

## Problem 1: (Automatic Differentiation) (0.5%)

Let  $f = f_4 \circ f_3 \circ f_2 \circ f_1$  be a differentiable function defined as

$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\mathbf{x} \mapsto f_4(f_3(f_2(f_1(\mathbf{x}))))$$

where  $f_1 : \mathbb{R}^n \rightarrow \mathbb{R}^{n_1}$ ,  $f_2 : \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$ ,  $f_3 : \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$ , and  $f_4 : \mathbb{R}^{n_3} \rightarrow \mathbb{R}$ . We assume  $n \gg 1$  and  $n \geq n_1 \geq n_2 \geq n_3 > 1$ .

1. (0.2%) Write down the derivative of  $f$  with respect to  $\mathbf{x}$  using **chain rule**.

Next, express **forward-mode** and **reverse-mode** auto-differentiation respectively with the formula you just wrote down. The reverse-mode auto-differentiation is also called **backpropagation**. In this problem, please "assign variables" when using chain rule. For example, suppose  $y = g(f(\mathbf{x}))$ . Let  $w_1 = f(\mathbf{x})$  and  $y = w_2 = g(w_1)$ . Then the partial derivative of  $y$  with respect to  $\mathbf{x}$  is given by

$$\frac{\partial y}{\partial \mathbf{x}} = \frac{\partial y}{\partial w_2} \frac{\partial w_2}{\partial w_1} \frac{\partial w_1}{\partial \mathbf{x}} = 1 \cdot \frac{\partial w_2}{\partial w_1} \frac{\partial w_1}{\partial \mathbf{x}} = \frac{\partial w_2}{\partial w_1} \frac{\partial w_1}{\partial \mathbf{x}} \left( = \frac{\partial y}{\partial w_1} \frac{\partial w_1}{\partial \mathbf{x}} \right).$$

$$1. w_1 = f_1(\mathbf{x}), \text{ 其中 } w_1 \in \mathbb{R}^{n_1}$$

$$w_2 = f_2(w_1), \text{ 其中 } w_2 \in \mathbb{R}^{n_2}$$

$$w_3 = f_3(w_2), \text{ 其中 } w_3 \in \mathbb{R}^{n_3}$$

$$w_4 = f_4(w_3) = f(\mathbf{x}), \text{ 其中 } w_4 \in \mathbb{R}$$

$$\text{chain rule: } \frac{\partial f}{\partial \mathbf{x}} = \frac{\partial w_4}{\partial \mathbf{x}} = \frac{\partial w_4}{\partial w_3} \frac{\partial w_3}{\partial w_2} \frac{\partial w_2}{\partial w_1} \frac{\partial w_1}{\partial \mathbf{x}} \neq$$

$$\frac{\partial w_4}{\partial w_3} : 1 \times n_3, \frac{\partial w_3}{\partial w_2} = n_3 \times n_2, \frac{\partial w_2}{\partial w_1} = n_2 \times n_1, \frac{\partial w_1}{\partial \mathbf{x}} = n_1 \times n$$

### Forward Mode

$$w_1 = f_1(\mathbf{x}), w_1 \in \mathbb{R}^{n_1}$$

$$\text{方向導數: } \dot{w}_1 = \frac{\partial w_1}{\partial \mathbf{x}} \dot{\mathbf{x}}$$

$$w_2 = f_2(w_1), w_2 \in \mathbb{R}^{n_2}$$

$$\text{方向導數: } \dot{w}_2 = \frac{\partial w_2}{\partial w_1} \dot{w}_1$$

$$w_3 = f_3(w_2), w_3 \in \mathbb{R}^{n_3}$$

$$\text{方向導數: } \dot{w}_3 = \frac{\partial w_3}{\partial w_2} \dot{w}_2$$

$$w_4 = f_4(w_3), w_4 \in \mathbb{R}$$

$$\text{方向導數: } \dot{w}_4 = \frac{\partial w_4}{\partial w_3} \dot{w}_3$$

$$\Rightarrow \dot{w}_4 \text{ 是 } \mathbf{x} \text{ 的方向導數: } \dot{w}_4 = \frac{\partial f}{\partial \mathbf{x}} \dot{\mathbf{x}}$$

### Reverse Mode

$$\text{令 } \bar{w}_4 = \frac{\partial w_4}{\partial w_4} = 1 \Rightarrow \text{output 對自身的梯度}$$

$$\bar{w}_3 = \left( \frac{\partial w_4}{\partial w_3} \right)^T \bar{w}_4$$

$$\bar{w}_2 = \left( \frac{\partial w_3}{\partial w_2} \right)^T \bar{w}_3$$

$$\bar{w}_1 = \left( \frac{\partial w_2}{\partial w_1} \right)^T \bar{w}_2$$

$$\bar{\mathbf{x}} = \left( \frac{\partial w_1}{\partial \mathbf{x}} \right)^T \bar{w}_1$$

$$\Rightarrow \bar{\mathbf{x}} \text{ 是 } \mathbf{x} \text{ 的梯度 } \bar{\mathbf{x}} = \left( \frac{\partial f}{\partial \mathbf{x}} \right)^T$$

2. (0.3%) Which mode is better in terms of efficiency? Why? For the second question, you need to prove your answer. Note that if you only answer the first question, you will only get 0.01%. Hint: you may write down the size of Jacobian matrices and see how many scalar multiplications ( $a \times b$ ,  $a, b \in \mathbb{R}$ ) are performed in the forward-mode and reverse-mode auto-differentiation respectively.

