

$$1. \quad \begin{array}{c} N/k \\ \overbrace{\textcircled{0} \textcircled{0} \textcircled{0} \dots}^k \end{array}$$

total of $\binom{k}{2}$ combinations, therefore $\binom{k}{2}$ one-to-one models.

each one-to-one takes $\geq \frac{N}{K}$ values.

therefore the training process took $\binom{k}{2} \cdot \binom{N}{K}^3 = \binom{k(k-1)}{2} \cdot \frac{8N^3}{K^3} = 4aN^3 \frac{k(k-1)}{K^3}$

$$2. \quad \phi_a(x_n) = \vec{x}_n = \begin{bmatrix} 1 \\ x_n \\ x_n^2 \\ \vdots \\ x_n^{d-1} \end{bmatrix}, \text{ let } \vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\tilde{w}^T \phi_a(x_n) = \tilde{w}^T \vec{x}_n = w_0 + w_1 x_n + w_2 x_n^2 + \dots + w_{d-1} x_n^{d-1}$$

$$\phi_a(\vec{x}) = \vec{V}$$

$$\text{for all } \{(x_n, y_n)\}_{n=1}^N, \quad g(\vec{x}) = \tilde{w}^T \phi_a(\vec{x}) = \tilde{w}^T \vec{V} = \begin{pmatrix} w_0 + w_1 x_1 + w_2 x_2 + \dots + w_{d-1} x_1^{d-1} \\ \vdots \\ w_0 + w_1 x_N + w_2 x_N^2 + \dots + w_{d-1} x_N^{d-1} \end{pmatrix}$$

Now, we look at $\vec{V} \tilde{w} = \vec{y}$, since \vec{V} is an $N \times N$ matrix,

In order for this equation to exist a solution, $\det(\vec{V}) \neq 0$, which is \vec{V} is invertible.

Since we know $\det(\vec{V}) = \prod_{1 \leq m < n \leq N} (x_m - x_n)$, and for all $a_1 \sim a_N$ are different, $\det(\vec{V}) \neq 0$

then we know $\tilde{w} = \vec{V}^{-1} \vec{y}$, there exist \tilde{w} s.t. $\vec{V} \tilde{w} = \vec{y}$

$$E_{in}(g) = \sum_{n=1}^N (w_0 + x_n w_1 + x_n^2 w_2 + \dots + x_n^{d-1} w_{d-1} - y_n)^2$$

$$\text{Since } \tilde{w} \vec{V}^{-1} = \vec{y}$$

we get

$$E_{in}(g) = \sum_{n=1}^N (y_n - y_n)^2 = \sum_{n=1}^N 0 = 0$$

$$3, \text{ let } \vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix}, \vec{y} = \begin{pmatrix} \vec{x}_1 + b_1 \\ \vec{x}_2 + b_2 \\ \vdots \\ \vec{x}_N + b_N \end{pmatrix}$$

the $\phi(\vec{x}) = \begin{pmatrix} 1 & 0 \\ 0 & 1 & \dots \\ \dots & \dots & N \end{pmatrix}_{N \times N}$ = an Identity matrix with size $N \times N = I_N$

Now perform linear regression on I_N

However, since I_N is an Identity matrix with size $N \times N$, it's invertible.

$$\text{Hence } \tilde{w} = I_N^{-1} \vec{y} = \vec{y}$$

$$\text{Therefore } E_{in} = (\phi(\vec{x}) \tilde{w} - \vec{y})(\phi(\vec{x}) \tilde{w} - \vec{y})^T = 0. \quad \times$$

Now we check E_{out} .

First, suppose we get x from the uniform distribution (but $x \neq x_1, \dots, x_N$) where $y = x + b$

$$\text{then } \phi(x) = \vec{x} = [0, 0, 0, \dots, 0]_{1 \times N}, g(x) = 0, \text{ error} = (0 - x - b)^2 = x^2 + 2xb + b^2$$

Next, we know the possibility for getting a specific value from a distribution

$$\text{is } 0, \text{ which is to say } \int_x^x \frac{1}{2} dx = 0$$

Therefore the possibility for getting $x_1 \sim x_N$ is 0.

$$\text{Therefore } E_{out} = E(x^2 + xt + t^2) = \underbrace{E(x^2) + E(xt) + E(t^2)}_{= \frac{4}{3}} = \frac{4}{3} \quad \times$$

$$\left\{ \begin{array}{l} 1. \text{Var}[x] = E[x^2] - (E[x])^2 \Rightarrow E[x^2] = \text{Var}[x] + (E[x])^2 = \frac{1}{12}(z^2) + 0^2 = \frac{1}{3} \\ 2. E[t^2] = \text{Var}[t] + (E[t])^2 = 1 \\ 3. E[xt] = E[x]E[t] = 0 \end{array} \right.$$

