) \right] \quad \text{Q.E.D.} \end{aligned}$$

## MLHW Q4.

- Let  $X$  be observed variable  $\Rightarrow$  外在特徵
- $Z$  be latent variable ( $X$  分別屬於哪一群)  $\Rightarrow$  內在性格 (不易觀察)
- Want to do maximum likelihood estimation  $\theta^* = \underset{\theta \in \Theta}{\operatorname{argmax}} p(X; \theta)$  where  $\{p(\cdot; \theta) : \theta \in \Theta\}$  is a collection of probability models
- EM algorithms

Initial  $\theta^{(0)}$ , then iterate 1, 2, 3 ... as follow:

▷ E-Step: Derive

$$Q(\theta | \theta^{(t)}) = \sum_z \underbrace{P(z|X; \theta^{(t)})}_{\text{根據當前的 probability}} \log P(X, z; \theta)$$

▷ M-Step: Find  
 $\theta^{(t+1)} = \underset{\theta \in \Theta}{\operatorname{argmax}} Q(\theta | \theta^{(t)})$

model 和 外在特徵  $\Rightarrow$   
猜測 內在性格

Example:

If  $X = \{x_1, x_2, \dots, x_N\}$ ,  $Z = \{z_1, z_2, \dots, z_N\}$ ,  $P_\theta(X, Z; \theta) = \prod_{i=1}^N P(x_i, z_i; \theta) \quad \forall \theta \in \Theta$

$\uparrow$   $N$  training data       $\uparrow$  corresponding latent variable       $\uparrow$  data are i.i.d generated

Then,  $Q(\theta | \theta^{(t)}) = \mathbb{E}_{\substack{Z \sim P(\cdot | X; \theta^{(t)})}} [\log P(X, Z; \theta)]$

$$= \sum_{i=1}^N \mathbb{E}_{\substack{z_i \sim P(\cdot | x_i; \theta^{(t)})}} [\log P(x_i, z_i; \theta)] \quad \text{by independence}$$

Denote the parameter estimate at  $t$  as:

$$\theta^{(t)} = \{(\pi_k^{(t)}, \mu_k^{(t)}, \Sigma_k^{(t)})\}_{k=1}^K = \sum_{i=1}^N \mathbb{E}_{\substack{z_i \sim P(\cdot | x_i; \theta^{(t)})}} [\log P(x_i, z_i; \theta)] \quad (\text{因為是獨立的，所以推測只要知道 } x_i \text{ 就好})$$

要算的東西

Apply to GMM

→ data 是從  $k$ -th gaussian distribution 產生且外在特徵是  $X$

In GMM with  $K$  gaussian distributions,  $P(X, z=k; \theta) = \pi_k N(X; \mu_k, \Sigma_k)$

where  $\theta = \{(\pi_k, \mu_k, \Sigma_k)\}_{k=1}^K$

▷  $P(z_i=k | x_i; \theta^{(t)}) = \frac{P(z_i=k, x_i; \theta^{(t)})}{\sum_{j=1}^K P(z_i=j, x_i; \theta^{(t)})}$

貝氏定理. → 從  $k$ -th gaussian dist. 發生的機率  
(sampling from  $k$ -th G.D 的 P.)

$$= \frac{\pi_k^{(t)} N(x_i; \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{j=1}^K \pi_j^{(t)} N(x_i; \mu_j^{(t)}, \Sigma_j^{(t)})} = S_{ik}^{(t)} \quad (\text{可以算出固定的值})$$

where  $N(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu))$  is the p.d.f. of

Gaussian distribution, with mean  $\mu$ , and covariance matrix  $\Sigma$

$$\begin{aligned} \triangleright \log P(x_i, z_i=k; \theta^{(t)}) &= \log \left( \pi_k \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} \exp \left( -\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right) \right) \\ &= \log \frac{\pi_k}{\sqrt{(2\pi)^m |\Sigma|}} - \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \quad (\text{A very simple function w.r.t. } \theta \text{ as desired}) \end{aligned}$$

• E-Step:

$$\begin{aligned} Q(\theta | \theta^{(t)}) &= \sum_{i=1}^N \mathbb{E}_{z_i \sim p(\cdot | x_i; \theta^{(t)})} [\log P(x_i, z_i=k; \theta)] \\ &= \sum_{i=1}^N \sum_{k=1}^K p(z_i=k | x_i; \theta^{(t)}) \log P(x_i, z_i=k; \theta) \\ &= \sum_{i=1}^N \sum_{k=1}^K \delta_{ik}^{(t)} \underbrace{\left[ \log \frac{\pi_k^{(t)}}{\sqrt{(2\pi)^m |\Sigma|}} - \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right]}_{\text{red line}} \end{aligned}$$

• M-Step: To maximize  $Q(\theta | \theta^{(t)})$ , we compute  $\nabla_{\mu_k} Q(\theta | \theta^{(t)})$  and  $\nabla_{\Sigma_k} Q(\theta | \theta^{(t)})$

$$\nabla_{\mu_k} Q(\theta | \theta^{(t)}) = 0 \Rightarrow \text{Take } \mu_k^* = \frac{\sum_{i=1}^N \delta_{ik}^{(t)} x_i}{\sum_{i=1}^N \delta_{ik}^{(t)}}, \quad \nabla_{\Sigma_k} Q(\theta | \theta^{(t)}) = 0 \Rightarrow \text{Take } \Sigma_k^* = \frac{\delta_{ik}^{(t)} (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^N \delta_{ik}^{(t)}}$$

