

1. In OVD, each classifier is trained on data from 2 classes.

For K-class classification, there are $\binom{K}{2} = \frac{K(K-1)}{2}$

Each binary classifier is trained on $\frac{N}{K}$ binary classifier

Given for binary classification, we need CPU time $\propto N^3$.

\Rightarrow For K-class classification, each binary classification has data size $2\left(\frac{N}{K}\right) \Rightarrow \text{CPU time} = \alpha \left(2\left(\frac{N}{K}\right)\right)^3$

$$\therefore \text{Total CPU time} = \frac{K(K-1)}{2} \cdot \alpha \left(2\left(\frac{N}{K}\right)\right)^3$$

$$= 4\alpha N^3 \frac{K(K-1)}{K^3}$$

$$= \frac{4\alpha N^3 (K-1)}{K^2}$$

#

$$\{x_n, y_n\}, g(x) = \tilde{w}^T \Phi_Q(x), z_n = \Phi_Q(x_n)$$

run linear regression on $\{(z_n, y_n)\}_{n=1}^N$ find \tilde{w} , prove there's

some Q s.t. $E_{\text{in}} = 0$

$\because X_n$ is 1-D data \Rightarrow after Φ_Q , we can simply get

$$z_n = (1, x_n, x_n^2, x_n^3, \dots, x_n^Q)$$

(Q+1 dimension data)

For linear regression by squared error,

$$E_{in}(g) = \frac{1}{N} \sum_{n=1}^N (\tilde{w}^T \tilde{x}_n - y_n)^2$$

$$= \frac{1}{N} \sum_{n=1}^N (z_n^T \tilde{w} - y_n)^2$$

$$= \frac{1}{N} \left\| \begin{bmatrix} z_1^T \tilde{w} - y_1 \\ z_2^T \tilde{w} - y_2 \\ \vdots \\ z_N^T \tilde{w} - y_N \end{bmatrix} \right\|^2$$

$$= \frac{1}{N} \left\| \begin{bmatrix} -z_1^T \\ -z_2^T \\ \vdots \\ -z_N^T \end{bmatrix} \tilde{w} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \right\|^2$$

$(N \times (Q+1)) \quad (Q+1) \times 1 \quad N \times 1$

$$= \frac{1}{N} \| Z w - y \|^2$$

$$= \frac{1}{N} (\tilde{w}^T Z^T Z \tilde{w} - 2 \tilde{w}^T Z^T y + y^T y)$$

To find min,

$$\nabla E_{in}(\tilde{w}_{Lin}) = 0 \Rightarrow \nabla E_{in}(w) = \frac{1}{N} (2 Z^T Z w - 2 Z^T y) = 0$$

$$\therefore \text{for } w_{in} = (Z^T Z)^{-1} Z^T y$$

We know the dimension of Z is $(N \times (Q+1)) \Rightarrow$ We pick

' $Q = N-1$ ' which can form a $N \times N$ matrix. (Vandermonde Matrix)

To detect if $(Z^T Z)$ is invertible, we check $\det(Z^T Z)$

$$\therefore \det(Z^T Z) = \det(Z^T) \det(Z)$$

$$= (\det(Z))^2 \quad (\because \det(Z^T) = \det(Z))$$

Vandermonde matrix

$$\Downarrow = \left(\begin{array}{cccc} 1 & x_1 & x_1^2 & \dots & x_1^{N-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{N-1} \\ & & \vdots & & \\ 1 & x_N & x_N^2 & \dots & x_N^{N-1} \end{array} \right)^2$$

$$= \left(\prod_{1 \leq n < m \leq N} (x_m - x_n) \right)^2$$

" Given that Assume all $\{x_n\}_{n=1}^N$ are different

$$\therefore \det(Z^T Z) = \left(\prod_{1 \leq n < m \leq N} (x_m - x_n) \right)^2 \neq 0 \quad \therefore (Z^T Z)^{-1} \text{ exists}$$

" We can find unique $W_{lin} = (Z^T Z)^{-1} Z^T y$ to form

$$g(x) = \tilde{W}_{lin}^T \tilde{\Phi}(x) \text{ such that } E_{in}(g) = \frac{1}{N} \sum_{n=1}^N (W_{lin}^T z_n - y_n) = 0$$

as we pick $Q = N-1$ perfectly $\#$.