(since x & t are independent)

$$4. \text{ For } (X_h^T X_h)_{11} = \sum_{n=1}^N X_{1n}^2 + \sum_{n=1}^N (X_{1n} + \epsilon_n)^2 = \sum_{n=1}^N 2X_{1n}^2 + 2X_{1n}\epsilon_n + \epsilon_n^2$$

$$E[(X_h^T X_h)_{11}] = \sum_{n=1}^N 2E[X_{1n}^2] + 2E[X_{1n}\epsilon_n] + E[\epsilon_n^2]$$

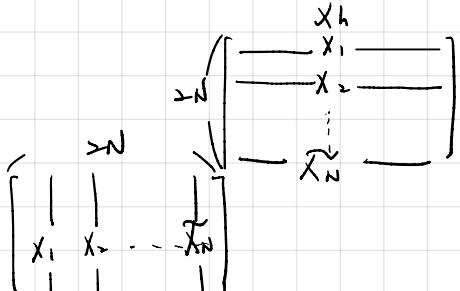
$$\epsilon_n \sim U[-\delta, \delta], \quad X_{1n} = x_{1n}$$

$$\text{Then } E[X_{1n}^2] = x_{1n}^2$$

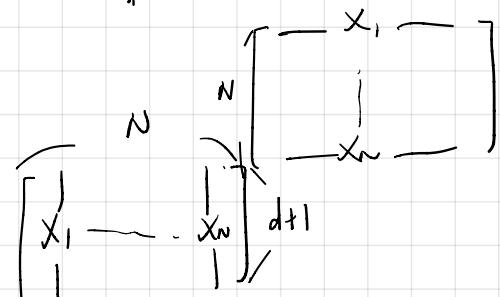
$$\left\{ \begin{array}{l} 2E[X_{1n}\epsilon_n] = 2E[X_{1n}]E[\epsilon_n] = 0 \\ (\text{since } x_{1n} \text{ and } \epsilon_n \text{ are independent}) \end{array} \right.$$

$$E[\epsilon_n^2] = \text{Var}[\epsilon_n] + (E[\epsilon_n])^2 = \frac{1}{12}(\approx \delta)^2 = \frac{1}{3}\delta^2$$

$$\text{We get } E[(X_h^T X_h)_{11}] = \sum_{n=1}^N 2x_{1n}^2 + \frac{1}{3}\delta^2 = \frac{N}{3}\delta^2 + \sum_{n=1}^N 2x_{1n}^2$$



X_h^T



X^T

$$\text{Then } E[X_h^T X_h] = X^T X + X^T X + \begin{pmatrix} \frac{N}{3}\delta^2 & & -\frac{N}{3}\delta^2 \\ & \ddots & \\ & & \frac{N}{3}\delta^2 \end{pmatrix}_{(d+1) \times (d+1)}$$

5.

$$E_{\text{aug}}(\vec{w}) = E_{\text{in}}(\vec{w}) + \frac{\lambda}{N} \vec{w}^T \vec{w}$$

$$\Rightarrow E_{\text{aug}}(\vec{w}) = E_{\text{in}}(\vec{w}) + \frac{2\lambda}{N} \vec{w}^T \vec{w}$$

$$\vec{w}_{t+1} \leftarrow \vec{w}_t - \eta \circ E_{\text{aug}}(\vec{w}_t)$$

$$\Rightarrow \vec{w}_t - (\eta \circ E_{\text{in}}(\vec{w}_t)) + \frac{2\lambda}{N} \eta \vec{w}_t$$

$$\not\equiv \left(1 - \frac{2\lambda\eta}{N}\right) \vec{w}_t - \eta \circ E_{\text{in}}(\vec{w}_t)$$

$$\not\equiv \left(1 - \frac{2\lambda\eta}{N}\right) \left(\vec{w}_t - \underbrace{\frac{n}{1 - \frac{2\lambda\eta}{N}} \circ E_{\text{in}}(\vec{w}_t)}_{\downarrow} \right)$$

$$\frac{n}{N - 2\lambda\eta}$$

$$\text{then } \alpha = 1 - \frac{2\lambda\eta}{N}, \quad \beta = \frac{N}{N - 2\lambda\eta}$$

XX

6.

$$\text{let } f(w) = \frac{1}{N} \sum_{n=1}^N (w \cdot x_n - y_n)^2$$

In order to maximize $f(w)$ under $g(w) = w^2 \leq C$, we find lagrange multiplier such that

$$\nabla f = \lambda \nabla g(w)$$

$$\Rightarrow \frac{1}{N} \sum_{n=1}^N \nabla (w \cdot x_n - y_n) (x_n) = \nabla w \lambda \Rightarrow \left(\sum_{n=1}^N w x_n \right) - N w \lambda = \sum_{n=1}^N x_n y_n \Rightarrow w^* = \frac{\sum_{n=1}^N x_n y_n}{\left(\sum_{n=1}^N x_n \right) - N \lambda}$$

