# LAB-4

# Logistic Regression

## Lab-4
## Logistic Regression

1) Consider a binary classificat$^n$ where we want to predict whether a student will pass or fail based on their study hours. The logistic regression model has been trained and learned parameters are $a_0 = -5$
$a_1 = 0.8$.

a) Write Logistic regres eq$^n$ for this probl

$$z = a_0 + a_1 x$$
$$= -5 + 0.8 x$$

$$y = \frac{1}{1 + e^{(-z)}}$$

b) Calculate the probability that a student who studies for 7 hours will pass

$$z = a_0 + a_1 x$$
$$= -5 + 0.8 (7)$$
$$= -5 + 5.6$$
$$= 0.6$$

$$y = \frac{1}{1 + e^{-(0.6)}}$$
$$= 0.64$$

Threshold = 0.5    ∴ 0.64 > 0.5
The student pass.

c) Determine the predicted class for this student based on threshold of 0.5

2) Consider $z = [2, 1, 0]$ of 3 classes.
Apply softmax fun' to find probability value of 3 classes

$$S = S_2 = \frac{c^2}{c_2 + c_1 + c_0} = \frac{c^2}{2 + 1 + 0} = 0.66$$

$$S_1 = \frac{c^1}{c_2 + c_1 + c_2} = 0.244$$

$$S_0 = \frac{c^0}{c^2 + c^1 + c^0} = 0.090$$

1) for dataset file 'HR_comma_sep.csv'

i) which variable did you identify as having a direct & clear impact on employee retention? Why?

   Left has direct impact whether employee stayed (0) or left (1)

ii) what was the accuracy of Logistic reg model? is it good accuracy? why or why not?

   Accuracy $= 0.76 \approx 76\%$
   it is a decent accuracy 70-80%
   if it is more than 80% then it is good accuracy.

2) for zoo dataset

i) Did you perform any data preprocessing steps? if yes, what are they and why they necessary.

Yes, we use label encoder for converting categorical encoder data into numerical form.

It is necessary for further classification

ii) Were there any musing or inconsistent value in the dataset: how do handle them?

There are inconsistent value in the Animal names column as there is not standardizat of formate & units

iii) what does the confusion matrix tell you about the performance of your model?

all the Ne are in diagonal theryare No misidentificat of classes.

iv) No class type are misclassified

## Lab-4
### KNN (K- nearest neighbour)

| Person | Age | Salaryk | Target | distance | Rank |
|--------|-----|---------|--------|----------|------|
| A | 18 | 50 | N | $\sqrt{914}$ | 5 |
| B | 23 | 55 | N | $\sqrt{2169}$ | 4 |
| C | 24 | 70 | N | $\sqrt{1021}$ | 2 |
| D | 41 | 60 | Y | 40.44 | 3 |
| E | 43 | 70 | Y | 31.04 | 1 |
| F | 38 | 40 | Y | 60.07 | 6 |
| X | 35 | 100 | ? | | |

$\rightarrow \sqrt{(18-23)^2 + (50-55)^2}$

$= \sqrt{3^2 + 60^2}$

$= \sqrt{9 + 3600} = 60.07$

$\rightarrow \sqrt{(43-35)^2 + (70+40)^2}$

$= \sqrt{8^2 + 30^2}$

$= \sqrt{964} = 31.04$

$\rightarrow \sqrt{(41-35)^2 + (60-100)^2}$

$= \sqrt{6^2 + 40^2} = 40.44$

$$K = 3 \rightarrow N \ Y \ Y$$
$$= N$$

For the Iris dataset :
Testing :

→ how to choose K value?

testing multiple K values & computing their accuracy. The accuracy & error rate for K are compared and the most optimal K is selected

here K = 3

diabetes dataset

∴ what is the purpose of feature scaling? how to perform it?

feature scaling ensures all features contribute equally to the nearest neighbours Scaling is done so that features like glucose / age don't dominate the one with the smaller ranges.