

P₁ P₂ P₃ P₄

Q1. D = { (0, 1, 1), (0, 2, 1), (4, 0, -1), (5, 0, -1) }

Initial seed points $\Rightarrow C_1 = (0, 0)$

$$C_2 = (3, 0)$$

Distance of P₁ from C₁: $\sqrt{(0-0)^2 + (1-0)^2} = 1 \leftarrow$
 from C₂: $\sqrt{(0-3)^2 + (1-0)^2} = \sqrt{10}$

P₂ from C₁: $\sqrt{(0-0)^2 + (2-0)^2} = 2 \leftarrow$
 from C₂: $\sqrt{(0-3)^2 + (2-0)^2} = \sqrt{13}$

2

P₃ from C₁: 4

P₃ from C₂: 1

P₄ from C₁: 5

from C₂: 2 \leftarrow

b) Pts. closest to C₁: {P₁, P₂}

New value of C₁ after update: $\frac{(0, 1) + (0, 2)}{2} = (0, 1.5)$

Pts. closest to C₂: {P₃, P₄}

Updated value of C₂: $\frac{(4, 0) + (5, 0)}{2} = (4.5, 0)$ 1

c) Purity of clustering:

Cluster 1: # points in class $y=1$: 2

points in class $y=-1$: 0

total point / size of cluster = 2

0.5

$$P_{1,1} = \frac{2}{2} = 1$$

$$P_{1,-1} = 0$$

purity of cluster 1 = $\max(P_{yy}, P_{i,-i}) = 1.$

Cluster 2:

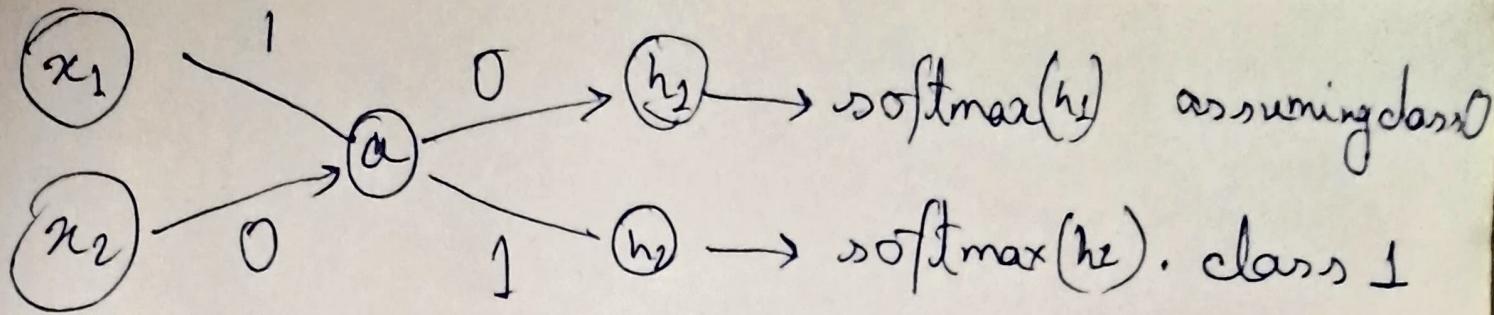
$$\begin{aligned} & \# \text{pt. in} \\ & \# \text{pt. in} \\ & 0.5 \end{aligned}$$

class $y=1 : 0$

class $y=-1 : 2$

Similarly purity of cluster 2 = 1

$$\begin{aligned} \text{Purity of clustering} &= \frac{1}{2} \times (\text{purity}_{\text{cluster } 1} + \text{purity}_{\text{cluster } 2}) \\ &= \frac{1}{2} \times (1+1) = 1. \quad 0.5 \end{aligned}$$



$$h_1 = (1 \cdot x_1 + 0 \cdot x_2) + 0 = 0$$

$$h_2 = (1 \cdot x_1 + 0 \cdot x_2) + 1 = x_1$$

for data \$(1, 1, 1)\$, \$x_1 = 1 - x_0 = 1\$, \$y = 1\$

$$\text{softmax}(h_1) = \frac{e^0}{e^0 + e^{x_1}} = \frac{1}{1+e} = 0.27$$

$$\text{softmax}(h_2) = \frac{e^{x_1}}{e^0 + e^{x_1}} = \frac{e}{1+e} = 0.73$$

a) Cross entropy for multiclass \rightarrow

$$= - \sum_{i=0}^C y_i \log(p_i) \quad y_i \text{ is ground truth}$$

here for \$[0, 1]\$ $L = - \sum_{i=0}^1 y_i \log(p_i)$

$$L = -0 \cdot \log(0.27) + 1 \cdot \log(0.73)$$

$$L = 0.13$$

c.) $a = \theta_0 x_1 + \theta_1 x_2$
 $= x_1$

So far batch normalization for 2 data points $(1, 1, 0), (-1, -1, 0)$ we calculate the

$$\text{mean } \mu = \frac{1-1}{2} = 0 \quad \textcircled{1} \quad \textcircled{2}$$

$$\text{variance} = \frac{(1-0)^2 + (-1-0)^2}{2-1} = 2 \quad \textcircled{1}$$

d.) MSE loss for $(1, 1, 0)$

$$h_1 = 0, \quad h_2 = x_1 = 1 \quad \textcircled{0.5}$$

$$\begin{aligned} \text{MSE loss} &= (h_1 - x_1)^2 + (h_2 - x_2)^2 \\ &= (0 - 1)^2 + (1 - 1)^2 \\ &= 1 \end{aligned} \quad \textcircled{1}$$

