

# Softmax Regression / Multinomial Regression

Date \_\_\_\_\_  
Page \_\_\_\_\_

→ for multiclass (more than 2) containing dataset.

CPA IQ Placement

Yes

No

Optout (Not sit in a placement)

No of class = 3

For above case simple logistic Regression is not work, make it worse on this we have to applied some logic,

- ① One hot encoding.
- ② Devide & conquer.

We transform above dataset into below

CPA	IQ	Yes	No	Optout
5	80	0	1	0
10	70	1	0	0
9	80	0	0	1

This is nothing but one hot encoding / Nominal encoding.

For that we use softmax function,

$$\sigma(t_i) = \frac{e^{t_i}}{\sum_{j=1}^K e^{t_j}}$$

let,

$$y_{cs} \Rightarrow 1$$

$$No \Rightarrow 2$$

$$Optout \Rightarrow 3$$

For yes,

$$\sigma(t_1) = \frac{e^{t_1}}{e^{t_1} + e^{t_2} + e^{t_3}}$$

For no,

$$\sigma(t_2) = \frac{e^{t_2}}{e^{t_1} + e^{t_2} + e^{t_3}}$$

For optout,

$$\sigma(t_3) = \frac{e^{t_3}}{e^{t_1} + e^{t_2} + e^{t_3}}$$

Here,  $K$  is the NO of classes.

$\text{yes} + \text{no} + \text{optout}$

$$= \frac{e^{t_1}}{e^{t_1} + e^{t_2} + e^{t_3}} + \frac{e^{t_2}}{e^{t_1} + e^{t_2} + e^{t_3}} + \frac{e^{t_3}}{e^{t_1} + e^{t_2} + e^{t_3}}$$

$$= \frac{e^{t_1 + t_2 + t_3}}{e^{t_1 + t_2 + t_3}}$$

$$= 1$$

Note:-

$$\therefore o(t)_1 + o(t)_2 + o(t)_3 + \dots + o(t)_k = 1$$

Big question is,

what is the value of  $o(t)$ ?

let's derive.

