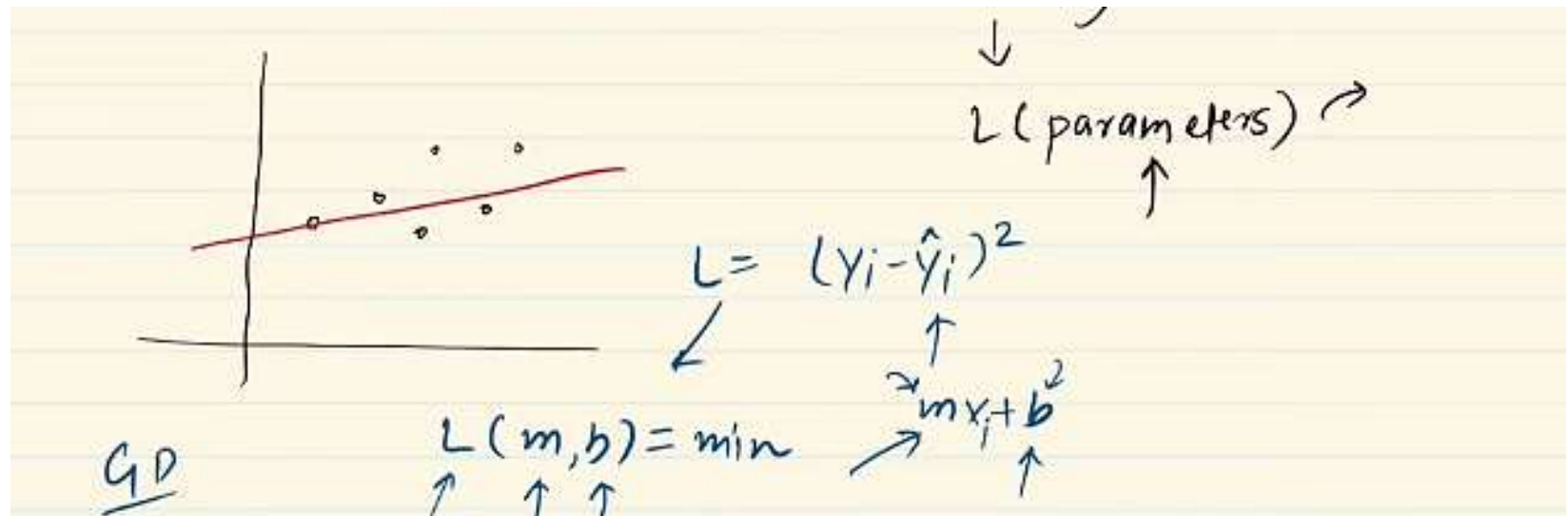Loss function is a method of evaluating how well your algorithm is modelling your dataset.

- High means Poor Model
- Low means great Model

$\downarrow$

$L\,(parameters)\;\nearrow$

$\uparrow$

$$L = (y_i - \hat{y}_i)^2$$

$\nwarrow$

$\uparrow$

$L(m,b) = min \qquad \nearrow mx_i + b^2$

GD

$\uparrow \quad \uparrow \quad \uparrow \qquad\qquad \uparrow$

Why is Loss function important?

[You can't improve what you can't measure.]
Peter Drucker

Loss Function in Deep Learning?    Backprop

cgpa | iq | package (lpa)

| cgpa | iq | package(lpa) |
|------|------|------|
| 7.1 | 83 | 3.2 |
| 8.5 | 91 | 4.5 |
| 6.3 | 102 | 6.1 |
| 5.1 | 87 | 2.7 |

# Loss Function in Deep Learning?

cgpa | iq | package (lpa)

| cgpa | iq | package (lpa) |
|------|-----|---------------|
| 7.1  | 83  | 3.2           |
| 8.5  | 91  | 4.5           |
| 6.3  | 102 | 6.1           |
| 5.1  | 87  | 2.7           |

**Backprop**

random w,b

7.1 cgpa

83 iq

$w$

$b_{11}$

$b_{12}$

$b_3$

$b_{21}$

$\hat{y}_i$ (3.7)

forward prop

Loss Function in Deep Learning?