Since $(w^*)^2 = C \Rightarrow \frac{\sum_{n=1}^N x_n y_n}{\left(\sum_{n=1}^N x_n \right) - N \lambda} = \sqrt{C} \Rightarrow \frac{\sum_{n=1}^N x_n y_n}{\sqrt{C}} = \left[\left(\sum_{n=1}^N x_n^2 \right) - N \lambda \right]$

$$\Rightarrow \lambda = \frac{-\sum_{n=1}^N x_n y_n}{N \sqrt{C}} + \frac{1}{N} \sum_{n=1}^N x_n^2$$

therefore $\alpha = \frac{-\sum_{n=1}^N x_n y_n}{N}$

$$\beta = \frac{1}{N} \sum_{n=1}^N x_n^2$$

**

1.

$$E(\tilde{w}) = \frac{1}{N} \sum_{n=1}^N (\tilde{w}^T V \vec{x}_n - y_n)^2 + \frac{\lambda}{N} \|\tilde{w}\|$$

$$\text{let } \tilde{w}^T V = \tilde{w}^T = [V_{11} w_1, V_{22} w_2, \dots, V_{dd} w_d]$$

$$\Rightarrow \tilde{w}^T = \tilde{w}^T V^{-1} \quad (V \text{ is invertible, since } V \text{ is diagonal matrix storing positive values})$$

Since V is a diagonal matrix, we know $V^{-1} = \begin{bmatrix} \frac{1}{V_{11}} & & & \\ & \frac{1}{V_{22}} & & 0 \\ & & \ddots & \\ 0 & & & \frac{1}{V_{dd}} \end{bmatrix}$

$$\text{we also have } \tilde{w} = \begin{bmatrix} V_{11} w_1 \\ V_{22} w_2 \\ \vdots \\ V_{dd} w_d \end{bmatrix} = V \tilde{w} \Rightarrow \tilde{w} = V^{-1} \tilde{w}$$

$$\text{then } E(\tilde{w}) = \frac{1}{N} \sum_{n=1}^N (w^T \vec{x}_n - y_n)^2 + \frac{\lambda}{N} \|V^{-1} \tilde{w}\|^2$$

$$\text{then } \Omega(\tilde{w}) = V^{-1} \tilde{w}$$

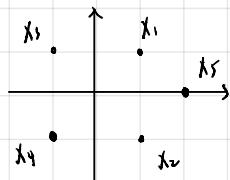
$$\text{and } \tilde{w} = V \tilde{w}$$

8.

- ① We remove one positive example, A minority will always return positive ($\because N-1 < N$) therefore $e=0$, and this will always be the same for all N positive examples.
- ② Then we remove one negative sample, A minority will always return negative ($\because N-1 < N$) therefore $e=0$, and this will always be true for all negative examples.

Hence $E_{\text{loop}}(A_{\text{minority}}) = 0$

9.



Situation $E_{\text{in}}=0 \Rightarrow x_1 \sim x_5$ different to the rest. $\Rightarrow 2^2$ combinations

x_1, x_5 same

x_2, x_5 same

x_1, x_3 same

x_3, x_4 same

x_2, x_4 same

all the same

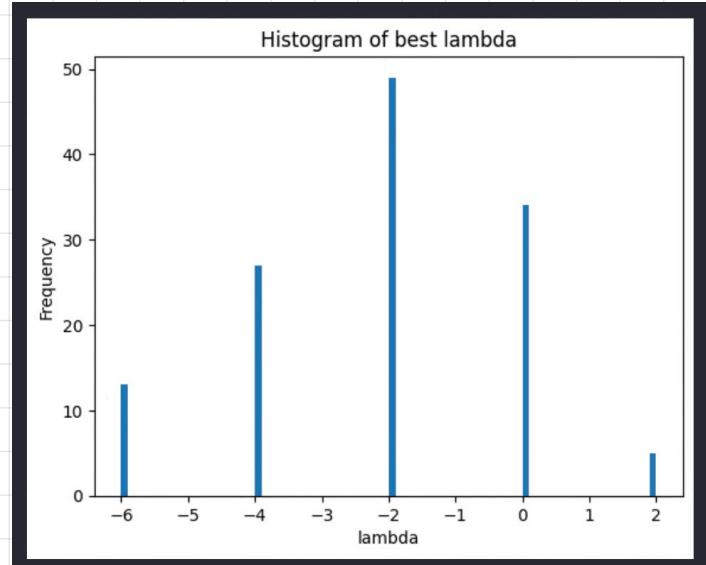
$$E_{\text{in}} = 1 \neq (x_1, x_5), (x_4, x_5), (x_2, x_3), (x_1, x_4), (x_1, x_2) \text{ same} \Rightarrow 10 \text{ combination}$$