Through this problem, we want students to tell the difference between the forward mode and the reverse mode, and to understand why the reverse mode is usually preferred for optimization over the forward mode. Some useful references:

- Automatic Differentiation in Machine Learning: a Survey
- CSC321 Lecture 10: Automatic Differentiation
- Wikipedia

## 2.

pf:

forward Mode:

$$\frac{\partial w_1}{\partial x} \bar{x} : \frac{\partial w_1}{\partial x} : n_1 \times n \text{ 矩阵}$$

complexity:  $O(n_1 n)$

$$w_2 = \frac{\partial w_2}{\partial w_1} \bar{w}_1 : \frac{\partial w_2}{\partial w_1} : n_2 \times n_1 \text{ 矩阵}$$

complexity:  $O(n_2 n_1)$

$$w_3 = \frac{\partial w_3}{\partial w_2} \bar{w}_2 : \frac{\partial w_3}{\partial w_2} : n_3 \times n_2 \text{ 矩阵}$$

complexity:  $O(n_3 n_2)$

$$w_4 = \frac{\partial w_4}{\partial w_3} \bar{w}_3 : \frac{\partial w_4}{\partial w_3} : 1 \times n_3 \text{ 的行向量}$$

complexity:  $O(n_3)$

$$O(n_1 n) + O(n_2 n_1) + O(n_3 n_2) + O(n_3)$$

當函數輸入維度遠大於輸出維度 ( $n >> 1$ )，reverse mode 效率好

1. forward mode 需要對每個輸入維度獨立計算導數，共要  $n$  次傳播

2. reverse mode 只需一次反向傳播即可得所有輸入維度的導數

由上向的推導，forward mode 需逐層進行矩陣向量乘法，reverse mode 進行矩陣向量的持置乘法，由於輸出是 scalar，reverse mode 的計算複雜度不隨輸入維度線性增長

Reverse Mode:

$$\bar{w}_3 = \left( \frac{\partial w_4}{\partial w_3} \right)^T \bar{w}_4 \Rightarrow \frac{\partial w_4}{\partial w_3} \text{ 是 } 1 \times n_3 \text{ 行向量} \Rightarrow O(n_3)$$

$$\bar{w}_2 = \left( \frac{\partial w_3}{\partial w_2} \right)^T \bar{w}_3 \Rightarrow \frac{\partial w_3}{\partial w_2} \text{ 是 } n_3 \times n_2 \text{ 矩阵} \Rightarrow O(n_3 n_2)$$

$$\bar{w}_1 = \left( \frac{\partial w_2}{\partial w_1} \right)^T \bar{w}_2 \Rightarrow \frac{\partial w_2}{\partial w_1} \text{ 是 } n_2 \times n_1 \text{ 矩阵} \Rightarrow O(n_2 n_1)$$

$$\bar{x} = \left( \frac{\partial w_1}{\partial x} \right)^T \bar{w}_1 \Rightarrow \frac{\partial w_1}{\partial x} \text{ 是 } n_1 \times n \text{ 矩阵} \Rightarrow O(nn_1)$$

$$\Rightarrow O(n_3) + O(n_2 n_3) + O(n_1 n_2) + O(nn_1)$$

## Problem 2 (Batch Normalization)(1%)

Batch normalization is a popular trick for training deep networks nowadays, which aims to preserve the distribution within hidden layers and avoids vanishing gradient issue. The algorithm can be written as below (see Algorithm 1):

---

### Algorithm 1 Batch Normalization

---

**Input** Feature from data points over a mini-batch  $B = (x_i)_{i=1}^m$   
**Output**  $y_i = BN_{\gamma, \beta}(x_i)$

- 1: **procedure** BATCHNORMALIZE( $B, \gamma, \beta$ )
- 2:    $\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$  ▷ mini-batch mean
- 3:    $\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$  ▷ mini-batch variance
- 4:   **for**  $i \leftarrow 1$  to  $m$  **do**
- 5:      $\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$  ▷ normalize
- 6:      $y_i \leftarrow \gamma \hat{x}_i + \beta$  ▷ scale and shift
- 7:   **end for**
- 8:   **return**
- 9: **end procedure**

---

During training we need to backpropagate the gradient of loss  $\ell$  through this transformation, as well as compute the gradients with respect to the parameters  $\gamma, \beta$ . Towards this end, please write down the close form expressions for  $\frac{\partial \ell}{\partial x_i}$ ,  $\frac{\partial \ell}{\partial \gamma}, \frac{\partial \ell}{\partial \beta}$  in terms of  $x_i, \mu_B, \sigma_B^2, \hat{x}_i, y_i$  (given by the forward pass) and  $\frac{\partial \ell}{\partial y_i}$  (given by the backward pass).