Backprop

cgpa | iq | package (lpa)

| 7.1 | 83 | 3.2 |
|-----|-----|-----|
| 8.5 | 91 | 4.5 |
| 6.3 | 102 | 6.1 |
| 5.1 | 87 | 2.7 |

7.1 cgpa

83 iq

random
w, b

w

$b_{11}$

$b_{12}$

$b_3$

w, b

$b_{21}$

$\hat{y}_i$

(3.7)

forward
prop

$L = (y_i - \hat{y}_i)^2$

$(3.2 - 3.7)^2$

Loss Function in Deep Learning?

Backprop

random w,b → gradient descent

✓ cgpa | iq | package (lpa)

| 7.1 | 83 | 3.2 |
|-----|-----|-----|
| 8.5 | 91 | 4.5 |
| 6.3 | 102 | 6.1 |
| 5.1 | 87 | 2.7 |
| .. | .. | .. |
| .. | .. | .. |
| .. | .. | .. |

7.1 cgpa

83 iq

w,b

$\hat{y}_i$ (3.7)

→ forward prop

→ $L = (y_i - \hat{y}_i)^2$

(3.2 - 3.7)²

4.3

# Loss functions in DL

**Object detection**
focal loss

**Regression**
- mse
- mae
- huber loss

**Classification**
- binary cross entropy
- categorical cross entro
- hinge loss

**Autoencoders**
- KL divergence

**GAN**
- discriminator loss
- minmax gan loss

**Embedding**
Triplet loss

Loss Function vs Cost Function

| cgpa | iq | (y$_i$) package | $\hat{y}_i$ Prediction |
|------|-----|---------|------------|
| 6.3 | 100 | 6.3 | 6.1 |
| 7.1 | 91 | 4.1 | 4 |
| 8.5 | 83 | 3.5 | 3.7 |
| 9.2 | 102 | 7.2 | 7 |

cgpa

iq

$\hat{y}_i$

$$\boxed{\text{Loss Function vs Cost Function}}$$

cgpa | iq | (y_i) package | $\hat{y}_i$ Prediction

| cgpa | iq | (y_i) package | $\hat{y}_i$ Prediction |
|------|-----|---------------|------------------------|
| 6.3 | 100 | 6.3 | 6.1 |
| 7.1 | 91 | 4.1 | 4 |
| 8.5 | 83 | 3.5 | 3.7 |
| 9.2 | 102 | 7.2 | 7 |

Loss function → single training eg

$y_1 = 6.3 \qquad \hat{y}_i = 6.1$

$(y_i - \hat{y}_i)^2$

$(6.3 - 6.1)^2 =$

$$\boxed{\text{Loss Function Vs Cost Function}}$$

cgpa • 

→ iq •

$\hat{y}_i$

| cgpa | iq | package $(y_i)$ | Prediction $\hat{y}_i$ |
|------|-----|--------|------------|
| 6.3 | 100 | 6.3 | 6.1 |
| 7.1 | 91 | 4.1 | 4 |
| 8.5 | 83 | 3.5 | 3.7 |
| 9.2 | 102 | 7.2 | 7 |

batch

Error funch.

Loss function → single training eg

$$y_1 = 6.3 \qquad \hat{y}_1 = 6.1$$

$$(y_i - \hat{y}_i)^2$$

$$(6.3 - 6.1)^2 =$$

Cost function

$$\frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

$$\frac{1}{4}\left[(6.1-6.3)^2 + (4.1-4)^2 + (3.5-3.7)^2 + (7.2-7)^2\right] = CF$$

1. Mean Squared Error (MSE)

   Squared loss    L2 Loss

| $y$ cgpa | iq | $y_i)$ package | $\hat{y}_i$ Prediction |
|---|---|---|---|
| 6.3 | 100 | 6.3 | 6.1 |
| 7.1 | 91 | 4.1 | 4 |
| 8.5 | 83 | 3.5 | 3.7 |
| 9.2 | 102 | 7.2 | 7 |