Maximum  $Q(\theta | \theta^{(t)})$  over  $\pi_1, \pi_2, \dots, \pi_K$  with constraint  $\pi_1 + \dots + \pi_K = 1$

$$\Rightarrow \text{Take } \pi_k^* = \frac{1}{N} \sum_{i=1}^N \delta_{ik}^{(t)}$$

代入之後  $\nabla_{\mu_k} Q(\theta | \theta^{(t)}) = 0$

代入後使

$$\frac{\partial}{\partial a_{ij}^k} Q(\theta | \theta^{(t)}) = 0$$

$$\nabla_{\mu_k} Q(\theta | \theta^{(t)}) = \sum_{i=1}^N \underline{\delta_{ik}^{(t)}} \underline{\Sigma_k^{-1}} (x_i - \mu_k) = \sum_{i=1}^N \sum_{l=1}^N \delta_{ik}^{(t)} (x_i - \mu_k)$$

因  $\delta_{ik}^{(t)}$  和  $\Sigma_k$  無關故也和  $\mu_k$  無關

$$\text{Let } \Sigma_k^{-1} = \begin{bmatrix} a_{11}^k & a_{12}^k & \dots & a_{1m}^k \\ a_{21}^k & a_{22}^k & \dots & a_{2m}^k \\ \vdots & \ddots & \ddots & \vdots \\ a_{m1}^k & a_{m2}^k & \dots & a_{mm}^k \end{bmatrix} \quad (\text{要 optimize } \Sigma_k \text{ 就是 optimize } \Sigma_k^{-1})$$

$$\frac{\partial}{\partial a_{ij}^k} Q(\theta | \theta^{(t)}) = \frac{1}{2} \sum_{l=1}^N \delta_{lk}^{(t)} e_j^T [\Sigma_k - (x_l - \mu_k)(x_l - \mu_k)^T] e_i = 0 \quad \forall i, j$$

$$\frac{\partial}{\partial \pi_k} (Q(\theta | \theta^{(t)}) - \lambda \sum_{j=1}^K \pi_j) = \sum_{i=1}^N \delta_{ik}^{(t)} \cdot \frac{1}{\pi_k} - \lambda \quad (\text{有 constraint 的條件算極值用 Lagrange Multiplier})$$

$$\text{Take } \pi_k^* = \frac{1}{\lambda} \sum_{i=1}^N \delta_{ik}^{(t)} \quad \because \pi_1^* + \dots + \pi_K^* = 1 \quad \therefore 1 = \frac{1}{\lambda} \sum_{k=1}^K \sum_{i=1}^N \delta_{ik}^{(t)} = \frac{N}{\lambda} \quad \therefore \lambda = N$$

$$\begin{aligned}
\nabla_{\mu_k} Q(\theta | \theta^{(t)}) &= \frac{\partial}{\partial \mu_k} \sum_{i=1}^N \sum_{k=1}^K \delta_{ik}^{(t)} \left[ \underbrace{\log \frac{\pi_k^{(t)}}{\sqrt{(2\pi)^m |\Sigma_k|}}}_{\text{const}} - \frac{1}{2} (\underline{x}_i - \underline{\mu}_k)^T \Sigma_k^{-1} (\underline{x}_i - \underline{\mu}_k) \right] \\
&= \sum_{i=1}^N \delta_{ik}^{(t)} \frac{\partial}{\partial \mu_k} - \frac{1}{2} \left( \underline{x}_i^T \Sigma_k^{-1} \underline{x}_i - \mu_k^T \Sigma_k^{-1} \underline{x}_i - \underline{x}_i^T \Sigma_k^{-1} \mu_k + \mu_k^T \Sigma_k^{-1} \mu_k \right) \\
&= \sum_{i=1}^N \delta_{ik}^{(t)} - \frac{1}{2} \left( -\Sigma_k^{-1} \underline{x}_i - (\underline{x}_i^T \Sigma_k^{-1})^T + [\Sigma_k^{-1} (\Sigma_k^{-1})^T] \mu_k \right) \\
&= \sum_{i=1}^N \delta_{ik}^{(t)} - \frac{1}{2} \left( -\Sigma_k^{-1} \underline{x}_i - \Sigma_k^{-1} \underline{x}_i + 2 \Sigma_k^{-1} \mu_k \right) = \sum_{i=1}^N \delta_{ik}^{(t)} \left( \Sigma_k^{-1} \underline{x}_i - \Sigma_k^{-1} \mu_k \right) \\
&= \sum_{i=1}^N \delta_{ik}^{(t)} \Sigma_k^{-1} (\underline{x}_i - \mu_k) \quad \text{Q.E.D.}
\end{aligned}$$

$$\begin{aligned}
\nabla_{\pi_k} \left( Q(\theta | \theta^{(t)}) - \lambda \sum_{j=1}^K \pi_j \right) &= \frac{\partial}{\partial \pi_k} \sum_{i=1}^N \sum_{k=1}^K \delta_{ik}^{(t)} \left[ \log \pi_k^{(t)} - \underbrace{\log \sqrt{(2\pi)^m |\Sigma_k|}}_{\text{const}} - \frac{1}{2} (\underline{x}_i - \underline{\mu}_k)^T \Sigma_k^{-1} (\underline{x}_i - \underline{\mu}_k) \right] \\
&= \sum_{i=1}^N \delta_{ik}^{(t)} \frac{1}{\pi_k} - \lambda \cdot 1
\end{aligned}$$