- Hint: You may first write down the close form expressions of  $\frac{\partial \ell}{\partial \hat{x}_i}, \frac{\partial \ell}{\partial \sigma_B^2}, \frac{\partial \ell}{\partial \mu_B}$ , and then use them to compute  $\frac{\partial \ell}{\partial x_i}, \frac{\partial \ell}{\partial \gamma}, \frac{\partial \ell}{\partial \beta}$ .

$$\frac{\partial \ell}{\partial y_i} \quad (\text{已知})$$

$$\text{根据 } y_i = \gamma \hat{x}_i + \beta$$

$$\begin{aligned} \frac{\partial \ell}{\partial \gamma} &= \sum_{i=1}^m \frac{\partial \ell}{\partial y_i} \cdot \frac{\partial y_i}{\partial \gamma} = \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \hat{x}_i \\ \frac{\partial \ell}{\partial \beta} &= \sum_{i=1}^m \frac{\partial \ell}{\partial y_i} \cdot \frac{\partial y_i}{\partial \beta} = \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \end{aligned}$$

$$\frac{\partial \ell}{\partial \hat{x}_i} : \frac{\partial \ell}{\partial \hat{x}_i} = \frac{\partial \ell}{\partial y_i} \cdot \frac{\partial y_i}{\partial \hat{x}_i} = \frac{\partial \ell}{\partial y_i} \cdot \gamma$$

$$\frac{\partial \hat{x}_i}{\partial \sigma_B^2} :$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} = (x_i - \mu_B) \cdot (\sigma_B^2 + \epsilon)^{-1/2}$$

$$\begin{aligned} \frac{\partial \hat{x}_i}{\partial \sigma_B^2} &= (x_i - \mu_B) \cdot \left( -\frac{1}{2} (\sigma_B^2 + \epsilon)^{-3/2} \right) \\ &= -\frac{1}{2} (x_i - \mu_B) (\sigma_B^2 + \epsilon)^{-3/2} \end{aligned}$$

$$\frac{\partial \hat{x}_i}{\partial \mu_B} : \frac{\partial \hat{x}_i}{\partial \mu_B} = \frac{\partial}{\partial \mu_B} \left( \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \right) = -\frac{1}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\frac{\partial \ell}{\partial \sigma_B^2} : \frac{\partial \ell}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial \sigma_B^2} = -\frac{1}{2} \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} (x_i - \mu_B) (\sigma_B^2 + \epsilon)^{-3/2}$$

$$\frac{\partial \ell}{\partial \mu_B} : \frac{\partial \ell}{\partial \mu_B} = \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial \mu_B} = -\frac{1}{\sqrt{\sigma_B^2 + \epsilon}} \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i}$$

$$\frac{\partial \ell}{\partial x_i} = \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial x_i} + \frac{\partial \ell}{\partial \mu_B} \cdot \frac{\partial \mu_B}{\partial x_i} + \frac{\partial \ell}{\partial \sigma_B^2} \cdot \frac{\partial \sigma_B^2}{\partial x_i}$$

$$\frac{\partial \hat{x}_i}{\partial x_i} = \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} \quad \frac{\partial \mu_B}{\partial x_i} = \frac{1}{m} \quad \frac{\partial \sigma_B^2}{\partial x_i} = \frac{2(x_i - \mu_B)}{m}$$

$$\Rightarrow \frac{\partial \ell}{\partial x_i} = \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial \ell}{\partial \mu_B} \cdot \frac{1}{m} + \frac{\partial \ell}{\partial \sigma_B^2} \cdot \frac{2(x_i - \mu_B)}{m}$$

$$= \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} - \left( \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} \right) \cdot \frac{1}{m} - \left( \frac{1}{2} \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} (x_i - \mu_B) (\sigma_B^2 + \epsilon)^{-3/2} \right) \cdot \frac{2(x_i - \mu_B)}{m}$$

$$= \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} \left( \frac{\partial \ell}{\partial \hat{x}_i} - \frac{1}{m} \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{(x_i - \mu_B)}{(\sigma_B^2 + \epsilon)} \cdot \frac{1}{m} \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} (x_i - \mu_B) \right)$$

$$= \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} \left( \frac{\partial \ell}{\partial \hat{x}_i} - \frac{1}{m} \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} - \hat{x}_i \cdot \frac{1}{m} \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \hat{x}_i \right)$$

#

### Problem 3 (Multiclass AdaBoost)(1.5%)

Let  $\mathcal{X}$  be the input space,  $\mathcal{F}$  be a collection of multiclass classifiers that map from  $\mathcal{X}$  to  $[1, K]$ , where  $K$  denotes the number of classes. Let  $\{(x_i, \hat{y}_i)\}_{i=1}^m$  be the training data set, where  $x_i \in \mathcal{X}$  and  $\hat{y}_i \in [1, K]$ . Given  $T \in \mathbb{N}$ , suppose we want to find functions

$$g_{T+1}^k(x) = \sum_{t=1}^T \alpha_t f_t^k(x), \quad k \in [1, K]$$

where  $f_t \in \mathcal{F}$  and  $\alpha_t \in \mathbb{R}$  for all  $t \in [1, T]$ . Here for  $f \in \mathcal{F}$ , we denote  $f^k(x) = \mathbf{1}\{f(x) = k\}$ , where  $\mathbf{1}(\cdot)$  is an indicator function, as the  $k$ 'th element in the one-hot representation of  $f(x) \in [1, K]$ . The aggregated classifier  $h : \mathcal{X} \rightarrow [1, K]$  is defined as

$$x \mapsto \operatorname{argmax}_{1 \leq k \leq K} g_{T+1}^k(x)$$

Please apply gradient boosting to show how the functions  $f_t$  and coefficients  $\alpha_t$  are computed with an aim to minimize the following loss function