## 1. Mean Squared Error (MSE)

Squared loss    L2 Loss

$\downarrow (y_i - \hat{y}_i)^2$        Advan    DisAdg

$(true - predict)^2$

$(6.3 - 6.1)^2 = -$

$$\boxed{(y_i - \hat{y}_i)^2}$$

| cgpa | iq | package $y_i$ | Prediction $\hat{y}_i$ |
|------|-----|---------|------------|
| 6.3 | 100 | 6.3 | 6.1 |
| 7.1 | 91 | 4.1 | 4 |
| 8.5 | 83 | 3.5 | 3.7 |
| 9.2 | 102 | 7.2 | 7 |



cgpa → ... → $\hat{y}_i$

iz

forwar

## 1. Mean Squared Error (MSE)

Squared loss   L2 loss

$\downarrow \; (y_i - \hat{y}_i)^2$

Advan   DisAdg

$(true - predict)^2$

$(6.3 - 6.1)^2 = -$

$(y_i - \hat{y}_i)^2$

$\boxed{(y_i - \hat{y}_i)^2}$

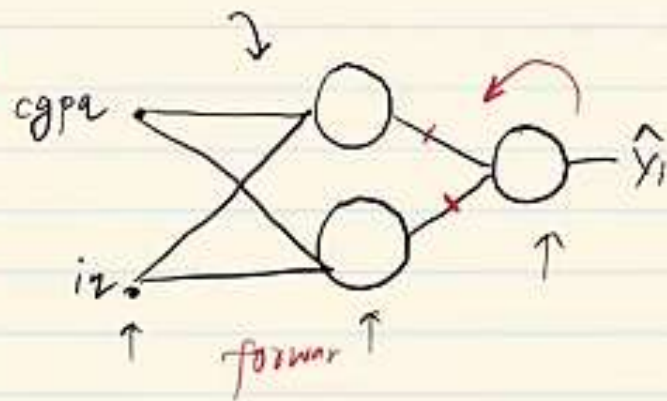| cgpa | Iq | package $y_i)$ | Prediction $\hat{y}_i$ | $y_i - \hat{y}_i$ |
|------|-----|---------|------------|-----------|
| 6.3 | 100 | 6.3 | 6.1 ← | 0.2 |
| 7.1 | 91 | 4.1 | 4 | 0.1 |
| 8.5 | 83 | 3.5 | 3.7 | -0.2 |
| 9.2 | 102 | 7.2 | 7 | 0.2 |

overall error

# 1. Mean Squared Error (MSE)

Squared loss    L2 Loss

$\downarrow (y_i - \hat{y}_i)^2$

$(true - predict)^2$

$(6.3 - 6.1)^2 = -$

$$\boxed{(y_i - \hat{y}_i)^2}$$

Advan    DISAdg

mae

$\boxed{(y_i - \hat{y}_i)^2}$

$\uparrow$ quadratic    punish

$\wedge^2$

true - predic $\longrightarrow$ 1 unit $\longrightarrow$ 1 unit

$\qquad$ 2 unit $\longrightarrow$ 4 uint

$\qquad$ 4 win $\longrightarrow$ 16 unit

magnify

| cgpa | lq | package | Prediction | $y_i - \hat{y}_i$ |
|------|-----|---------|------------|---------|
| 6.3 | 100 | 6.3 | 6.1 ← | 0.2 |
| 7.1 | 91 | 4.1 | 4 | → 0.1 |
| 8.5 | 83 | 3.5 | 3.7 | -0.2 |
| 9.2 | 102 | 7.2 | 7 | 0.2 |

overall error



cgpa

wb

iz

$\hat{y}_i$

drastic big

forwar

$(y_i - \hat{y}_i)^2$    $(w,b) \rightsquigarrow \textcircled{L}\ \underline{min}$

