

1.

臺灣大學

$$E(X) = \bar{X} \text{ (mean of } X)$$

$$\text{Var}(X) = E((X - \bar{X})^2) = \sigma_X^2$$

$$E(W) = 0 \text{ (mean of } W)$$

$$\text{Var}(W) = E((W - 0)^2) = \sigma_W^2$$

$$s_j^{(l)} = \sum_{i=1}^{d^{(l-1)}} w_{ij}^{(l)} x_i^{(l-1)}$$

$$= w_{j1}^{(l)} x_1^{(l-1)} + w_{j2}^{(l)} x_2^{(l-1)} + \dots + w_{jd}^{(l)} x_d^{(l-1)}$$

$$= \sum_{i=1}^{d^{(l-1)}} [w_{j1}^{(l)} x_i^{(l-1)}, w_{j2}^{(l)} x_i^{(l-1)}, \dots]$$

Because  $w_{ij}^{(l)}$  are independent to each other and to all  $x_i^{(l-1)}$ ,  
the covariance between  $x_1^{(l-1)}$  and  $w_{ij}^{(l)}$  would be 0.

$$\text{cov}(X, W) = E((X - \bar{X})(W - 0)) = E(X \cdot W) - \bar{X} \cdot 0 = 0$$

From above, we can know that the mean of  $s_j^{(l)}$  are  $E(X \cdot W)$ ,  
which equals to  $0 + \bar{X} \cdot 0$ , which is 0 and independent to each other.

2.

$$\text{Var}(s_j^{(l)})$$

$$= \text{Var}(\bar{X} \cdot W) \cdot d^{(l-1)}$$

$$= (E(\bar{X}^T X W W^T) - E(\bar{X} \cdot W)^2) \cdot d^{(l-1)}$$

$$= (\text{Var}(X) \text{Var}(W) + \text{Var}(X) \cdot E(W)^2 + \text{Var}(W) E(X)^2) \cdot d^{(l-1)}$$

$$= (\sigma_X^2 \sigma_W^2 + \sigma_X^2 \cdot 0 + \sigma_W^2 \cdot \bar{X}^2) \cdot d^{(l-1)}$$

$$= (\sigma_X^2 + \bar{X}^2) \sigma_W^2 d^{(l-1)}$$

3.

3.

We know that  $x_i^{(l-1)} = \max(s_i^{(l-1)}, 0)$

So,  $E[x_i^{(l-1)^2}] = \int_{-\infty}^{\infty} \max(s_i^{(l-1)}, 0)^2 p(s_i^{(l-1)}) ds_i^{(l-1)}$

$= \int_0^{\infty} (s_i^{(l-1)})^2 p(s_i^{(l-1)}) ds_i^{(l-1)}$

$(s_i^{(l-1)} \text{ are zero-mean, symmetric})$

$= \frac{1}{2} \int_{-\infty}^{\infty} (s_i^{(l-1)})^2 p(s_i^{(l-1)}) ds_i^{(l-1)}$

$= \frac{1}{2} E[(s_i^{(l-1)})^2]$  #

4.

4.

$\text{Var}(s_i^{(l)})$

$= d^{(l-1)} \text{Var}(w_{ij}^{(l)} x_i^{(l-1)})$

$= d^{(l-1)} \text{Var}(w_{ij}^{(l)}) E[(x_i^{(l-1)})^2]$

$(\text{by problem 3})$

$= d^{(l-1)} \text{Var}(w_{ij}^{(l)}) \cdot \frac{1}{2} E[(s_i^{(l-1)})^2]$

$(\text{by } s_i^{(l-1)} \text{ are zero-mean, symmetric random var})$

$= d^{(l-1)} \text{Var}(w_{ij}^{(l)}) \frac{1}{2} E[(s_i^{(l-1)} - E(s_i^{(l-1)}))^2]$

$\underbrace{\hspace{10em}}_{\text{Var}(s_i^{(l-1)})}$

$= d^{(l-1)} \underbrace{\text{Var}(w_{ij}^{(l)})}_{=\sigma_w^2} \frac{1}{2} \text{Var}(s_i^{(l-1)})$

$= \frac{d^{(l-1)}}{2} \sigma_w^2 \text{Var}(s_i^{(l-1)})$  #

5.