Sol^a ① We have,

$$L(y, F(x)) = e^{0.5(y - F(x))^2}$$

Differentiating w.r.t $F(x)$ we get-

$$\begin{aligned} \Rightarrow \frac{\partial L(y, F(x))}{\partial F(x)} &= - \left(e^{0.5(y - F(x))^2} \cdot 0.5 \cdot 2(y - F(x)) \right) \\ &= - \left(e^{0.5(y - F(x))^2} \cdot (y - F(x)) \right) [1.5] \\ &= - L(y, F(x)) \cdot r \end{aligned}$$

$$[0.25]$$

$$\text{or } - \frac{\partial L(y, F(x))}{\partial F(x)} = L(y, F(x)) \cdot r$$

where $r = y - F(x)$ is the residue.

In gradient boosting, the updated labels after each iteration are given by the negative gradient of the loss.

So, the updated labels are $L(y, F(x)) \cdot r$ [0.25]

Solⁿ ④ For a single sample x_i , we have the following pdf

$$p(x_i | \theta) = e^{-\theta x_i}$$

(a) ∵ Joint likelihood of observing n iid samples with the above pdf is given by

$$\begin{aligned} L(\theta) &= p(x_1, x_2, \dots, x_n | \theta) \\ &= \prod_{i=1}^n e^{-\theta x_i} = e^{-\theta \sum_{i=1}^n x_i} \end{aligned} \quad [1]$$

$$\Rightarrow \log L(\theta) = \log(e^{-\theta \sum_{i=1}^n x_i}) = -\theta \sum_{i=1}^n x_i \equiv l(\theta)$$

(b) To find MLE of θ , we need to differentiate $L(\theta)$ w.r.t θ and set it to 0.

$$\frac{\partial L(\theta)}{\partial \theta} = e^{-\theta \sum_{i=1}^n x_i} \left(-\sum_{i=1}^n x_i \right) = 0$$

We can see that this way we are not able to solve for θ . [1]

Therefore, to obtain the θ , that maximizes the likelihood we can apply gradient ascent here.

$$\theta = \theta + \eta \frac{\partial L(\theta)}{\partial \theta} \quad [1]$$

Q8) To make the hidden node maximally active, we determine the input x that maximizes the output σ of hidden node.

$h = \frac{1}{1 + e^{-z}}$ is a monotonically increasing function, so, we'd maximise z .

$z = w \cdot x$, to maximise z we make the angle b/w w & x as 0

So, $x = cw$, i.e x is scalar multiple of w .

To make the hidden node maximally active, we keep $\|x\|=1$ and optimise the dot product.

one can also use Lagrangian method

$$c = \frac{1}{\|w\|} x - 1 = \frac{w \cdot x}{\|w\|}$$

$$x = \frac{w}{\|w\|} = \frac{w}{\sqrt{w^T w}}$$

Q9) To fit a linear regression model b/w two points, we can ~~approximately~~ fit a line that minimises the error.

$$\text{error} = \sum (y - \hat{y})^2$$

$$= \sum [y - (wx + b)]^2$$

$$= (w+b)^2 + (2w+b-1)^2$$

$$= 5w^2 + 2b^2 + 6wb - 2b - 2w + 1$$

(e) To fit a linear regression model w/o two points, we can fit a line that minimises the error.

$$\text{error} = \sum (y - \hat{y})^2 \quad (\text{square})$$

$$= \sum [y - (w\mathbf{x} + b)]^2$$

$$= (w+b)^2 + (2w+b-1)^2$$

$$= 5w^2 + 2b^2 + 6wb - 2b - 2w + 1$$

$$\frac{\partial(\text{error})}{\partial w} = 0 \quad \text{and} \quad \frac{\partial(\text{error})}{\partial b} = 0$$

$$10w + 6b - 2 = 0$$

$$4b + 6w - 2 = 0$$

$$w=1, b=-1$$

$$\text{use } \mathbf{w} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$$

[we will add a dimension to \mathbf{x} with value 1 to get b in \mathbf{w}]

$$\mathbf{w} = \left(\begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\mathbf{w} = \left(\begin{bmatrix} 5 & 3 \\ 3 & 2 \end{bmatrix} \right)^{-1} \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \frac{1}{1} \begin{bmatrix} 2 & -3 \\ -3 & 5 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\text{so, } \boxed{\mathbf{w} = 1, b = -1}$$

(b) Compute w & b using D1

$$w = 1, b = -1$$

Compute w & b using D2

$$w = \left(\begin{bmatrix} -1 & -2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ -2 & 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} -1 & -2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$= \left(\begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 2 & 3 \\ 3 & 5 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$= \begin{bmatrix} -1 \\ -2 \end{bmatrix}$$

$$w = -1, b = -2$$

.5

$$E(\hat{f}(x)) = E(\hat{f}_{D1}(x)) + E(\hat{f}_{D2}(x))$$

$$w_{avg} = (1-1)/2$$

$$b_{avg} = (-1+2)/2$$

$$E(\hat{f}) = w_{avg}x + b_{avg}$$

$$\text{bias} = E(\hat{f}) - f$$

$$= \frac{E(-3)}{2} = -\frac{3}{2}$$

$$\text{bias} = E(\hat{f}(x)) - f(x)$$

$$= -\frac{3}{2} - x$$