# What is multiclass Classification?

| Gpa | iq | placements |
|-----|-----|-----------|
| 9 | 90 | |
| 6 | 60 | |
| 7 | 70 | |

Y
N ———→ 3 classes
Opt



x → opt
x → Y
x → N

---

# How to Logistec Regression handle Multiclass Classification problems?

* mostly use in binary classification

* But now use in multiclassification with some technique

(i) OVR (one vs Rest) → OVA (One vs All)

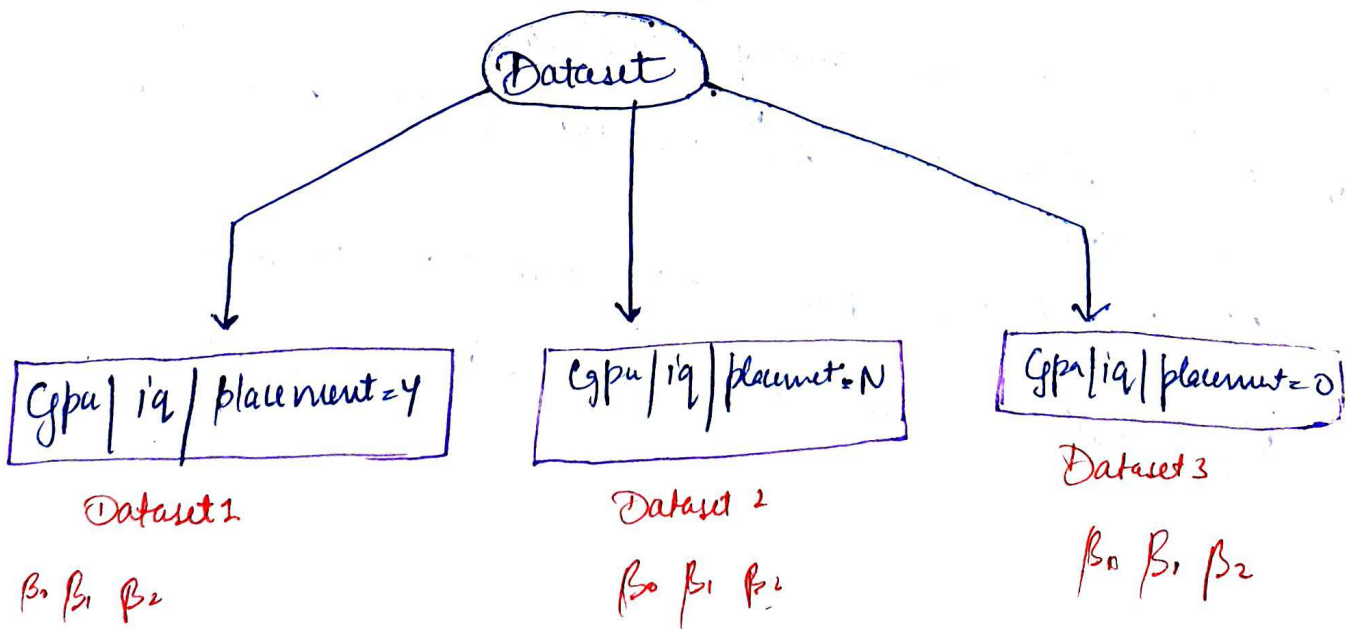(ii) Multinomial Logistic Regression → Soft Max Reg

# OVR Approach

train → prediction → code
↓

| Cgpa | iq | Placement |
|------|-----|-----------|
| — | — | Y |
| — | — | N |
| — | — | $\theta$ |

$(Y, N, \theta) \rightarrow$ Logistic Reg model

* No. of category = apply No. of time Logistic Reg model.

## Using One Hot Encoding :

| Cgpa | iq | placement = y | placement = N | placement = $\theta$ |
|------|-----|---------------|---------------|----------------------|
| — | — | 1 | 0 | 0 |
| — | — | 0 | 1 | 0 |
| — | — | 0 | 0 | 1 |

Dataset

| Cgpa | iq | placement = y |
|------|-----|---------------|

Dataset 1

$\beta_0 \ \beta_1 \ \beta_2$

| Cgpa | iq | placement = N |
|------|-----|---------------|

Dataset 2

$\beta_0 \ \beta_1 \ \beta_2$

| Cgpa | iq | placement = $\theta$ |
|------|-----|----------------------|

Dataset 3

$\beta_0 \ \beta_1 \ \beta_2$

\* Apply Normal Logistic Regression on all three dataset.