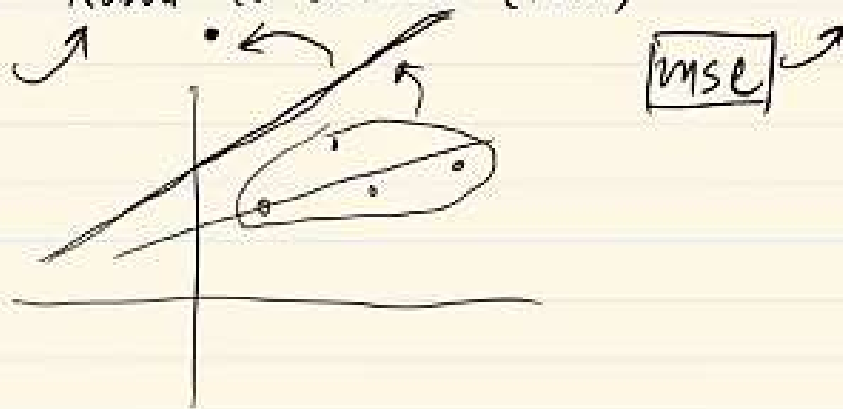**Advantages**

1) Easy to interpret ↗
2) Differentiable (GD)
3) 1 local minima

**Disadvantage**

1) Error unit (squared) → diff
2) Robust to Outliers (Not)

$\boxed{mse}$ ↗

# 1. Mean Squared Error (MSE)

Squared loss   L2 loss

$(y_i - \hat{y}_i)^2$

Advan    DisAdv

$(true - predict)^2$

mae

$(6.3 - 6.1)^2 = -$

$(y_i - \hat{y}_i)^2$

$\boxed{(y_i - \hat{y}_i)^2}$

quadratic   punish

$\wedge^2$

true - predict → 1 unit → 1 unit

magnify         2 unit → 4 uint

                4 unit → 16 unit

$(y_i - \hat{y}_i)^2$

$(w, b) \rightsquigarrow \bigcirc L \; \underline{min}$

| cgpa | iq | package | Prediction | $(y_i - \hat{y}_i)$ |
|------|-----|---------|-----------|---------|
| 6.3 | 100 | 6.3 | 6.1 | 0.2 |
| 7.1 | 91 | 4.1 | 4 | 0.1 |
| 8.5 | 83 | 3.5 | 3.7 | -0.2 |
| 9.2 | 102 | 7.2 | 7 | 0.2 |
| 8.1 | 100 | 50 | 8 | |

Overall error

cgpa   relu

iz.

former   activation   drastic

Sigmoid   d linear   big

$(42)^2$

W

2. Mean Absolute Error (MAE) $\rightarrow$ L1 loss

$$L = |y_i - \hat{y}_i|$$

$$C = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|^2$$

$2|{\sim}abs$

$|true - predict|$

punish

$2\,lpa$

$C\delta p|\ iq\,|\ pace\ lp2\ |(y_i - \hat{y}_i)|$

$\searrow lpa$

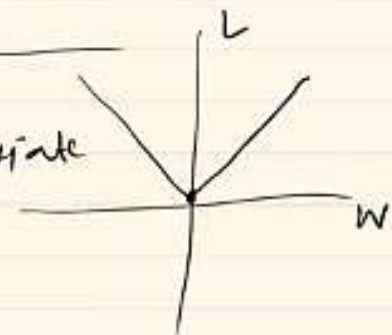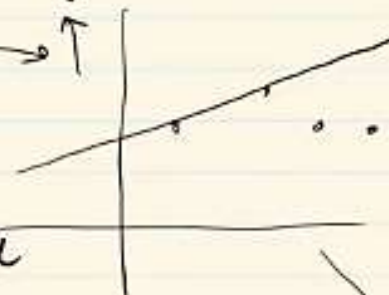**Advantages**

1) Intuitive and easy
2) Unit $\rightarrow$ same $-y$
3) Robust to outliers

**Disadvange**

1) Not diffrentiable
   $\searrow$ GD $\supset$ differentiate
   Subgradient $\searrow$
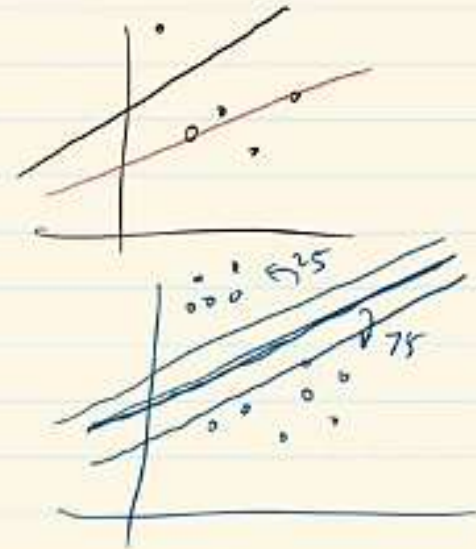
3. Huber Loss ✓