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

$$g_t^k(x) = 0, \quad \forall k \in [1, K], \quad \forall x \in \mathcal{X}$$

$$\text{negative gradient of the loss function: } r_{ik} = -\frac{\partial L}{\partial g_t^k(x_i)}$$

$$\frac{\partial L}{\partial g_t^k(x_i)} = \exp \left( \frac{1}{K-1} s_i \right) \cdot \left( \frac{1}{K-1} \delta_{k \neq \hat{y}_i} - \frac{k-1}{K-1} \delta_{k=\hat{y}_i} \right)$$

$$s_i = \sum_{j \neq \hat{y}_i} (g_t^j(x_i) - g_t^{\hat{y}_i}(x_i)), \quad \delta_p = \text{指示函数 } P=1 (\text{true}) \quad P=0 (\text{False})$$

$$k = \hat{y}_i: \quad r_{ij} = -(-\exp(\frac{1}{K-1} s_i)) = \exp(\frac{1}{K-1} s_i)$$

$$k \neq \hat{y}_i: \quad r_{ik} = -\left(\frac{1}{K-1} \exp\left(\frac{1}{K-1} s_i\right)\right) = -\frac{1}{K-1} \exp\left(\frac{1}{K-1} s_i\right)$$

$$w_i = \exp\left(\frac{1}{K-1} s_i\right), \quad \text{最小化加权错误率 } E_L = \frac{\sum_{i \in \mathcal{I}} w_i \cdot \mathbf{1}\{f_t(x_i) \neq \hat{y}_i\}}{\sum_{i \in \mathcal{I}} w_i}$$

找最优的  $\alpha_t$ , 使得沿着  $\alpha_t$  的方向, 损失函数最小

$$g_{T+1}^k(x) = g_T^k(x) + \alpha_t f_t^k(x), \quad L(\alpha_t) = \sum_{i=1}^m \exp\left(\frac{1}{K-1} (s_i + \alpha_t \cdot \mathbf{1}_{f_t(x_i) \neq \hat{y}_i})\right), \quad z_i = \sum_{k \neq \hat{y}_i} (f_t^k(x_i) - f_t^{\hat{y}_i}(x_i)) = \begin{cases} -(k-1), & f_t(x_i) = \hat{y}_i \\ 1, & f_t(x_i) \neq \hat{y}_i \end{cases}$$

$$\text{if } f_t(x_i) = \hat{y}_i \Rightarrow L_1 = \sum_{i: f_t(x_i) = \hat{y}_i} w_i \cdot \exp(-\alpha_t \cdot (k-1)/(K-1)) = \sum_{i: f_t(x_i) = \hat{y}_i} w_i \cdot \exp(-\alpha_t) \quad \Rightarrow L(\alpha_t) = L_1 + L_2$$

$$\text{if } f_t(x_i) \neq \hat{y}_i \Rightarrow L_2 = \sum_{i: f_t(x_i) \neq \hat{y}_i} w_i \cdot \exp(\alpha_t \cdot 1/(K-1)) = \sum_{i: f_t(x_i) \neq \hat{y}_i} w_i \cdot \exp(\frac{\alpha_t}{K-1})$$

$$\text{最小化 } L(\alpha_t) \text{ 使 } \frac{dL}{d\alpha_t} = -L_1 \exp(-\alpha_t) + \frac{1}{K-1} L_2 \exp(\frac{\alpha_t}{K-1}) = 0 \Rightarrow -L_1 \cdot \exp(-\alpha_t) + \frac{L_2}{K-1} \cdot \exp(\frac{\alpha_t}{K-1}) = 0 \quad \left( L_1: \sum_{i: f_t(x_i) = \hat{y}_i} w_i, L_2: \sum_{i: f_t(x_i) \neq \hat{y}_i} w_i \right)$$

$$\Rightarrow -(k-1)L_1 + L_2 \exp(\alpha_t \cdot \frac{1}{K-1}) = 0 \Rightarrow \exp(\alpha_t \cdot \frac{1}{K-1}) = \frac{(k-1)L_1}{L_2} \Rightarrow \alpha_t = \frac{k-1}{K} \ln\left(\frac{(k-1)L_1}{L_2}\right)$$

$$g_{T+1}^k(x) = g_T^k(x) + \alpha_t f_t^k(x)$$

$$\text{- Sample weights } w_i = \exp\left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} [g_T^k(x_i) - g_T^{\hat{y}_i}(x_i)]\right)$$

fit  $f_t$  to minimize weighted misclassification error using these weights.

$$\alpha_t = \frac{k-1}{K} \left[ \ln \left( (k-1) \sum_{i: f_t(x_i) = \hat{y}_i} w_i \right) - \ln \left( \sum_{i: f_t(x_i) \neq \hat{y}_i} w_i \right) \right]$$

## Problem 4 (AdaBoosting) (0.5%)

1. Consider training an AdaBoosting classifier using decision stumps on the data set illustrated in Figure 1:



Figure 1: AdaBoost Data set

- (a) Which examples will have their weights increased at the end of the first iteration? Circle them.  
 (b) How many iterations will it take to achieve zero training error? Justify your answers and show each iteration step.