et

$$st \quad g_{lr} + p_{xr}^0 g_{lr} + 8x_{lr}^0 w = 1 \quad (1)$$

$$st \quad g_{lr} + 2x_{lr}^0 w + p_{xr}^0 l_r = 1 \quad (2)$$

$$st \quad g_{lr} + 0.2x_{lr}^0 w + 8x_{lr}^0 w = 1 \quad (3)$$

rect,

CAPA IQ placement

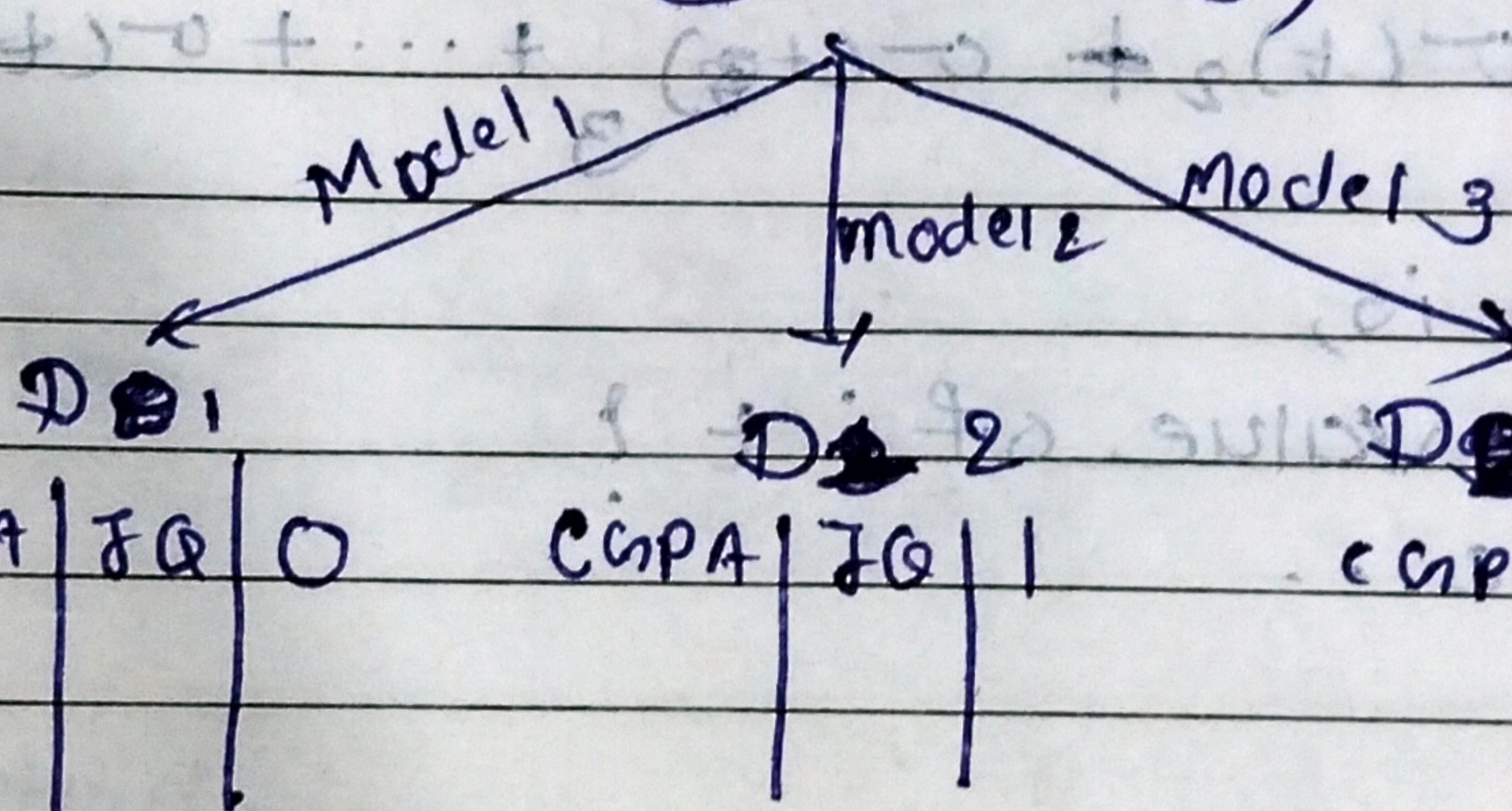
8	59	0
9	75	1
8.8	80	2

SO,

CAPA IQ placement 0  
placement 0 placement 1 placement 2

8	59	0	1	0	0
9	75	0	0	1	0
8.8	80	0	0	0	1

Dataset (D)



$$\text{coef}_1: w_1^0, w_2^0, w_3^0 \quad \text{coef}_2: w_1^1, w_2^1, w_3^1 \quad \text{coef}_3: w_1^2, w_2^2, w_3^2$$

$$① t_1 = w_1^0 x_8 + w_2^0 x_9 + w_3^0 \quad t_2 = w_1^1 x_8 + w_2^1 x_9 + w_3^1 \quad t_3 = w_1^2 x_8 + w_2^2 x_9 + w_3^2$$

$$② t_1 = w_1^0 x_9 + w_2^0 x_75 + w_3^0 \quad t_2 = \dots \quad t_3 = \dots$$

$$③ t_1 = w_1^0 x_{8.8} + w_2^0 x_{80} + w_3^0$$

→ In above we decide our dataset & considering one dataset at time we calculate weights & bias.

→ This is how we get the value of  $t_1, t_2, t_3$ .

→ Now we can apply gradient descent to obtain  $w_1^0, w_2^0, w_3^0, w_1^1, w_2^1, w_3^1, w_1^2, w_2^2, w_3^2$

→ evaluate them separate we can also merge them.