$$L = \begin{cases} \frac{1}{2}(y-\hat{y})^2 & \text{for } |y-\hat{y}| \le \delta \\ \delta|y-\hat{y}| - \frac{1}{2}\delta^2 & \text{otherwise} \end{cases}$$

mse — outliers ✓
mae — normal points ✓

$\downarrow$ mae  huber  $\uparrow$ mse
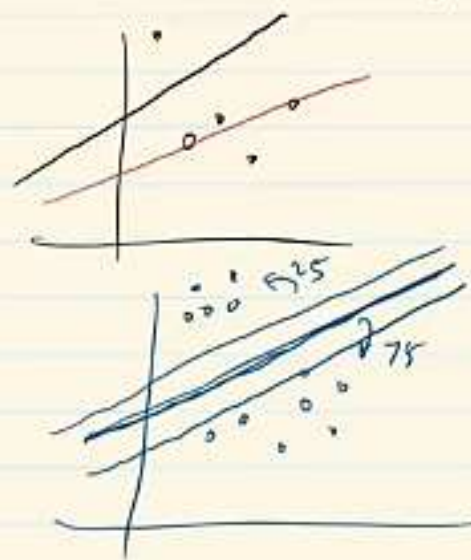
3. Huber Loss ✓

mae — normal point ✓

mse — outliers ✓

nyper

$$L = \begin{cases} \frac{1}{2}(y-\hat{y})^2 & \text{for } |y-\hat{y}| \leq \delta \quad \text{mse} \\ \delta|y-\hat{y}| - \frac{1}{2}\delta^2 & \text{otherwise} \quad \text{mae} \end{cases}$$

mse

δ25

δ75

4. Binary Cross Entropy

→ classification
→ Two classes

$\rightarrow 1 \quad 0$

| cgpa | iq | placement |
|------|-----|-----------|
| 8 | 80 | 1 |
| 7 | 70 | 0 |
| 6 | 60 | 0 |

Loss function $= -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$

$y \rightarrow$ actual value / target
$\hat{y} \rightarrow$ NN prediction

relu
sigmoid

hidden

Activation

{ sigmoid }

4. Binary Cross Entropy

→ classification
→ Two classes

cgpa | iq | placement

→ 1 0

| cgpa | iq | placement |
|---|---|---|
| 8 | 80 | 1 |
| 7 | 70 | 0 |
| 6 | 60 | 0 |

0.73

fP

relu   0-1
sigmoid

Activation

hidden

ζ sigmoid ζ

Loss function $= -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$

$y \rightarrow$ actual value / target
$\hat{y} -$ NN prediction

Cost function $= -\dfrac{1}{n}\left[\displaystyle\sum_{i=1}^{n} y_i \log \hat{y_i} + (1-y_i)\log(1-\hat{y_i})\right]$

4. **Binary Cross Entropy**

$\rightarrow$ classification

$\rightarrow$ Two classes

100 days

cgpa | iq | placement

| 8 | 80 |
|---|---|
| 7 | 70 |
| 6 | 60 |

1
0
0

$\rightarrow$ 1 0     GD     fP

$\rightarrow$ relu     0-1
sigmoid

b

w

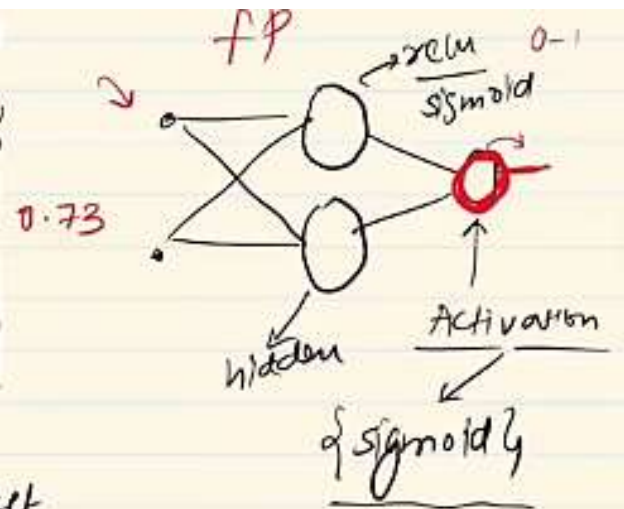$\begin{bmatrix} 0.73 \\ 0.25 \end{bmatrix}$