3. x is uniformly sampled from $[-1, 1]$

$y = x + \varepsilon$, ε is independently sampled from Gaussian distribution with mean 0 and Variance 1.

$$\therefore (\Phi(x))_n = z_n = \begin{bmatrix} 1 \\ x - x_n \end{bmatrix}$$

" all x_n are different in training data set, z will have exact one '1' and rest will be '0'.

∴ By squared error,

$$E_{in}(g) = \frac{1}{N} \sum_{n=1}^N (\tilde{w}^T \tilde{z}_n - y_n)^2$$

∴ For $E_{in}(g)$,

Since the transform creates a one-hot encoded vector for each training sample, the linear regression model will perfectly predict output y for each training sample x .

$$\therefore E_{in}(g) = 0$$

For $E_{out}(g)$,

For any new sample x' , $\Phi(x')$ will be a vector of zeros,

(∵ x' is different from all training samples x_n)

$$\therefore g(x') = \tilde{w}^T \Phi(x') = 0 \quad (\text{for all } x' \text{ not in the training set.})$$

By $y' = x' + \epsilon$, ∵ $g(x') = 0 \Rightarrow$ The squared error for new sample is

$$(y' - g(x'))^2 = (x' + \epsilon)^2$$

$$\therefore E_{out}(g) = E((x' + \epsilon)^2)$$

$$= E(x'^2) + E(2x'\epsilon) + E(\epsilon^2)$$

$E(x')$ is uniform distribution over $[-1, 1] \Rightarrow E(x'^2) = \frac{1}{12} (1 - (-1)^2) = \frac{1}{3}$

$E(x'\epsilon) = E(x')E(\epsilon) = 0$ (∵ mean of ϵ is 0)

$E(\epsilon^2) = 1$ (variance of ϵ is 1) ∴ $E_{out}(g) = 1 + \frac{1}{3} = \frac{4}{3}$

$$4. \quad X = \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \dots & x_N \\ | & | & & | \end{bmatrix}^T = \begin{bmatrix} \text{---} x_1 \text{---} \\ \text{---} x_2 \text{---} \\ \vdots \\ \text{---} x_N \text{---} \end{bmatrix}_{N \times (d+1)}$$

$$X_h = \begin{bmatrix} | & | & & | & | & | & & | \\ x_1 & x_2 & \dots & x_N & \tilde{x}_1 & \tilde{x}_2 & \dots & \tilde{x}_N \\ | & | & & | & | & | & & | \end{bmatrix}^T = \begin{bmatrix} \text{---} x_1 \text{---} \\ \text{---} x_2 \text{---} \\ \vdots \\ \text{---} x_N \text{---} \\ \text{---} \tilde{x}_1 \text{---} \\ \text{---} \tilde{x}_2 \text{---} \\ \vdots \\ \text{---} \tilde{x}_N \text{---} \end{bmatrix}_{(2N) \times (d+1)}$$

For

$$X_h^T X_h = \left(\begin{bmatrix} | & | & & | & | & | & & | \\ x_1 & x_2 & \dots & x_N & \tilde{x}_1 & \tilde{x}_2 & \dots & \tilde{x}_N \\ | & | & & | & | & | & & | \end{bmatrix} \begin{bmatrix} \text{---} x_1 \text{---} \\ \text{---} x_2 \text{---} \\ \vdots \\ \text{---} x_N \text{---} \\ \text{---} \tilde{x}_1 \text{---} \\ \text{---} \tilde{x}_2 \text{---} \\ \vdots \\ \text{---} \tilde{x}_N \text{---} \end{bmatrix} \right)_{(d+1) \times (d+1)}$$

$$= \begin{bmatrix} | \\ x_1 \\ | \end{bmatrix} \begin{bmatrix} \text{---} x_1 \text{---} \end{bmatrix} + \begin{bmatrix} | \\ x_2 \\ | \end{bmatrix} \begin{bmatrix} \text{---} x_2 \text{---} \end{bmatrix} + \dots + \begin{bmatrix} | \\ \tilde{x}_N \\ | \end{bmatrix} \begin{bmatrix} \text{---} \tilde{x}_N \text{---} \end{bmatrix}$$

$$= X^T X + \begin{bmatrix} | \\ \tilde{x}_1 \\ | \end{bmatrix} \begin{bmatrix} \text{---} \tilde{x}_1 \text{---} \end{bmatrix} + \begin{bmatrix} | \\ \tilde{x}_2 \\ | \end{bmatrix} \begin{bmatrix} \text{---} \tilde{x}_2 \text{---} \end{bmatrix} + \dots + \begin{bmatrix} | \\ \tilde{x}_N \\ | \end{bmatrix} \begin{bmatrix} \text{---} \tilde{x}_N \text{---} \end{bmatrix}$$

$$= \underbrace{X^T X}_{\textcircled{1}} + \underbrace{\tilde{X}^T \tilde{X}}_{\textcircled{2}}$$

