# LAB - 03
## LOGUTIC REGRESSION

Q1. Consider a binary classification problem where we want to predict whether a student will pass or fail.

Given $a_0 = -5$ (intercept)

$a_1 = 0.8$ (coefficient)

(a) Write the logistic regression equation for this problem.

$$sigmoid(z) = P(x) = \frac{1}{1 + e^{-z}}$$

$$z = a_0 + a_1 x$$

$$z = -5 + 0.8x \qquad P(x) = \frac{1}{1 + e^{-(-5 + 0.8x)}}$$

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

$$Probability (x | pass) = \frac{1}{1 + e^{5 - 0.8(7)}} = \frac{1}{1 + e^{-0.6}}$$

$$= 0.6457$$

(c) Determine the predicted class (pass or fail) for this student based on a threshold.

If threshold = 0.5

$$P(x = 7) \quad \Leftrightarrow > 0.5$$

Student will ~~fail~~ [pass.]

Q2. Consider $z = [2, 1, 0]$ for three classes. Apply softmax function

$$\text{softmax}(z_k) = \frac{e^{z_k}}{\sum_{j=1}^{k} e^{z_j}}$$

$$\text{softmax}(z_1) = \frac{e^2}{e^2 + e^1 + e^0} = \frac{7.39}{7.39 + 2.72 + 1} \approx 0.665$$

$$\text{softmax}(z_2) = \frac{e^1}{e^2 + e^1 + e^0} = \frac{2.72}{7.39 + 2.72 + 1} \approx 0.244$$

$$\text{softmax}(z_3) = \frac{e^0}{e^2 + e^1 + e^0} = \frac{1}{7.39 + 2.72 + 1} \approx 0.095$$

Probabilities for three classes = $66.5\%$, $24.4\%$, $9.1\%$.

**q.** HR_comma_sep. csv

(i) which variables did you identify as having a direct and clear impact on employee retention ? why ?

variables such as satisfaction - level , average - monthly - hours , number - project , time - spend - company , salary .

- satisfaction level : strong negative correlation
- average - monthly - hours : strong positive correlation
- number - project : positive correlation
- time - spend - company : strong positive correlation
- salary : strong negative correlation

(ii) what was the accuracy of your logistic regression model ?

The accuracy was around (0.79) (79%)

**e.** Zoo. csv

(i) Did you perform any data preprocessing steps ? If yes, what were they ?

① Handling categorial data.

class — type → class _mapping

② splitting dataset
80% training, 20% testing

③ Drop irrelevant columns

Drop animal_name, class as it is irrelevant.

④ Feature Scaling.

standardscaler () to standardize numerical features
mean = 0
variance = 1

(ii) Were there any missing or inconsistent values in the dataset?

Soln    No missing values.

(iii) What does the confusion matrix tell you about the performance of your

The confusion matrix represents how well the model classifies differ classes.

Diagonal values : True Correct classification

Other values : Incorrect classification

(iv) Which class types were most frequently misclassified?

Soln    Class 3, 5, and 6 were misclassified.

Possible reasons
①  Feature similarity
②  Small sample size