$$\begin{aligned}
 & \text{Var}(s_i^{(l)}) \\
 &= d^{(l-1)} \text{Var}(w_{ij}^{(l)}) E[\chi_i^{(l-1)^2}] \\
 &= \text{Var}(s_i^{(l-1)})
 \end{aligned}$$

$$\begin{aligned}
 & E[\chi_i^{(l-1)^2}] \\
 &= \int_{-\infty}^{\infty} \max(s_i^{(l-1)}, \alpha s_i^{(l-1)})^2 p(s_i^{(l-1)}) ds_i^{(l-1)} \\
 &= \frac{1+\alpha}{2} \int_{-\infty}^{\infty} (s_i^{(l-1)})^2 p(s_i^{(l-1)}) ds_i^{(l-1)} \\
 &= \frac{1+\alpha}{2} E[(s_i^{(l-1)})^2] = \frac{1+\alpha}{2} \text{Var}(s_i^{(l-1)})
 \end{aligned}$$

$$\begin{aligned}
 & d^{(l-1)} \text{Var}(w_{ij}^{(l)}) E[\chi_i^{(l-1)^2}] \\
 &= d^{(l-1)} \sigma_w^2 \frac{1+\alpha}{2} \text{Var}(s_i^{(l-1)}) = \text{Var}(s_i^{(l-1)})
 \end{aligned}$$

$$\frac{d^{(l-1)} \sigma_w^2 (1+\alpha)}{2} = 1$$

$$1+\alpha = \frac{2}{d^{(l-1)} \sigma_w^2} \quad \alpha = \frac{2}{d^{(l-1)} \sigma_w^2} - 1$$

6.

李慶學

$$\begin{aligned}
 V_1 &= \beta V_0 + (1-\beta) \Delta_1 = \beta(1-\beta) \Delta_1 \\
 V_2 &= \beta V_1 + (1-\beta) \Delta_2 = \beta^2(1-\beta) \Delta_1 + (1-\beta) \Delta_2 \\
 V_3 &= \beta V_2 + (1-\beta) \Delta_3 = \beta^3(1-\beta) \Delta_1 + \beta^2(1-\beta) \Delta_2 + (1-\beta) \Delta_3 \\
 V_4 &= \beta V_3 + (1-\beta) \Delta_4 = \beta^4(1-\beta) \Delta_1 + \beta^3(1-\beta) \Delta_2 + \beta^2(1-\beta) \Delta_3 + (1-\beta) \Delta_4
 \end{aligned}$$

At  $t=T$ , the formula results in

$$\beta V_0 V_T = \sum_{t=1}^T \beta^t (1-\beta) \Delta_t$$

$$\alpha_t = \beta^{T-t} (1-\beta)$$



7.

$$\begin{aligned}
 & \alpha_1 \leq \frac{1}{2} \\
 \Rightarrow & \beta^{T-1}(1-\beta) \leq \frac{1}{2} \longrightarrow \text{smallest } T \\
 & \ln(\beta^{T-1}(1-\beta)) \leq \ln \frac{1}{2} \\
 & \ln \beta^{T-1} + \ln(1-\beta) \leq \ln 1 - \ln 2 \\
 & \ln \beta^{T-1} - \ln \beta + \ln(1-\beta) \leq \ln 1 - \ln 2 \\
 & T \ln \beta - \ln \beta + \ln(1-\beta) \leq -\ln 2 \\
 & T \ln \beta \leq \ln \beta - \ln(1-\beta) - \ln 2 \\
 & T \ln \beta \leq \ln \left( \frac{\beta}{2(1-\beta)} \right) \\
 & \text{We know that } 0 < \beta < 1, \text{ so } \ln \beta < 0 \\
 & T \ln \beta \leq \ln \left( \frac{\beta}{2(1-\beta)} \right) \\
 & T \geq \frac{\ln \left( \frac{\beta}{2(1-\beta)} \right)}{\ln \beta}, \text{ the smallest } T \text{ is } \frac{\ln \left( \frac{\beta}{2(1-\beta)} \right)}{\ln \beta}
 \end{aligned}$$

8.

