

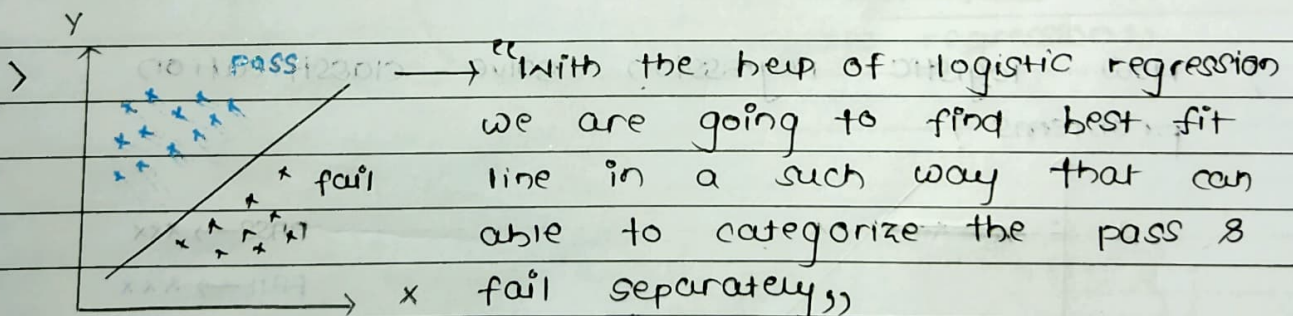
LOGISTIC REGRESSION

- > Logistic regression is supervised learning technique. It is used to predicting the categorical dependent variable using given set of independent variable.
- > used to solve classification problems
 - Binary classification
 - multiclass classification

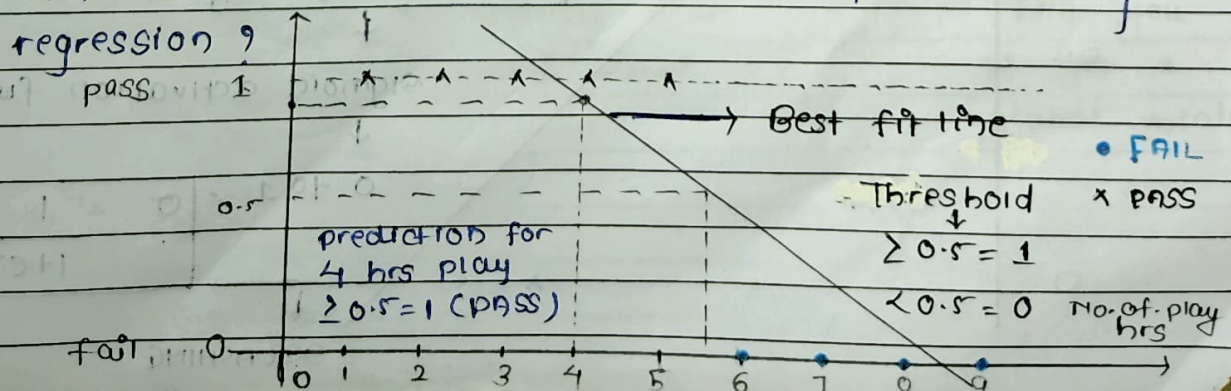
- > Binary classification → output = 2 categories
- > multiclass classification → output = greater than 2 categories.

> eg: Dataset

No. of play hours	pass / fail (y)
9 hrs	fail 0
8 hrs	fail 0
7 hrs	fail 0
6 hrs	fail 0
5 hrs	pass 1
4 hrs	pass 1
3 hrs	pass 1
2 hrs	pass 1



> can we solve this classification problem using regression?



if the student playing for 4 hours a new predicted datapoint ≥ 10.5 which means that student is pass.

In case of outlier we will get new best fit line that reduce our performance and result into wrong prediction.

Why we cannot used linear regression for classification?

Ans: Linear regression deals with continuous values not with discrete values.