| (a) 在第一次迭代結束時，只有被錯誤分類的樣本權重增加，初始時所有樣本的權重相等  
 第一次迭代中，最佳AdaBoost 將所有樣本預測為藍色叉號，可最小化錯誤分類的樣本數量。  
 只有紅色圓圈的權重會增加

(b) 4個藍叉： $+1$  (樣本1~4)

1個紅圓： $-1$  (樣本5)

初始權重： $0.2$

第一次迭代：

$h_1$  純分類器：所有樣本都預測為 $+1$

$$E_1 = w_5^{(1)} = 0.2$$

$$\alpha_1 = \frac{1}{2} \ln\left(\frac{1-E_1}{E_1}\right) = \frac{1}{2} \ln\left(\frac{0.8}{0.2}\right) = 0.693$$

$$\text{更新：錯誤 } w_5^{(2)} = w_5^{(1)} \times e^{\alpha_1} = 0.2 \times e^{+0.693} = 0.4$$

$$\text{正確 } w_i^{(2)} = w_i^{(1)} \times e^{-\alpha_1} = 0.2 \times e^{-0.693} = 0.1$$

$$\text{總權重: } 0.4 + 4 \times 0.1 = 0.8$$

$$\text{樣本1-4: } w_i^{(4)} = \frac{0.1}{0.8} = 0.125 \quad \text{樣本5: } w_5^{(4)} = \frac{0.4}{0.8} = 0.5$$

第二次迭代：

$h_2$ : 紅圓正確分類為 $-1$ ，誤分2個藍叉 (樣本1, 2, 4)

$$E_2 = w_1^{(4)} + w_2^{(4)} = 0.125 + 0.125 = 0.25$$

$$\alpha_2 = \frac{1}{2} \ln\left(\frac{1-E_2}{E_2}\right) = \frac{1}{2} \ln 3 = 0.5493$$

$$\text{更新：錯誤 } w_5^{(3)} = w_5^{(2)} \times e^{\alpha_2} = 0.5 \times e^{0.5493} = 0.732$$

正確：

$$\text{樣本3, 8, 4: } w_i^{(3)} = w_i^{(2)} \times e^{-\alpha_2} = 0.125 \times e^{-0.5493} = 0.0722$$

$$\text{樣本5: } w_5^{(3)} = w_5^{(2)} \times e^{-\alpha_2} = 0.5 \times e^{-0.5493} = 0.2877$$

$$\text{總權重: } 2 \times 0.0722 + 2 \times 0.2877 + 0.25 = 0.8661$$

$$\text{樣本1, 2: } w_i^{(2)} = \frac{0.0722}{0.8661} = 0.25 \quad \text{樣本3, 8, 4: } w_i^{(2)} = \frac{0.2877}{0.8661} = 0.3283$$

$$\text{樣本5: } w_5^{(2)} = \frac{0.25}{0.8661} = 0.2833$$

第三次迭代：

$h_3$ : 單純關注樣本1和2，可能誤分樣本5

$$E_3 = w_5^{(3)} = 0.2833$$

$$\alpha_3 = \alpha_2 + \frac{1}{2} \ln\left(\frac{1-E_3}{E_3}\right) = \frac{1}{2} \ln\left(\frac{0.6666}{0.2833}\right) = \frac{1}{2} \ln 2 = 0.3466$$

$$\text{更新：錯誤 (樣本5): } w_5^{(4)} = w_5^{(3)} \times e^{\alpha_3} = 0.2833 \times e^{0.3466} = 0.4715$$

正確：

$$\text{樣本1, 2: } w_i^{(4)} = 0.125 \times e^{-0.3466} = 0.1768$$

$$\text{樣本3, 8, 4: } w_i^{(4)} = 0.0723 \times e^{-0.3466} = 0.0589$$

$$\text{總權重: } 0.4715 + 2 \times 0.1768 + 2 \times 0.0589 = 0.9429$$

$$\text{新權重分佈: 樣本1, 2: } w_i^{(4)} = \frac{0.1768}{0.9429} = 0.1877$$

$$\text{樣本3, 8, 4: } w_i^{(4)} = \frac{0.0589}{0.9429} = 0.0625$$

$$\text{樣本5: } w_5^{(4)} = \frac{0.4715}{0.9429} = 0.5$$

$$\text{最終強分類器: } H(x) = \text{sign}(w_1 h_1(x) + w_2 h_2(x) + w_3 h_3(x))$$

樣本1, 2 (+1):

$$h_1(x) = +1 \quad h_2(x) = -1 \quad (\text{誤分類}), h_3(x) = +1$$

$$w_1(+1) + w_2(-1) + w_3(+1) = 0.4924 > 0 \quad H(x) = +1 \rightarrow \text{分類正確}$$

樣本3, 8, 4 (+1):

$$h_1(x) = +1, h_2(x) = +1, h_3(x) = -1$$

$$w_1(+1) + w_2(+1) + w_3(-1) = 0.8958 > 0 \quad H(x) = +1 \rightarrow \text{分類正確}$$

樣本5 (-1):

$$h_1(x) = +1 \quad (\text{誤分類}), h_2(x) = -1, h_3(x) = -1$$

$$w_1(+1) + w_2(-1) + w_3(-1) = -0.2268 < 0 \quad H(x) = -1 \rightarrow \text{分類正確}$$

