

Lab 4

Logistic regression

- 1) Consider a binary classification problem where we want to predict whether students will pass or fail based on their study hours. The logistic regression model has been trained on data and learned the parameters as $a_0 = -5$ and $a_1 = 0.8$.

- 2) Write a logistic regression equation

$$P(y=1|x) = \frac{1}{1 + e^{-(a_0 + a_1 x)}}$$

where x = study hours

$$P(y=1|x) = \frac{1}{1 + e^{-(a_0 + a_1 x)}}$$

probability

$$P(y=1|x=7) = \frac{1}{1 + e^{-(a_0 + a_1 x)}}$$

$$= 0.6956$$

Given threshold = 0.5

$$P(\text{Student}) = 0.6956$$

$P > \text{threshold}$,

Student is in pass class

2) Given classes $\Rightarrow [2, 1, 0]$
 $C^2 \Rightarrow [C^2, e^1, e^0]$
 $\Rightarrow [7.389, 2.718, 1]$
 $E^2 \Rightarrow [11.107]$
 Probabilities $\Rightarrow [0.665, 0.219, 0.09]$

For dataset file "HR - Comm - attr.csv"

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

Key variables are
 \Rightarrow Spesification level \Rightarrow Strong negative correlation

\Rightarrow Time spent in company \Rightarrow Positive correlation

\Rightarrow Work accident \Rightarrow Negative correlation

\Rightarrow Salary \Rightarrow Lower salary leads to more leavers

\Rightarrow Department \Rightarrow Retention varies for departments

\Rightarrow Model accuracy \Rightarrow Accuracy is 78.58% which is concerning and indicates if on imbalanced dataset, many leaves are missed.

Class ~~represents~~ ^{is} Big and misclassified
once ~~it~~ ^{it} could be because
a) Limited amount of data
b) Feature similarity in Big and Good class

Yes

Class	Feature	Value	Label
0	81	81	0
0	22	22	0
0	80	80	0
0	02	02	0
0	05	05	0
0	08	08	0
0	01	01	0

Class	Feature	Value	Label
1	18	18	1
1	22	22	1
1	80	80	1
1	02	02	1
1	05	05	1
1	08	08	1
1	01	01	1