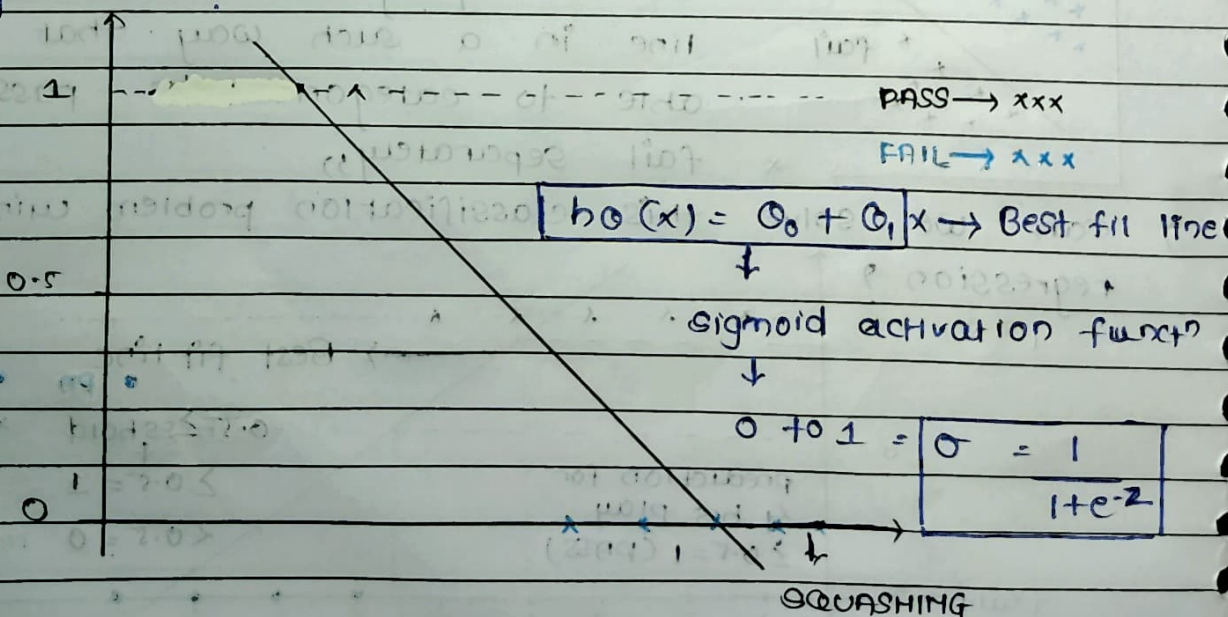
(2) Shift in threshold in case of outliers that can affect the performance of model result into wrong prediction.

3. The outcome comes to > 1 & < 0 .

to solve this problem we will use logistic regression.

$[0 \text{ to } 1] \Rightarrow$ Squashing technique.

How logistic regression solve classification problems?



$$\therefore h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}}$$

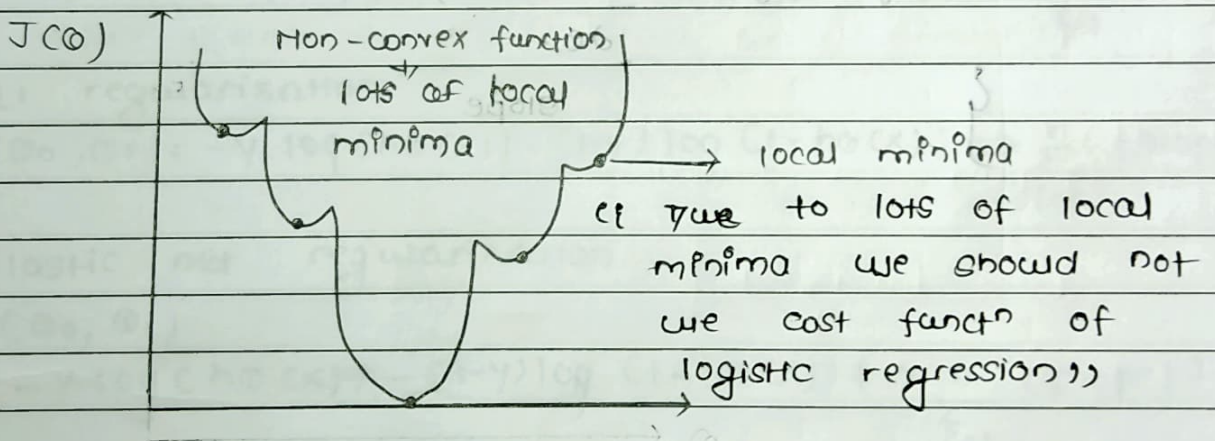
squashing the best fit line.

$$\therefore h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}} \Rightarrow \text{logistic regression}$$

> Cost function for logistic regression

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$

$$\therefore h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}} \quad \left. \vphantom{\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}}} \right\} \text{value 0 to 1}$$



> log loss function

$$J(\theta_0, \theta_1) = -y_i \log(h_{\theta}(x_i)) - (1 - y_i) \log(1 - h_{\theta}(x_i))$$

$y_i \rightarrow$ real or actual data

$h_{\theta}(x_i) \rightarrow$ predicted value = \hat{y}

"help you to get a global minima"

if $y_i = 1$

$$J(\theta_0, \theta_1) = \sum -\log(h_{\theta}(x_i)) \quad \text{if } y = 1$$

$$J(\theta_0, \theta_1) = \sum -\log(1 - h_{\theta}(x_i)) \quad \text{if } y = 0$$