0.13

hidden

[Activation]

{ sigmoid }

0.12

Loss function $= \boxed{-y \log(\hat{y}) - (1-y)\log(1-\hat{y})}$

$y \rightarrow$ actual value / target

$\hat{y} \rightarrow$ NN prediction

Keras

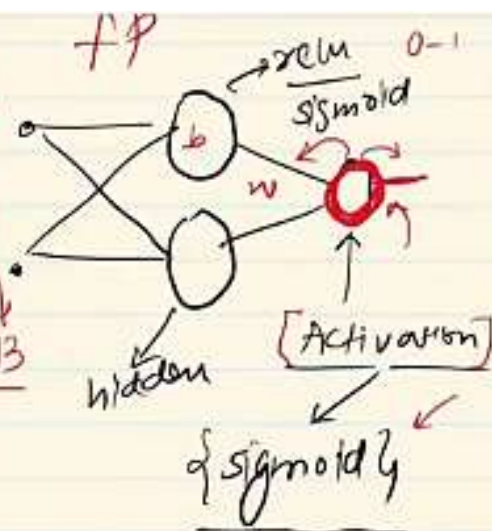0.12     $-(1-0)\log(1-0.25)$

$-1\log(0.75)$     $-0.12$

Cost function $= -\dfrac{1}{n}\left[\displaystyle\sum_{i=1}^{n} y_i \log \hat{y}_i + (1-y_i)\log(1-\hat{y}_i)\right]$

maximum likli

$-1\log(0.73)$

$-1 \times -0.13 = 0.13$

Logistic Reg

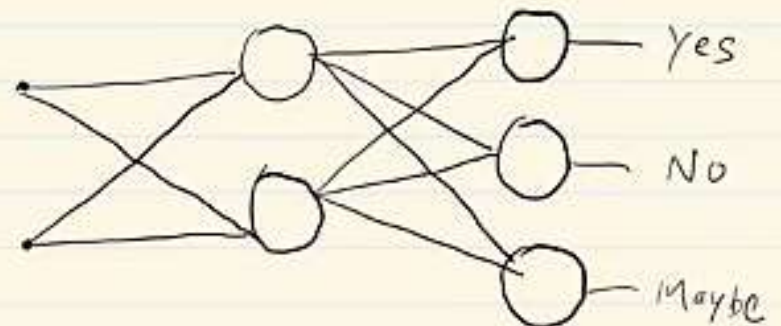5. Categorical Cross Entropy [used in Softmax Regression] ↓

→ [Multi-class] {classification}

$$L = -\sum_{j=1}^{k} Y_j \log(\hat{Y}_j)$$

where K is # classes in the data

| cgpa | iq | placed? | Yes | No | Maybe |
|------|-----|---------|-----|-----|-------|
| 8 | 80 | Yes 1 | 1 | 0 | 0 |
| 6 | 60 | No 2 | 0 | 1 | 0 |
| 7 | 70 | Maybe 3 | 0 | 0 | 1 |

Yes

No

Maybe

5. Categorical Cross Entropy [used in Softmox Regression] ↓

→ [Multi-class] {classification}



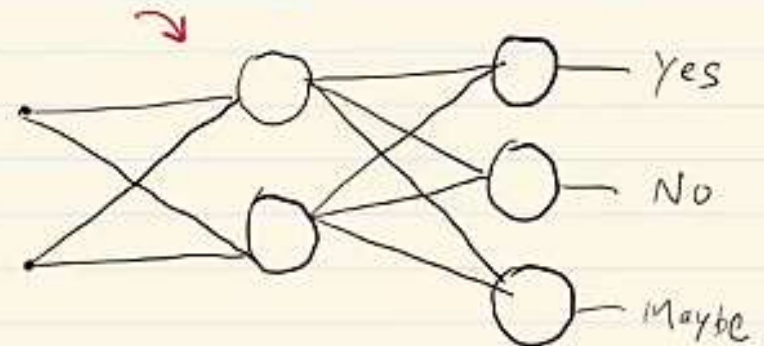$$L = - \sum_{j=1}^{\kappa} y_j \log (\hat{y}_j)$$