⇒ 3次迭代，組合而成的強分類器可正確分類所有訓練樣本，這確認

2. Suppose AdaBoost is run on  $N$  training examples, and suppose on each round that the weighted training error  $\epsilon_t$  of the  $t$ 'th weak hypothesis is at most  $1/2 - \gamma$ , for some number  $0 < \gamma < 1/2$ . After how many iterations,  $T$ , will the combined hypothesis be consistent with the  $N$  training examples, i.e., achieves zero training error? Your answer should only be expressed in terms of  $N$  and  $\gamma$ . (Hint: Recall that exponential loss is an upper bound for 0-1 loss. What is the training error when 1 example is misclassified?)

$$E_{\text{train}} = \frac{1}{N} \sum_{i=1}^N I(y_i \neq \text{sign}(F(x_i))) , F(x) = \sum_{t=1}^T \eta_t h_t(x) \text{ 為組合分類器, } \eta_t \text{ 是第 } t \text{ 次的分類器權重, } h_t(x) \text{ 是第 } t \text{ 次的分類器}$$

$$E_{\text{exp}} = \sum_{i=1}^N e^{-x_i F(x)} \quad 0-1 \text{ 損失上界: } I(y_i \neq \text{sign}(F(x_i))) \leq e^{-x_i F(x)} \Rightarrow E_{\text{train}} \leq \frac{1}{N} E_{\text{exp}}$$

$$\text{AdaBoost 由: } E_{\text{exp}} = \prod_{t=1}^T Z_t \quad Z_t = \epsilon_t^{1/(1-\epsilon_t)}$$

$$\text{由於 } \epsilon_t = \frac{1}{2} - \gamma \Rightarrow 1 - \epsilon_t = \frac{1}{2} + \gamma$$

$$Z_t \leq 2 \sqrt{(\frac{1}{2} - \gamma)(\frac{1}{2} + \gamma)} = 2 \sqrt{\frac{1}{4} - \gamma^2} = \sqrt{1 - 4\gamma^2} \Rightarrow \sqrt{1 - 4\gamma^2} \leq e^{-2\gamma}$$

$$\because 0 < \gamma < \frac{1}{2} \Rightarrow \sqrt{1 - 4\gamma^2} \leq e^{-2\gamma} \therefore Z_t \leq 2 \times \frac{1}{2} e^{-2\gamma} = e^{-2\gamma}$$

$$\therefore \text{每次 } Z_t \leq e^{-2\gamma} \therefore E_{\text{exp}} = \prod_{t=1}^T Z_t \leq (e^{-2\gamma})^T = e^{-2\gamma T}$$

$$E_{\text{train}} \leq \frac{1}{N} E_{\text{exp}} \leq \frac{1}{N} e^{-2\gamma T} \quad (E_{\text{train}} \leq \frac{1}{N}, \text{所有樣本都被正確分類})$$

$$e^{-2\gamma T} \leq 1 \Rightarrow -2\gamma T \leq -\ln N \Rightarrow T \geq \frac{\ln N}{2\gamma}$$

至少  $T \geq \frac{\ln N}{2\gamma}$  次迭代，才能在  $N$  個訓練樣本上達 0 到訓練誤差

## Problem 5 (Gradient Descent Convergence) (1.5%)

Suppose the function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable. Also,  $f$  is  $\beta$ -smoothness and  $\alpha$ -strongly convex.

$$\beta\text{-smoothness : } \beta > 0, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq \beta \|\mathbf{x} - \mathbf{y}\|_2$$

$$\alpha\text{-strongly convex : } \alpha > 0, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) \geq \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$

Then we propose a gradient descent algorithm

1. Find a initial  $\theta^0$ .
2. Let  $\theta^{n+1} = \theta^n - \eta \nabla_{\theta} f(\theta^n)$   $\forall n \geq 0$ , where  $\eta = \frac{1}{\beta}$ .

The following problems lead you to prove the gradient descent convergence.

1. Prove the property of  $\beta$ -smoothness function

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

- (a) Define  $g : \mathbb{R} \rightarrow \mathbb{R}, g(t) = f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))$ . Show that  $f(\mathbf{y}) - f(\mathbf{x}) = \int_0^1 g'(t) dt$ .
- (b) Show that  $g'(t) = \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^T(\mathbf{y} - \mathbf{x})$ .
- (c) Show that  $|f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})| \leq \int_0^1 |(\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}))^T(\mathbf{y} - \mathbf{x})| dt$ .
- (d) By Cauchy-Schwarz inequality and the definition of  $\beta$ -smoothness, show that  $|f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})| \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$ , hence we get

$$f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

$$1. (a). g(t) = f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))$$

$$\Rightarrow \begin{cases} g(0) = f(\mathbf{x} + 0(\mathbf{y} - \mathbf{x})) = f(\mathbf{x}) \\ g(1) = f(\mathbf{x} + 1(\mathbf{y} - \mathbf{x})) = f(\mathbf{y}) \end{cases} \Rightarrow f(\mathbf{y}) - f(\mathbf{x}) = g(1) - g(0) = \int_0^1 g'(t) dt$$

$$(b) \text{ 鏡式: } g'(t) = \frac{d}{dt} f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) = \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^T \cdot \frac{d}{dt} (\mathbf{x} + t(\mathbf{y} - \mathbf{x})) = \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^T(\mathbf{y} - \mathbf{x})$$

$$(c) \text{ 由(a),(b)得: } f(\mathbf{y}) - f(\mathbf{x}) = \int_0^1 \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^T(\mathbf{y} - \mathbf{x}) dt$$

$$\Rightarrow f(\mathbf{y}) - f(\mathbf{x}) = \int_0^1 [\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^T(\mathbf{y} - \mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})] dt$$

$$f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) = \int_0^1 [\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})]^T(\mathbf{y} - \mathbf{x}) dt$$

$$|f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})| \leq \int_0^1 \|[\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})]^T(\mathbf{y} - \mathbf{x})\| dt$$