> final aim: minimize cost function $J(\theta_0, \theta_1)$ by changing θ_0, θ_1

> Convergence algorithm:

Repeat untill convergence

{ learning rate α

$$\theta_j = \theta_j + \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

$\frac{\partial}{\partial \theta}$

slope

LOGISTIC REGRESSION WITH REGULARIZATION

- > Regularization \rightarrow reducing overfitting
- > Applying regularization techniques

$$J(\theta_0, \theta_1) = -y \log(h\theta(x)) - (1-y) \log(1-h\theta(x)) + L_2$$

$$J(\theta_0, \theta_1) = -y \log(h\theta(x)) - (1-y) \log(1-h\theta(x)) + L_1$$

\therefore L_2 regularization \rightarrow reducing overfitting

\therefore L_1 regularization \rightarrow feature selection.

$$J(\theta_0, \theta_1) = -y \log(h\theta(x)) - (1-y) \log(1-h\theta(x)) + L_1 + L_2$$

\therefore L_2 reg + L_1 reg \Rightarrow Elastic net

- > L_2 regularization

$$J(\theta_0, \theta_1) = -y \log(h\theta(x)) - (1-y) \log(1-h\theta(x)) + \lambda \sum_{i=1}^n (\text{slope})^2$$

- > L_1 regularization

$$J(\theta_0, \theta_1) = -y \log(h\theta(x)) - (1-y) \log(1-h\theta(x)) + \lambda \sum_{i=1}^n |\text{slope}|$$

- > Elastic net regularization

$$J(\theta_0, \theta_1)$$

$$= -y \log(h\theta(x)) - (1-y) \log(1-h\theta(x)) + \lambda \sum_{i=1}^n (\text{slope})^2 + \lambda \sum_{i=1}^n |\text{slope}|$$

PERFORMANCE METRICS, ACCURACY, PRECISION RECALL AND F-BETA.

1.) Confusion matrix

	1	0	Actual (Y)	eg: Dataset	predicted
1	3	2	x_1	x_2	y
0	1	1			\hat{y}
predicted (ŷ)					

> Notation

1 - P, 0 - N

1	TP	FP
0	FN	TN

$$\frac{TP + TN}{TP + FP + FN + TN} \rightarrow \text{Model Accuracy}$$

predicted = $\frac{3 + 1}{3 + 2 + 1 + 1} = 41.67\%$

TP, TN \rightarrow Correct match, FP, FN \rightarrow wrong match

> eg:

Dataset \rightarrow Imbalance dataset

1000 datapoint $\left\{ \begin{array}{l} 900 \rightarrow 1 \\ 100 \rightarrow 0 \end{array} \right\}$ Imbalance

confusion matrix

Trumb model $\rightarrow 1 \Rightarrow$ "we get

90% Accuracy 1 0 \rightarrow Actual

Accuracy 1	900	100
0	0	0

$$= \text{Accuracy} = 90\%$$

"if the accuracy is 90%, it is not sufficient model is

not good. To overcome this problem

we can use recall & precision,

2.) precision

> Measure the accuracy of positive prediction,

> $\text{precision} = \frac{TP}{TP+FP}$ } out of all the actual values how many are correctly predicted.

3.) Recall

> Number of correct post prediction made out of all post prediction that could have been made.

> $\text{Recall} = \frac{TP}{TP+FN}$ } out of all the predicted value how many are correctly predicted with actual value.

NOTE: > Whenever FP (False post) is important we will use precision. (FP) ↓ reducing.
 > Whenever FN (False neg-) is important we will use recall. (FN) ↓ reducing.

Usecase 1: Spam classification (precision) → FP is imp

Text → Model → Spam / Not Spam

1 0 → Actual

1	TP	FP
0	FN	TN

↓

predicted

TP(1)
 Mail (Spam) → Model → Spam } Accurate
 Mail (Not Spam) → Model → Not Spam }

TN(0)

Mail (Not Spam) → Model → Spam } wrong prediction
 ↓
 FP → (imp)

Mail (Spam) → Model → Not Spam } wrong prediction
 ↓
 FN

(Blunder)

Usecase 2: FN is important (Recall)

> To predict person has diabetic or not

	Diabetic	No diabetic
Diabetic	TP	FP
No diabetic	FN	TN

TP \Leftarrow Actual \rightarrow diabetic \rightarrow model \rightarrow diabetic \rightarrow correct

TN \Leftarrow Actual \rightarrow No diabetic \rightarrow model \rightarrow no diabetic \rightarrow correct