For ① $\Rightarrow E(X^T X) = X^T X$

For ② $\Rightarrow E(\tilde{X}^T \tilde{X}) = E((X + \epsilon)^T (X + \epsilon))$

$$= E(X^T X + \epsilon^T X + X^T \epsilon + \epsilon^T \epsilon)$$

For $E(X^T X + \varepsilon^T X + X^T \varepsilon + \varepsilon^T \varepsilon)$

p6

$$= X^T X + 0 + 0 + E(\varepsilon^T \varepsilon) \left(\because E(\varepsilon^T X) = E(X^T \varepsilon) = 0 \right)$$

For $E(\varepsilon^T \varepsilon)$,

' ε has components uniformly distributed in $[-\delta, \delta]$, its

variance $\text{Var}(\varepsilon_i) = \frac{1}{12} (\delta - (-\delta))^2 = \frac{\delta^2}{3}$ for each i

$$\therefore E(\varepsilon^T \varepsilon) = \frac{\delta^2}{3} I_{(d+1) \times (d+1)}$$

In conclusion,

$$\therefore E(X_h^T X_h) = E(X^T X) + E(\tilde{X}^T \tilde{X})$$

$$= X^T X + \left(X^T X + \frac{\delta^2}{3} I \right)$$

$$= 2X^T X + \frac{\delta^2}{3} I \quad \# \quad \left(I \text{ is identity matrix with size } (d+1) \times (d+1) \right)$$

5.

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

$$\nabla (E_{\text{aug}}(w)) = \nabla (E_{\text{in}}(w)) + \frac{\partial}{\partial w} \left(\frac{\lambda}{N} w^T w \right)$$

$$= \nabla E_{\text{in}}(w) + \frac{2\lambda}{N} w$$

By gradient descent algorithm

$$w_{t+1} \leftarrow w_t - \eta \nabla E_{\text{aug}}(w_t)$$

$$\Rightarrow w_{t+1} \leftarrow w_t - \eta \left(\frac{2\lambda}{N} w_t + \nabla E_{\text{in}}(w_t) \right)$$

$$\Rightarrow w_{t+1} \leftarrow \left(1 - \frac{2\lambda\eta}{N} \right) w_t - \eta \nabla E_{\text{in}}(w_t)$$

$$\Rightarrow w_{t+1} \leftarrow \left(1 - \frac{2\lambda\eta}{N} \right) \left(w_t - \frac{\eta N}{N - 2\lambda\eta} \nabla E_{\text{in}}(w_t) \right)$$

$$\therefore \alpha = 1 - \frac{2\lambda\eta}{N}, \beta = \frac{\eta N}{N - 2\lambda\eta}$$

p7

6. Find w^* by take gradient to $\min_{w \in \mathbb{R}} \frac{1}{N} \sum_{n=1}^N (w x_n - y_n)^2 + \frac{\lambda}{N} w^2$ #

$$E_{\text{avg}}(w) = \frac{1}{N} \sum_{n=1}^N (w^2 x_n^2 - 2w x_n y_n + y_n^2) + \frac{\lambda}{N} w^2$$

$$\text{Let } \nabla E_{\text{avg}}(w^*) = 0,$$

(Given X is 1-D data)

$$\nabla E_{\text{avg}}(w^*) = \frac{1}{N} \sum_{n=1}^N (2x_n^2 w^* - 2x_n y_n) + \frac{2\lambda}{N} w^* = 0$$

$$\Rightarrow \frac{2}{N} \sum_{n=1}^N (x_n^2 w^* - x_n y_n) + \frac{2\lambda}{N} w^* = 0$$

$$\Rightarrow w^* \left(\frac{2}{N} \sum_{n=1}^N x_n^2 + \frac{2\lambda}{N} \right) - \frac{2}{N} \sum_{n=1}^N x_n y_n = 0$$

$$\Rightarrow w^* = \frac{\sum_{n=1}^N y_n x_n}{\sum_{n=1}^N x_n^2 + \lambda}$$

① We check $(w^*)^2 = c \Rightarrow w^* = \pm \sqrt{c}$ (for every $\lambda > 0$)

if $w^* = \sqrt{c}$,

$$\frac{\sum_{n=1}^N y_n x_n}{\sum_{n=1}^N x_n^2 + \lambda} = \sqrt{c} \Rightarrow \sqrt{c} \left(\sum_{n=1}^N x_n^2 + \lambda \right) = \sum_{n=1}^N y_n x_n$$