example:-

Prediction $\{6.5, 65\}$ → $Y, N, \theta$

P(Y)     P(N)     P(\theta)

0.6      0.3      0.5

## Normalize:

$$P(Y) = \frac{0.6}{0.6+0.3+0.5} = 0.42 = 42\%.$$

$$P(N) = \frac{0.3}{0.6+0.3+0.5} = 0.23 = 23\%.$$

$$P(\theta) = \frac{0.5}{0.6+0.3+0.5} = 0.35 = 35\%.$$

$$P(Y) + P(N) + P(\theta) = 0.42 + 0.23 + 0.35 = 1$$

\* Highest Probablity is output like 0.42 is Highest so, Output is Y.

\* Not efficient with large dataset having high number of class.

code

# Soft Max function

$$\sigma(\vec{z})_i = \frac{\sigma^{z_i}}{\sum_{j=1}^{k} e^{z_i}}$$

$z_1 \qquad z_2 \qquad z_3$

$$\sigma(\vec{z_1}) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

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

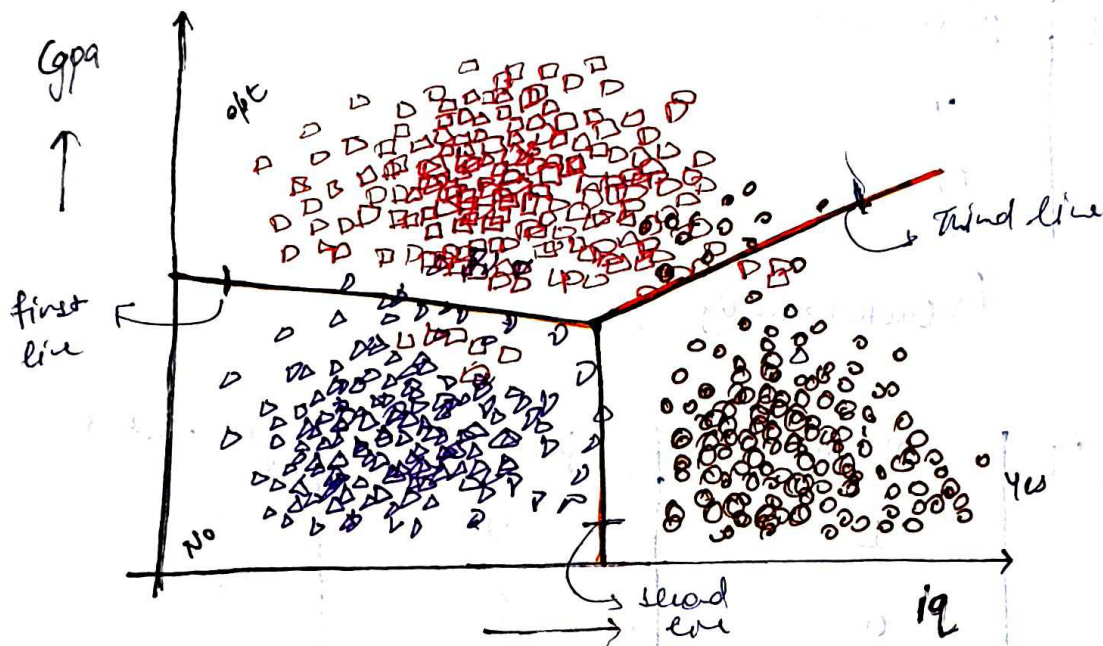$$e(\vec{z_3}) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\sigma(\vec{z_1}) + e(\vec{z_2}) + e(\vec{z_3}) = 1$$

$\Downarrow$

multi class

# Softmax Logistic Regression / Multinomial LR



for k classes, Softmax LR will draw k lines to create the decision Region

\# for Binary we have

$$Ax + By + C = 0$$
$$B_0 + B_1 x_1 + B_2 x_2 = 0$$

$$\boxed{\begin{array}{ccc} B_0 & B_1 & B_2 \end{array}}$$

\# for 3 lines we have $3 \times 3$ total 9 parameters.