- (d) Apply the Cauchy-Schwarz inequality to the integrand:

$$|[\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})]^T(\mathbf{y} - \mathbf{x})| \leq \|\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\|_2 \cdot \|\mathbf{y} - \mathbf{x}\|_2$$

$\beta$ -smoothness property:

$$\|\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\|_2 \leq \beta \|\mathbf{x} + t(\mathbf{y} - \mathbf{x}) - \mathbf{x}\|_2 = \beta t \|\mathbf{y} - \mathbf{x}\|_2$$

$$\Rightarrow |[\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})]^T(\mathbf{y} - \mathbf{x})| \leq \beta t \|\mathbf{y} - \mathbf{x}\|_2^2$$

$$|f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})| \leq \int_0^1 \beta t \|\mathbf{y} - \mathbf{x}\|_2^2 dt = \beta \|\mathbf{y} - \mathbf{x}\|_2^2 \int_0^1 t dt$$

$$\int_0^1 t dt = \left[ \frac{t^2}{2} \right]_0^1 = \frac{1}{2}$$

$$|f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})| \leq \beta \|\mathbf{y} - \mathbf{x}\|_2^2 \cdot \frac{1}{2} = \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

$$\text{左式是} \neq \text{右式: } f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

2. Let  $\mathbf{y} = \mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})$  and use 1., prove that

$$f\left(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})\right) - f(\mathbf{x}) \leq -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2$$

and

$$f(\mathbf{x}^*) - f(\mathbf{x}) \leq -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2,$$

where  $\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$ .

3. Show that  $\forall n \geq 0$ ,

$$\|\boldsymbol{\theta}^{n+1} - \boldsymbol{\theta}^*\|_2^2 = \|\boldsymbol{\theta}^n - \boldsymbol{\theta}^*\|_2^2 + \eta^2 \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n)\|_2^2 - 2\eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n)^T (\boldsymbol{\theta}^n - \boldsymbol{\theta}^*),$$

where  $\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$ .

$$2. \quad \mathbf{y} - \mathbf{x} = (\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})) - \mathbf{x} = -\frac{1}{\beta} \nabla f(\mathbf{x})$$

$$\|\nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x})\|_2 = \|\nabla f(\mathbf{x})^T (-\frac{1}{\beta} \nabla f(\mathbf{x}))\|_2 = \frac{1}{\beta} \|\nabla f(\mathbf{x})\|_2^2$$

$$\|\mathbf{y} - \mathbf{x}\|_2^2 = \left\| -\frac{1}{\beta} \nabla f(\mathbf{x}) \right\|_2^2 = \left( \frac{1}{\beta} \right)^2 \|\nabla f(\mathbf{x})\|_2^2 = \frac{1}{\beta^2} \|\nabla f(\mathbf{x})\|_2^2$$

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \left(-\frac{1}{\beta}\right) \|\nabla f(\mathbf{x})\|_2^2 + \frac{\beta}{2} \left(\frac{1}{\beta^2}\right) \|\nabla f(\mathbf{x})\|_2^2 = f(\mathbf{x}) - \frac{1}{\beta} \|\nabla f(\mathbf{x})\|_2^2 + \frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2 = f(\mathbf{x}) - \left(\frac{1}{\beta} - \frac{1}{2\beta}\right) \|\nabla f(\mathbf{x})\|_2^2 = f(\mathbf{x}) - \frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2$$

$$\text{得} \Rightarrow f(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})) - f(\mathbf{x}) \leq -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2$$

$f(\mathbf{x}^*) \leq f(\mathbf{y})$  :  $\mathbf{x}^*$  是  $f$  的全局最小值點，對任何  $\mathbf{y}$  都有  $f(\mathbf{x}^*) \leq f(\mathbf{y})$

$$\Rightarrow \mathbf{y} = \mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x}) : f(\mathbf{x}^*) \leq f(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x}))$$

$$\text{由前面的推導知: } f(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})) - f(\mathbf{x}) \leq -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2$$

$$\therefore f(\mathbf{x}^*) - f(\mathbf{x}) \leq f(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})) - f(\mathbf{x}) \leq -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2$$

$$3. \quad \text{梯度下降: } \boldsymbol{\theta}^{n+1} = \boldsymbol{\theta}^n - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n)$$

$$\boldsymbol{\theta}^{n+1} - \boldsymbol{\theta}^* = (\boldsymbol{\theta}^n - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n)) - \boldsymbol{\theta}^* = (\boldsymbol{\theta}^n - \boldsymbol{\theta}^*) - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n)$$