$$\begin{aligned}
 & \alpha_t = \frac{a_t}{\sum_{t=1}^T a_t} \\
 & = \frac{\beta^{T-t}(1-\beta)}{\sum_{t=1}^T \beta^{T-t}(1-\beta)} \\
 & = \frac{\beta^{T-1}(1-\beta)}{\beta^{T-1}(1-\beta) + \beta^{T-2}(1-\beta) + \dots + \beta^0(1-\beta)} \\
 & = \frac{\beta^{T-1}}{\beta^{T-1} + \beta^{T-2} + \dots + \beta^0} \\
 & = \frac{1}{\beta^{t-1} + \beta^{t-2} + \dots + \beta^{1-T}} \\
 & = \frac{1}{\sum_{k=1}^T \beta^{t-k}}
 \end{aligned}$$

9.

1.  $a_1' \leq \frac{1}{2}$   
 $\Rightarrow \sum_{k=1}^T \beta^{1-k} \leq \frac{1}{2}$   $\longrightarrow$  smallest  $T$

$\frac{\beta^{T+1}}{\sum_{k=1}^T \beta^{1-k}} \leq \frac{1}{2}$   
 $\sum_{k=1}^T \beta^{1-k} \geq 2$   
 $\beta \sum_{k=1}^T \beta^{-k} \geq 2$   
 $\beta \left( \frac{1}{\beta^1} + \frac{1}{\beta^2} + \dots + \frac{1}{\beta^T} \right) \geq 2$   
 $= \beta \sum_{n=0}^{T-1} \left( \frac{1}{\beta} \right) \left( \frac{1}{\beta} \right)^n \geq 2$   
 $\beta \left( \frac{1}{\beta} \right) \left( \frac{1 - \left( \frac{1}{\beta} \right)^T}{1 - \frac{1}{\beta}} \right) \geq 2$   $1 - \left( \frac{1}{\beta} \right)^T < 0$   
 $0 < \frac{1}{\beta} < 1$   
 $1 - \left( \frac{1}{\beta} \right)^T \leq 2 \left( 1 - \frac{1}{\beta} \right)$   
 $-\left( \frac{1}{\beta} \right)^T \leq 1 - \frac{2}{\beta}$   
 $-\left( \frac{1}{\beta} \right)^{T-1} \leq \beta - 2$   
 $(-1)(T-1) \ln \left( \frac{1}{\beta} \right) \leq \ln(\beta - 2)$   
 $(1-T) \left( \ln \frac{1}{\beta} - \ln \beta \right) \leq \ln(\beta - 2)$   
 $T - 1 \geq \frac{\ln(\beta - 2)}{\ln \beta}$   
 $T \geq \frac{\ln(\beta - 2)}{\ln \beta} + 1$

the smallest  $T$  is  $\frac{\ln(\beta - 2)}{\ln \beta} + 1$

10.

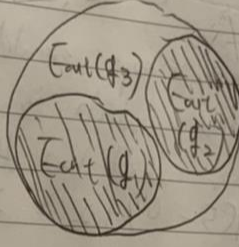
10.  
 $T = E_p[\|y - X(w \odot p)\|_2^2]$   
 $d_w(T) = E_p[\frac{\partial}{\partial w} (y - X(w \odot p)) (-Xp)] = 0$   
 $\Rightarrow E_p[(X(w \odot p) - y)(Xp)] = 0$   
 $E_p[Xw \odot Xp] - y E_p[Xp] = 0$   
 $E_p[0.5w]^T X^T X \cdot 0.5 = 0.5 y^T X$   
 $w = 2 y^T X^{-1}$

$p$  is binary vector, whose expectation is 0.5. (between 0 and 1)

11.

11. The best possible  $E_{\text{out}}(G)$  is 0. (0.08 + 0.16 + 0.32 = 0.56 < 1)

The worst possible  $E_{\text{out}}(G)$  is the coloring part below.



Because  $E_{\text{out}}(q_3) > E_{\text{out}}(q_1) + E_{\text{out}}(q_2)$

$\begin{array}{ccc} \text{||} & & \text{||} \\ 0.32 & & 0.08 + 0.16 = 0.24 \end{array}$

So, the worst possible  $E_{\text{out}}(G)$  is 0.24

$0 \leq E_{\text{out}}(G) \leq 0.24$

12.