Then the expectation $= \frac{10}{32} \times 1 + \frac{22}{32} \times 0 = \frac{5}{16}$

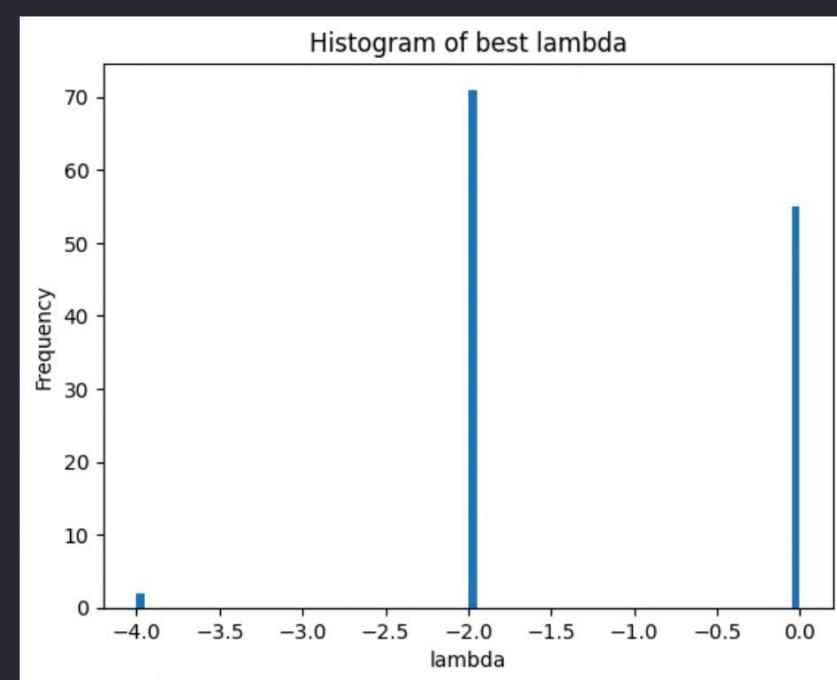
| 0 .

```
Training with lambda: 1e-06  
(200, 84) (200,  
Accuracy = 96% (192/200) (classification)  
  
Training with lambda: 0.0001  
(200, 84) (200,  
Accuracy = 92% (184/200) (classification)  
  
Training with lambda: 0.01  
(200, 84) (200,  
Accuracy = 91% (182/200) (classification)  
  
Training with lambda: 1.0  
(200, 84) (200,  
Accuracy = 87.5% (175/200) (classification)  
  
Training with lambda: 100.0  
(200, 84) (200,  
Accuracy = 80.5% (161/200) (classification)  
  
Best lambda: 1e-06  
Best lambda(log_10): -6.0  
Best accuracy: 96.0
```

| 1 .



| 2 .



From Q10, we get the best lamda is 1e-6.

From Q11, we can see that lamda value 1e-2 appeared the most.

From Q12, we see that lambda value with 1e-2 appeared the most.

Therefore we know from cross validation that, when lamda = 1e-2, the err would be minimize.

Moreover, the distribution for Q12 is not as wide as Q11, only three values are left.

The distribution for Q11 is more like normal distribution.

We know from this experiment, cross validation can limit the amount of "good" lamdas than simply separating our data set.

Through using folds to verify which lambda is good, we reduce the amount of good lambdas into 3 rather than 5.

13.

$$\min_w E_{in}(w) \text{ subject to } w^T w \leq C$$

is the same as finding Lagrange multiplier $\lambda > 0$ and \vec{w}_{reg}

$$\text{such that } \nabla E_{in}(\vec{w}_{reg}) + \frac{2\lambda}{N} w = 0 \Rightarrow -\frac{2\lambda}{N} \vec{w}_{reg} = \nabla E_{in}(\vec{w}_{reg})$$

Hence we know that \vec{w}_{reg} and $-\nabla E_{in}(\vec{w}_{reg})$ are only different with a multiplying factor.

Moreover $-\nabla E_{in}(\vec{w}_{reg})$ point towards w_{L2N} , if and only if $X^T X$ is reversible

$$\text{which is to say } X^T X = \alpha I$$

$$\text{therefore } \vec{w}_{reg} + (-\nabla E_{in}(\vec{w}_{reg})) t = \vec{w}_{L2N} \Rightarrow n \vec{w}_{reg} = \vec{w}_{L2N}$$

$$\text{and this is the form of } \vec{w}_u = \vec{w}_{L2N} \times \frac{\sqrt{C}}{\|w_{L2N}\|} \text{ when we set}$$

$$n = \frac{\sqrt{C}}{\|w_{L2N}\|}$$

Hence we prove w_u solves the L_2 -constrained regression problem $\min_w E_{in}(w)$ subject to $w^T w \leq C$ i.f.f. $X^T X = \alpha I$