We have loss function,

$$L = - \left( \sum_{i=1}^n y_i (\log \hat{y}_i) - \sum_{i=1}^n (1-y_i) (\log (1-\hat{y}_i)) \right)$$

For logistic regression, but we can modify it such a way that it can handle combined scenario also,

$$L = -\frac{1}{n} \left( \sum_{i=1}^n \sum_{j=1}^K y_{ij} \log(\hat{y}_{ij}) \right)$$

Means,

$\text{for } i = 0$

1<sup>st</sup> Row, Prediction of model 1 +

Prediction of model 2 +

Prediction of model 3.

$i = 1$

2<sup>nd</sup> Row,

Prediction of model 1

Prediction of model 2

Prediction of model 3

$k = \text{NO of class / model}$

→ Remember  $\hat{y}$  is never been 0 so  $\log(\hat{y})$  never been undefined because  $\hat{y}$ , the prediction is pass through the sigmoid function then store into  $\hat{y}$ .

→ In loss function of Logistic Regression, if in  $m$  classes, <sup>actual value of dependent variable</sup> prediction is 1 then  $-\frac{1}{m} y_i (\log \hat{y}_i)$  is calculated & if prediction

actual value of dependent variable is 0 then  $-\frac{1}{m} (1 - y_i) (\log \hat{y}_i)$  is calculated. ~~but~~

$-\frac{1}{m} (1 - y_i) (\log \hat{y}_i)$  is calculated. ~~but~~

→ But here as we conquer all the classes we only calculate  $-\frac{1}{m} y_i (\log \hat{y}_i)$  as another class is in another model.

$$L = -\frac{1}{n} \left( \sum_{i=1}^n \sum_{j=1}^K y_j^i (\log(\hat{y}_j^i)) \right)$$

$n=3$

$K=3$

$m=2$

means,  $\text{dof}=2$  & model=3 / class=3

$x_1$	$x_2$	$y(0/1/2)$
5	6	0
5	6	1
5	6	2

↓ conversion,

$x_1$	$x_2$	$y$	$x_1 \cdot x_2$
5	6	0	30
5	6	1	-30
5	6	2	0

$$\sqrt{=} f^{-1}(y)$$

$$L = -\frac{1}{n} \left( \sum_{i=1}^n \sum_{j=1}^3 y_j^i (\log(\hat{y}_j^i)) \right)$$

$$= -\left( y_1^1 (\log \hat{y}_1^1) + y_2^1 (\log \hat{y}_2^1) + y_3^1 (\log \hat{y}_3^1) \right)$$

$$= -\left( 1 (\log \hat{y}_1^1) + 0 (\log \hat{y}_2^1) + 0 (\log \hat{y}_3^1) \right)$$

$$L = -\log \hat{y}_1^1$$

Above is just an example for understanding how the loss function work!

Now for  $w_1^0, w_2^0, w_3^0, w_1^1, w_2^1, w_3^1, \dots$   
 we differentiate L with  $w_1^0, w_2^0, w_3^0, \dots$

Because in gradient decent

$$w = w - \eta \frac{dL}{dw}$$

$$w = \begin{bmatrix} w_0^{(0)} & w_1^{(0)} & w_2^{(0)} \\ w_0^{(1)} & w_1^{(1)} & w_2^{(1)} \\ w_0^{(2)} & w_1^{(2)} & w_2^{(2)} \end{bmatrix}$$

$$((\hat{e}_{POI})_0 + (\hat{e}_{POI})_1 + (\hat{e}_{POI})_2) +$$

$$(\hat{e}_{POI})_0 + (\hat{e}_{POI})_1 + (\hat{e}_{POI})_2 +$$

$$(\hat{e}_{POI})_0 + (\hat{e}_{POI})_1 + (\hat{e}_{POI})_2 +$$

$$(\hat{e}_{POI})_0 + (\hat{e}_{POI})_1 + (\hat{e}_{POI})_2 +$$