3 for each line

\# assumption for 2D data ( 2 feature)

for 3 (feature) = $3 \times 4$ = 12 parameter

category

$\rightarrow B_0 \ B_1 \ B_2 \ B_3$

training → Prediction → code
          ↓

| Cgpa | iq | placement |
|------|-----|-----------|
| -    | -   | Y         |
| -    | -   | N         |
| -    | -   | O         |

↳ Onehot encoding

| Cgpa | iq | placement Y | placement N | placement O |
|------|-----|-------------|-------------|-------------|
| -    | -   | 1           | 0           | 0           |
| -    | -   | 0           | 1           | 0           |
| -    | -   | 0           | 0           | 1           |

## Loss function :

general case →

$$L = -\frac{1}{n} \sum_{i=1}^{n} \sum_{k=1}^{k} Y_i^k \log \hat{Y}_i^k$$  → SoftMax LR minimize

↖ No. of classes

Binary :-

$$L = -\frac{1}{n} \sum_{i=1}^{n} Y_i \log(\hat{Y}_i) + (1-Y_i) \log(1-\hat{Y}_i)$$

↳ Special case only for $k=2$

| cgpa | iq | placement=Y | placement=N | placement=0 |
|------|-----|------|------|------|
| — | — | 1 | 0 | 0 |
| — | — | 0 | 1 | 0 |
| — | — | 0 | 0 | 1 |

let assume we have only 1 row

$$n=1$$

So, $\sum\limits_{i=1}^{n} = 1$

$$L = \sum_{k=1}^{k} y^k \log \hat{y}^k$$

$$L = y^{yes} \log \hat{y}^{yes} + y^{No \to 0} \log \hat{y}^{No} + y^{\theta \to 0} \log \hat{y}^{\theta}$$

$$= (1) \log \hat{y}^{yes} \quad + \quad 0 + 0$$

**For 2nd row**

$$L = 0 + (1) \log \hat{y}^{No} + 0$$

**For 3rd row**

$$L = 0 + 0 + (1) \log \hat{y}^{\theta}$$

\* In every row, we want those class which value has 1.

$$L = \log \hat{y}_1^{yes} + \log \hat{y}_2^{No} + \log \hat{y}_3^{\theta}$$

$\hookrightarrow$ 1 row      $\hookrightarrow$ 2 row      $\hookrightarrow$ 3 row

$\hookrightarrow \beta_0, \beta_1, \beta_2 \to$ 1 line    $\hookrightarrow \beta_0, \beta_1, \beta_2 \to$ 2nd line    $\hookrightarrow \beta_0, \beta_1, \beta_2 \to$ 3nd line

$\hat{y}_1^{yes} = ?$      $\hat{y}_2^{No} = ?$      $\hat{y}_3^{\theta} = ?$

Sigmoid $\hookrightarrow (0-1)$

$\boxed{\hat{y}_i}$ $\longrightarrow$ output of Logistic Regression

$\hat{y}_i = \sigma(z_i)$

$$\boxed{z_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}}$$

$\to$ for Normal Logistic R.

$$\hat{y}_1^{yes} = \sigma(z_i) = \frac{e^{z_{yes}}}{e^{z_{yes}} + e^{z_{no}} + e^{z_{opt}}}$$

$\hookrightarrow$ Softmax

$\underbrace{\quad\quad\quad}_{1 \text{ row}}$

$$z_{yes} = \underbrace{\beta_0 + \underbrace{\beta_1 x_{i1} + \beta_2 x_{i2}}}_{}$$

coefficient and intersept of first line

$$z_{No} = \underbrace{\beta_0 + \underbrace{\beta_1 x_{i2} + \beta_2 x_{i2}}}_{}$$

coefficient and intercept of second line

$$z_{opt} \quad \underbrace{\beta_0 + \underbrace{\beta_1 x_{i1} + \beta_2 x_{i2}}}_{}$$

coefficient and intercept of Third line

$$\hat{y}_{\alpha}^{no} = \sigma(z_i) = \frac{e^{z_{\text{no}}}}{e^{z_{\text{yes}}} + e^{z_{no}} + e^{z_{opt}}}$$

$$\hat{y}_{\beta}^{op} = \sigma(z_i) = \frac{e^{z_{opt}}}{e^{z_{\text{yes}}} + e^{z_{no}} + e^{z_{opt}}}$$

$$\begin{cases} \begin{bmatrix} \beta_0^1 & \beta_1^1 & \beta_2^1 \\ \beta_0^2 & \beta_1^2 & \beta_2^2 \\ \beta_0^3 & \beta_1^3 & \beta_2^3 \end{bmatrix} \end{cases} \begin{array}{l} \rightarrow \text{1st line} \\ \rightarrow \text{2nd line} \\ \rightarrow \text{3rd line} \end{array}$$