$$\|\boldsymbol{\theta}^{n+1} - \boldsymbol{\theta}^*\|_2^2 = \|(\boldsymbol{\theta}^n - \boldsymbol{\theta}^*) - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n)\|_2^2 = (\boldsymbol{\theta}^n - \boldsymbol{\theta}^* - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n))^T (\boldsymbol{\theta}^n - \boldsymbol{\theta}^* - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n))$$

$$= (\boldsymbol{\theta}^n - \boldsymbol{\theta}^*)^T (\boldsymbol{\theta}^n - \boldsymbol{\theta}^*) - 2\eta (\boldsymbol{\theta}^n - \boldsymbol{\theta}^*)^T \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n) + \eta^2 (\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n))^T \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n)$$

$$= \|\boldsymbol{\theta}^n - \boldsymbol{\theta}^*\|_2^2 - 2\eta (\boldsymbol{\theta}^n - \boldsymbol{\theta}^*)^T \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n) + \eta^2 \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n)\|_2^2$$

$$\text{得} \quad \|\boldsymbol{\theta}^{n+1} - \boldsymbol{\theta}^*\|_2^2 = \|\boldsymbol{\theta}^n - \boldsymbol{\theta}^*\|_2^2 + \eta^2 \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n)\|_2^2 - 2\eta (\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n))^T (\boldsymbol{\theta}^n - \boldsymbol{\theta}^*)$$

4. Use 2. and the definition of  $\alpha$ -strongly convex to prove  $\forall n \geq 0$

$$\|\theta^{n+1} - \theta^*\|_2^2 \leq (1 - \frac{\alpha}{\beta})\|\theta^n - \theta^*\|_2^2,$$

where  $\theta^* = \arg \min_{\theta} f(\theta)$ .

5. Use the above inequality to show that  $\theta^n$  will converge to  $\theta^*$  when  $n$  goes to infinity.

4. 梯度下降:  $\theta^{n+1} = \theta^n - \eta \nabla f(\theta^n)$      $\eta = \frac{\alpha}{\beta}$

$$\|\theta^{n+1} - \theta^*\|_2^2 = \|\theta^n - \theta^*\|_2^2 + \eta^2 \|\nabla f(\theta^n)\|_2^2 - 2\eta \nabla f(\theta^n)^T (\theta^n - \theta^*)$$

$\alpha$ -strongly convex:

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{\alpha}{2} \|y - x\|_2^2 \quad \xrightarrow{\text{由 } \nabla f(\theta^*) = 0 \text{ 得 } f(\theta^*) \geq f(\theta) + \frac{\alpha}{2} \|\theta - \theta^*\|_2^2} \\ \Rightarrow f(\theta^n) - f(\theta^*) \geq \frac{\alpha}{2} \|\theta^n - \theta^*\|_2^2$$

$$\beta\text{-smoothness: } \|\nabla f(\theta^n)\|_2^2 \leq 2\beta(f(\theta^n) - f(\theta^*)) \quad \because \nabla f(\theta^*) = 0 \Rightarrow f(\theta^{n+1}) \leq f(\theta^n) - \frac{1}{2\beta} \|\nabla f(\theta^n)\|_2^2$$

$$\|\nabla f(\theta^n)\|_2^2 \geq 2\alpha(f(\theta^n) - f(\theta^*))$$

$$f(\theta^{n+1}) - f(\theta^*) \leq f(\theta^n) - f(\theta^*) - \frac{1}{2\beta} \|\nabla f(\theta^n)\|_2^2 \leq f(\theta^n) - f(\theta^*) - \frac{\alpha}{\beta}(f(\theta^n) - f(\theta^*)) = ((1 - \frac{\alpha}{\beta})(f(\theta^n) - f(\theta^*)))$$

$$\because f(\theta^n) - f(\theta^*) \geq \frac{\alpha}{2} \|\theta^n - \theta^*\|_2^2 \quad \therefore \|\theta^n - \theta^*\|_2^2 \leq \frac{2}{\alpha} (f(\theta^n) - f(\theta^*))$$

$$\Rightarrow \|\theta^{n+1} - \theta^*\|_2^2 \leq \frac{\alpha}{2} (f(\theta^{n+1}) - f(\theta^*)) \leq \frac{2}{\alpha} ((1 - \frac{\alpha}{\beta})(f(\theta^n) - f(\theta^*))) \leq (1 - \frac{\alpha}{\beta}) \|\theta^n - \theta^*\|_2^2$$

$$\forall n \geq 0: \|\theta^{n+1} - \theta^*\|_2^2 \leq (1 - \frac{\alpha}{\beta}) \|\theta^n - \theta^*\|_2^2$$

$$5. \|\theta_{n+1} - \theta^*\|_2^2 \leq (1 - \frac{\alpha}{\beta}) \|\theta_n - \theta^*\|_2^2 \quad \because 0 < \frac{\alpha}{\beta} \leq 1, 0 \leq g = 1 - \frac{\alpha}{\beta} < 1$$

$$\|\theta_n - \theta^*\|_2^2 \leq g^n \|\theta^0 - \theta^*\|_2^2$$

$$n \rightarrow \infty, g^n \rightarrow 0 \quad \lim_{n \rightarrow \infty} \|\theta_n - \theta^*\|_2^2 = 0$$

$\theta_n$  converges to  $\theta^*$