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

② if $w^* = -\sqrt{c}$

$$\frac{\sum_{n=1}^N y_n x_n}{\sum_{n=1}^N x_n^2 + \lambda} = -\sqrt{c} \Rightarrow -\sqrt{c} \left(\sum_{n=1}^N x_n^2 + \lambda \right) = \sum_{n=1}^N y_n x_n \quad (\text{when } w^* = \sqrt{c}) \quad \#$$

$$\therefore (\alpha, \beta) = \left(\frac{\sum_{n=1}^N y_n x_n}{\sum_{n=1}^N x_n^2 + \lambda}, -\frac{\sum_{n=1}^N x_n^2}{\sum_{n=1}^N x_n^2 + \lambda} \right)$$

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

p8

$$\therefore (\alpha, \beta) = \left(-\sum_{n=1}^N y_n x_n, \sum_{n=1}^N x_n^2 \right)$$

$$7. \min_{\tilde{w} \in \mathbb{R}^{d+1}} \frac{1}{N} \sum_{n=1}^N (\tilde{w}^T \Phi(x_n) - y_n)^2 + \frac{\lambda}{N} \|\tilde{w}\|_1 \quad (\text{when } w^* = -\sqrt{C}) \quad \#$$

'V' is diagonal matrix which will scaling \hat{x} -th term in 'X'

$$\frac{1}{N} \sum_{n=1}^N (\tilde{w}^T V x_n - y_n)^2 + \frac{\lambda}{N} \|\tilde{w}\|_1 \rightarrow \text{which means } (|w_1| + |w_2| + \dots + |w_d|)$$

Comparing to regularized linear regression on original data

$$\frac{1}{N} \sum_{n=1}^N (w^T x_n - y_n)^2 + \frac{\lambda}{N} \mathcal{R}(w)$$

We observe that $w^T = \tilde{w}^T V$, Suppose V is a diagonal matrix which can be written as $V = \text{diag}(v_1, v_2, \dots, v_d, \dots, v_{d+1})$

$$\therefore \mathcal{R}(w) = \sum_{\hat{x}=1}^{d+1} \left| \frac{w_{\hat{x}}}{v_{\hat{x}}} \right| \quad \#$$

Furthermore, the relation between w and \tilde{w} , we have

$$w^T = \tilde{w}^T V \Rightarrow (w^T)^T = (\tilde{w}^T V)^T$$

$$\therefore w = V^T \tilde{w} \quad \#$$

8.

Divided into '2' cases,

Case 1. When positive example is left out.

∴ Now, we have $N-1$ positive and N negative examples

∴ For A_{minority} , we always predict the positive because we have N negative examples (majority), which is more than $N-1$ positive example (minority) (Given that A_{minority} always predict minority.)

$$E_1 = \frac{N}{2N-1}$$

Case 2. When negative example is left out

∴ Now, we have $N-1$ negative, N positive examples,

Similarly, A_{minority} will always predict negative

$$E_2 = \frac{N}{2N-1} \quad \left(\because N-1 < N \right)$$

(minority for negative) (majority for positive)

$$\begin{aligned} \therefore E_{\text{loss}}(A_{\text{minority}}) &= \frac{1}{2N} \left(N \cdot \frac{N}{2N-1} + N \cdot \frac{N}{2N-1} \right) \\ &= \frac{N}{2N-1} \end{aligned}$$

case 1 case 2.

Consider expectation of binomial distribution

$$\therefore E = \frac{1}{32} (n \cdot p) \Big|_{(n=5, p=\frac{1}{2})}$$

$$= \frac{1}{32} \cdot 5 \cdot \frac{1}{2} = \frac{5}{64}$$

13.