$\rightarrow$ Start with 3 random line then start
with Gradient descent

$$\beta_0^1 = \beta_0^1 - \eta \frac{\partial L}{\partial \beta_0^1}$$

$$\frac{\partial L}{\partial \beta_0^3} = \frac{-1}{n} \sum_{i=1}^{n} \sum_{k=1}^{k} y_i \log (\hat{y}_i^k)$$

$\left.\begin{array}{l} \text{doing this for 9} \\ \text{different lines} \\ \text{in 1000 times} \end{array}\right\}$

* after 1000 times   we get a value which less
Loss and good best fit line .

Let assume $k = 2$ $(0, 1)$ then our Loss function will be binary cross entropy

$$= -\frac{1}{n} \sum_{i=1}^{n} y_i^1 \log(\hat{y}_i^1) + y_i^0 \log(\hat{y}_i^0)$$

$\boxed{y_i^1 \quad y_i^0}$ $\rightarrow$ you like Placement $\hat{N}_i\hat{y}$

$$1 \Rightarrow \quad y_i^1 = 1 \quad y_i^0 = 0$$

$$0 \Rightarrow \quad y_i^1 = 0 \quad y_i^0 = 1$$

So,

$$= -\frac{1}{n} \sum_{i=1}^{n} \overset{y}{y_i^1} \left( \log \overset{y}{\hat{y}_i^1} \right) + (1 - \overset{N}{y_i^1}) \log(1 - \overset{\tilde{N}}{\hat{y}_i^1})$$

$$\boxed{L = -\frac{1}{n} \sum_{i=1}^{n} y_i \log \hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)}$$

$\hookrightarrow$ Binary cross entropy

$\boxed{\text{prooved}}$

# When to Use what?

## Use One-vs-Rest (OVR) when:

1. Classes are Non-Mutually Exclusive: OVR is appropriate if an instance can belong to more than one class, as each classifier provides an independent probability for each class.

2. Dealing with Imbalanced Data: OVR might perform better when class distribution is highly imbalanced since each class gets a dedicated model.

## Use Multinomial Logistic Regression (SoftMax Regression) when:

1. Computational Efficiency is Required: Softmax Regression is generally more efficient for large datasets and a high number of classes.

2. Classes are Mutually Exclusive: SoftMax Regression is a good choice when each instance can only belong to one class. The SoftMax function provides a set of probabilities that sum to 1, fitting well with mutually exclusive classes.

3. Interpretability is important: The probabilities output by SoftMax Regression are more interpretable than those from OVR, as they always sum to 1. This can make model predictions easier to explain.

$\rho$ Gradient descent     $\dashrightarrow$ Predict $\{6.5, 65\}$

$$\begin{array}{c|ccc}
\text{line 1} \rightarrow & \beta_0^1 & \beta_1^1 & \beta_2^1 \\
\text{line 2} \rightarrow & \beta_0^2 & \beta_1^2 & \beta_2^2 \\
\text{line 3} \rightarrow & \beta_0^3 & \beta_1^3 & \beta_2^3
\end{array}$$

$Z_1 = \beta_0^1 + 6.5\,\beta_1^1 + 65\,\beta_2^1$

$Z_2 = \beta_0^2 + 6.5\,\beta_1^2 + 65\,\beta_2^2$

$Z_3 = \beta_0^3 + 65\,\beta_1^3 + 65\,\beta_2^3$

Softmax $\Rightarrow$ $\sigma(Z_1) = \dfrac{e^{Z_1}}{e^{Z_1} + e^{Z_2} + e^{Z_3}}$

$\sigma(Z_2) = \dfrac{e^{Z_2}}{e^{Z_1} + e^{Z_2} + e^{Z_3}}$       $\sigma(Z_3) = \dfrac{6^{Z_3}}{e^{Z_1} + e^{Z_2} + e^{Z_3}}$

* get probability and high prob. is outcome
    or predicted outcome

$Z_1 = \beta_0^1 + 6.5\,\beta_1^1 + 65\,\beta_2^1$