1 point →

where Ⓚ is # classes in the data
↳3

1 point

$$L = - y_1 \log(\hat{y}_1) - y_2 \log (\hat{y}_2) - y_3 \log (\hat{y}_3)$$

| cgpa | iq | placed? | Yes | No | Maybe |
|------|-----|---------|-----|-----|-------|
| 8 | 80 | Yes 1 | 1 | 0 | 0 |
| 6 | 60 | No 2 | 0 | 1 | 0 |
| 7 | 70 | Maybe 3 | 0 | 0 | 1 |

Yes

No

Maybe

5. Categorical Cross Entropy [used in Softmax Regression] ↙

→ [Multi-class] {classification}

$$L = - \sum_{j=1}^{K} y_j \log(\hat{y_j})$$

1 point →

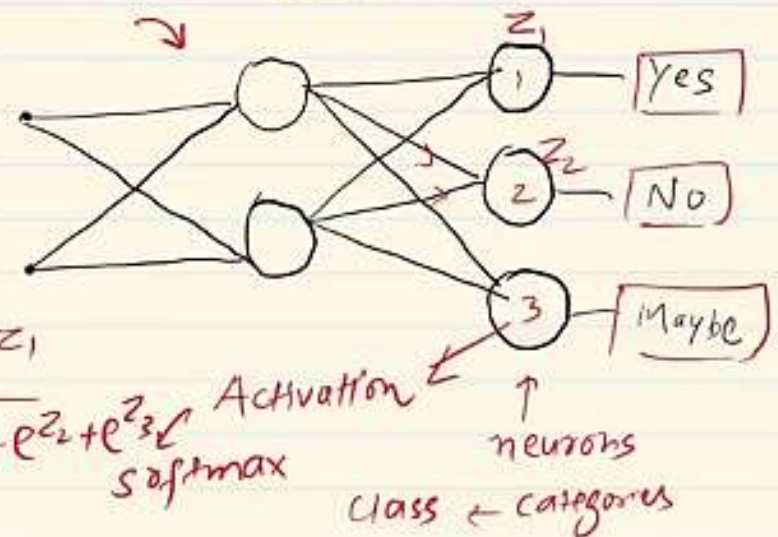| cgpa | iq | placed? | Yes | No | Maybe |
|------|-----|---------|-----|-----|-------|
| 8 | 80 | Yes 1 | 1 | 0 | 0 |
| 6 | 60 | No 2 | 0 | 1 | 0 |
| 7 | 70 | Maybe 3 | 0 | 0 | 1 |

where (K) is # classes in the data
↳3

1 point

$$L = - y_1 \log(\hat{y_1}) - y_2 \log(\hat{y_2}) - y_3 \log(\hat{y_3})$$

$$\frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_1}}$$

$$\frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$f(z) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

softmax



Activation

neurons

Class ← categories

$z_1$ — Yes
$z_2$ — No
3 — Maybe

**5. Categorical Cross Entropy** [used in Softmax Regression] ↓↙ → OHE

→ [Multi-class] {classification}

$$L = -\sum_{j=1}^{k} y_j \log(\hat{y}_j)$$

1 point

| cgpa | iq | placed? | Yes | No | Maybe | |
|------|-----|---------|-----|-----|-------|---|
| 8 | 80 | Yes 1 | 1 | 0 | 0 | ⊢ Yes |
| 6 | 60 | No 2 | 0 | 1 | 0 | ⊢ No |
| 7 | 70 | Maybe 3 | 0 | 0 | 1 | ⊢ Maybe |

where (K) is # classes in the data
↳3

1 point

$$L = -y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - y_3 \log(\hat{y}_3)$$

[0.2  0.3  0.5]

[1   0   0]

$$\frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_1}}$$

$$\frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$f(z) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$  Activation
(softmax)



$z_1$ 0.2  0-1  [Yes]

$z_2$ 0.3  0-1  [No] = 2

0.5  0-1  [Maybe]

forward prop

neurons

class ← categories

Reg → mse

Outlier — mae → huber loss

Classification → binary → bce

multi → CCE

SCE