Unconstrained linear regression solution is

$$W_{lin} = (X^T X)^{-1} X^T y$$

By method of Dr. Regularize, subject to $\|w\|^2 \leq C$

$$W_C = \frac{W_{lin}}{\|W_{lin}\|} \cdot \sqrt{C}$$

Assume $X^T X = \alpha I$,

$$W_{lin} = (X^T X)^{-1} X^T y = \frac{1}{\alpha} X^T y$$

(ridge problem)

Now, we consider the regularized problem, which minimize

$\|Xw - y\|^2 + \lambda \|w\|^2$, The solution is

$$W_{ridge} = (X^T X + \lambda I)^{-1} X^T y$$

With $X^T X = \alpha I$, W_{ridge} will be $(\alpha I + \lambda I)^{-1} X^T y$

$$= \frac{1}{\alpha + \lambda} X^T y$$

When $\|W_{lin}\|^2 > C \Rightarrow \left\| \left(\frac{1}{\alpha} \right)^2 (X^T y)^T (X^T y) \right\| > C$,

By equating $\frac{W_{lin}}{\|W_{lin}\|} \sqrt{C} = W_{ridge} \Rightarrow \left\| \frac{\sqrt{C}}{\alpha \|W_{lin}\|} X^T y \right\| = \frac{1}{\alpha + \lambda} X^T y \Rightarrow \frac{\sqrt{C}}{\alpha \|W_{lin}\|} = \frac{1}{\alpha + \lambda}$

$\Rightarrow \lambda = \frac{\alpha \|W_{lin}\|}{\sqrt{C}} - \alpha$ \therefore We still can find λ to solve C-constrained

linear regression if $X^T X = \alpha I$.

P11

For 'Only if' part

prove: if scaling W_{lin} is equivalent to solving the C-constrained problem, then $X^T X = \alpha I$

By equating $\frac{W_{lin}}{\|W_{lin}\|} \sqrt{C} = W_{ridge}$

$$\frac{W_{lin}}{\|W_{lin}\|} \cdot \sqrt{C} = (X^T X + \lambda I)^{-1} X^T y$$

Substituting $W_{lin} = (X^T X)^{-1} X^T y$, we have

$$\frac{(X^T X)^{-1} X^T y}{\|(X^T X)^{-1} X^T y\|} \sqrt{C} = (X^T X + \lambda I)^{-1} X^T y$$

$$\Rightarrow \left(\frac{X^T X}{\|(X^T X)^{-1} X^T y\|} \sqrt{C} \right)^{-1} X^T y = (X^T X + \lambda I)^{-1} X^T y$$

$$\Rightarrow \underbrace{\left(\frac{\sqrt{C}}{\|(X^T X)^{-1} X^T y\|} X^T X \right)^{-1} X^T y}_{(\text{constant})} = (X^T X + \lambda I)^{-1} X^T y$$

If we want to equate above, we must let $X^T X = \alpha I$

$$\therefore \frac{\sqrt{C}}{\|(X^T X)^{-1} X^T y\|} \alpha I = \alpha I + \lambda I \Rightarrow \frac{\alpha \sqrt{C}}{\|W_{lin}\|} = \alpha + \lambda$$

$$\Rightarrow \lambda = \frac{\alpha \sqrt{C}}{\|W_{lin}\|} - \alpha$$

\therefore we can find such λ to equating $\frac{W_{lin}}{\|W_{lin}\|} \sqrt{C} = W_{ridge}$

$\therefore W_{lin}$ is equivalent to solving the C-constrained only if $(X^T X) = \alpha I$ #