12. If a point is an error, then there should be at least  $\frac{K+1}{2}$  wrong classifiers.

So, the upper bound of  $E_{\text{out}}(G)$  would be the sum of each error divided by  $\frac{K+1}{2}$ , which is  $\frac{\sum_{k=1}^K e_k}{\frac{K+1}{2}} = \frac{2}{K+1} \sum_{k=1}^K e_k$

13.

13.

The probability that an example is sampled at least once equals to  $N(1-p)$  (an example is not sampled at all).

$p$  (not sampled at all)

$$= \left(1 - \frac{1}{N}\right)^N = \left(1 - \frac{1}{N}\right)^{PN}$$

While  $N$  is very large,

$$\left(1 - \frac{1}{N}\right)^{PN} = \left[\left(1 - \frac{1}{N}\right)^N\right]^P \Rightarrow$$

$$= (e^{-1})^P$$

$$= e^{-P}$$

The approximately of the examples are sampled at least once is  $N(1-p)$  (not sampled at all)

$$= N - N \cdot e^{-P}$$

$$= N - (e^{-P} \cdot N)$$


---

Let  $y = \lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}\right)^N$

$$\ln y = \lim_{N \rightarrow \infty} N \ln \left(1 - \frac{1}{N}\right)$$

$$\ln y = \lim_{N \rightarrow \infty} \frac{\ln \left(1 - \frac{1}{N}\right)}{N^{-1}}$$

(with l'Hospital's rule)

$$= \lim_{N \rightarrow \infty} \frac{\left(-\frac{1}{N^2}\right)(0 + N^{-2})}{-N^{-2}}$$

$$= \lim_{N \rightarrow \infty} -\left(\frac{1}{1 - \frac{1}{N}}\right)$$

$$= \lim_{N \rightarrow \infty} -\left(\frac{N}{N-1}\right)$$

$$= -1$$

$$y = e^{-1}$$

14.

14.

For each dimension, the number of stamp should be between  $[L, R]$ , which is  $R-L$ .

So, the number of decision stamps should be  $2d(R-L)$

However, there exist the situation which is all positive and all negative, so we +2. The answer would be  $2d(R-L) + 2$ , which is  $2 \cdot 4 \cdot (5 - 0) + 2 = 42$



15.

15.

$$K_{ds}(x, x') = \text{sign}(x_1 - \theta_1) \text{sign}(x'_1 - \theta_1) + \text{sign}(x_2 - \theta_2) \text{sign}(x'_2 - \theta_2) + \dots$$

$$\dots + \text{sign}(x_{191} - \theta_{191}) \text{sign}(x'_{191} - \theta_{191})$$

if  $\theta_k$  is between  $x_k$  and  $x'_k$ ,  $\text{sign}(x_k - \theta_k) \text{sign}(x'_k - \theta_k) = -1$   
 otherwise  $= 1$

There are  $|x_k - x'_k| + 1$  integers between  $x_k$  and  $x'_k$ , so  
 there are  $\underbrace{2(|x - x'| + 1)}_{= 2|x - x'|}$  results that equals  $-1$

We know that  $|G| = 2d(R-L) + 2$ , so the number of  
 $\text{sign}(x_k - \theta_k) \text{sign}(x'_k - \theta_k) = 1$  is  $2d(R-L) + 2 - 2|x - x'|$

So,  $K_{ds}(x, x') = \underbrace{2d(R-L) + 2}_{+1} - \underbrace{2|x - x'|}_{-1}$   
 $= 2d(R-L) + 2 - 4|x - x'|$

16.

16.

With infinite input vectors,  
 $K_{ds}(x, x')$  should be  
 $2d(R-L) + 2 - 4|x - x'|$

17.

The lecture I like most is activation in deep learning. Before I studied in this class, the most part that interest me is about neural networks. For me, face recognition is a very awesome thing to me. I feel surprised to predict a value or class with only just lots of data but not logics. Although I find at last that there is lots of derivation behind the models, it is still a cool thing. The part introduced the activation in different layer, which is the point of a neural network model.

18.

The lecture I like least is dual support vector machine. In that part, there is lots of Lagrange dual problem derivation, and I think the part was only mentioned quickly. To be honest, I've never learned about that part before, so I really couldn't understand it very well.