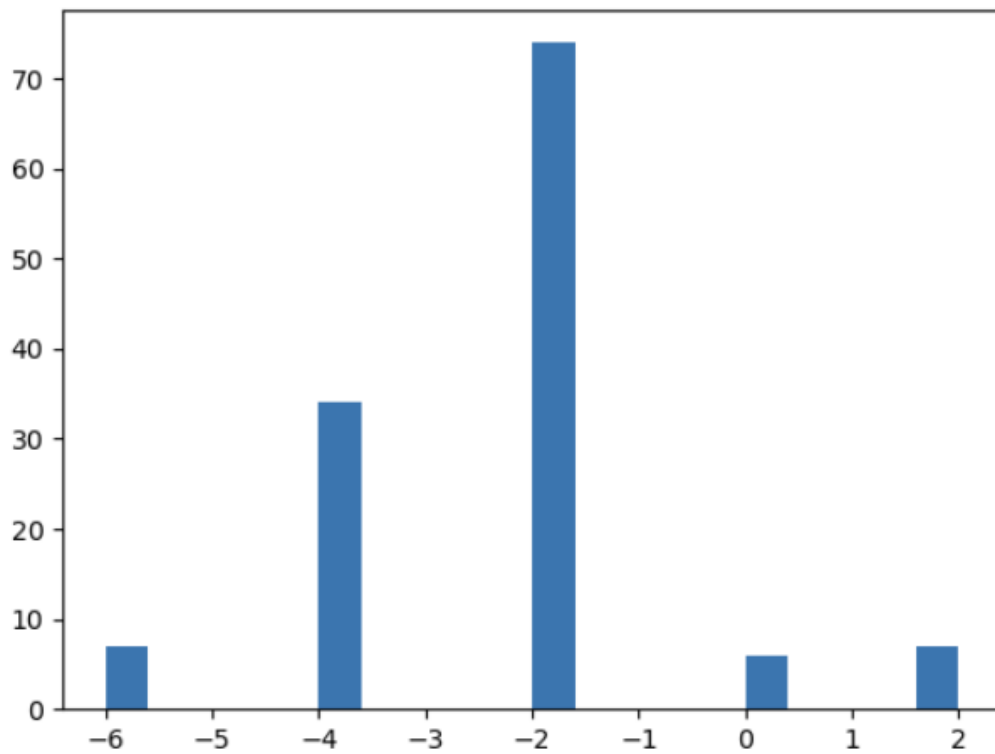
HTML hw4 solution

姓名：謝銘倫, 系級：電機三, 學號：B10502166

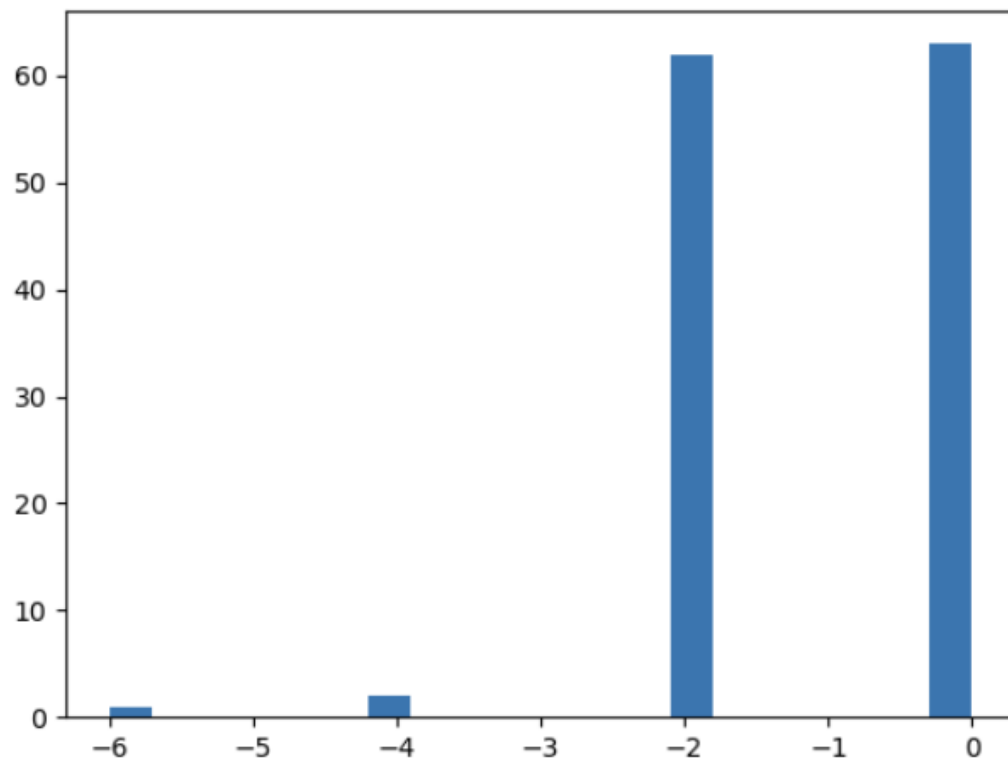
10.

ans: $\log_{10}(\lambda^*) = -4$

11.



12.



Findings compare with p11:

1. In problem 11, we have $\log_{10}(\lambda^*)$ more on -2 and -4
2. In problem 11, we have $\log_{10}(\lambda^*)$ more on 0 and -2
3. In problem 11, we use simple validation data, which is vulnerable to the only one training data set. if the $\log_{10}(\lambda^*)$ generate from training data is not well enough, it will cause bad E_{in} in validation data.

On the othre hand, the v-fold using different training data set with each fold, which is more safer to get $\log_{10}(\lambda^*)$, which will more authentically show the real $\log_{10}(\lambda^*)$ for data set.