FP \Leftarrow Actual \rightarrow No diabetic \rightarrow model \rightarrow diabetic \rightarrow wrong

FN \Leftarrow Actual \rightarrow diabetic \rightarrow model \rightarrow No diabetic \rightarrow wrong

(Blunder)

4.) F-Beta Score

$$\therefore F\text{-Beta Score} = (1 + \beta^2) \frac{\text{precision} \times \text{Recall}}{\text{precision} + \text{Recall}}$$

i.) If FP and FN are both important

$$\beta = 1$$

$$F1 \text{ score} = \left\{ (1+1) \frac{P \times R}{P+R} \right\} \Rightarrow \text{Harmonic Mean}$$

ii.) If FP is more important than FN

$$\beta = 0.5$$

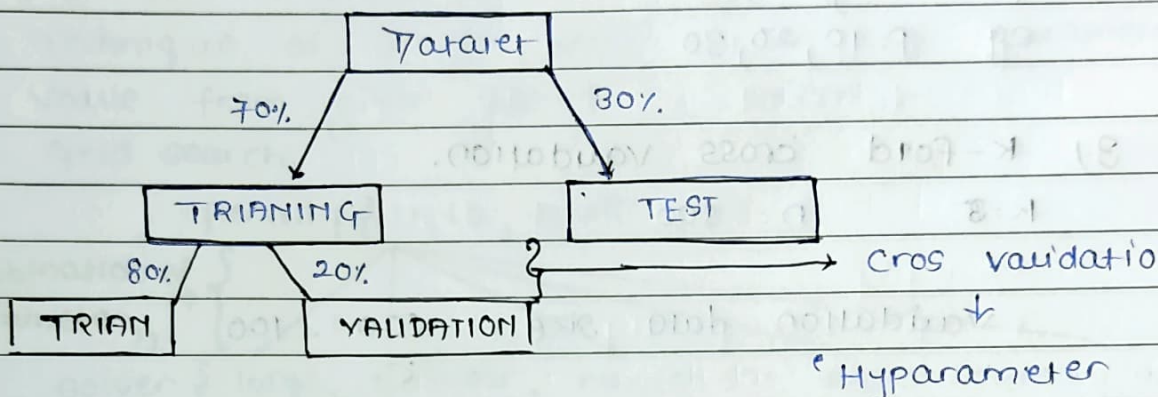
$$\left\{ F_{0.5} \text{ Score} = (1 + 0.25) \frac{P \times R}{P+R} \right\}$$

iii.) If FN is more important than FP

$$\beta = 2 : \left\{ F_2 = (1+4) \frac{P \times R}{P+R} \right\}$$

CROSS VALIDATION AND ITS TYPES

eg: consider dataset with 1000 datapoints



random state = 100
 = 42
 → 85% accuracy } getting different accuracy because we changing the random state value
 75%
 65%

TRAINING ⇒ CV=5	Exp1	Train	validation	Acc1
	Exp2	Train	validation	Acc2
	⋮			
	Exp	Train	validation	Acc

mean of accuracy

TYPE OF CROSS VALIDATION:

1.) Leave one out cross validation (LOOCV)

training data → 500 datapoints
 validation ↓

Exp1 =	1	499 training data	⇒ Acc1 = 85%	} Avg of all accuracy
	validation			
Exp2 =	1	2	498 training data	⇒ Acc2 = 80%

Disadvantage:

- ① Time complexity is huge
- ② Model tends to overfit. ⇒ Train Acc ↑
 New data ↓

2.) leave p out cross validation

p \Rightarrow hyperparameter

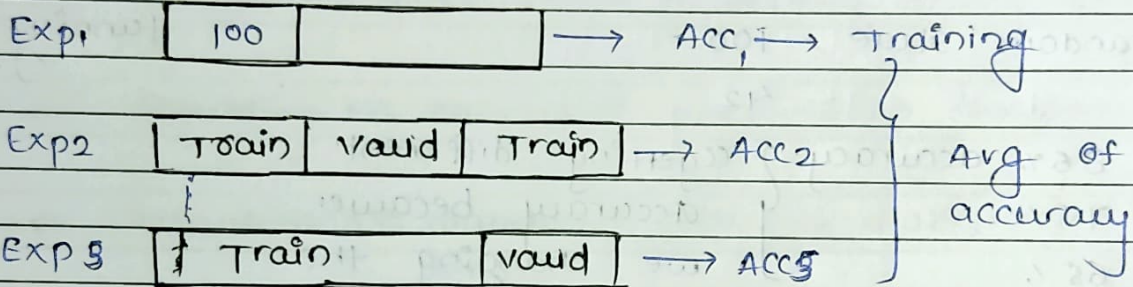
eg: p=10, 20, 30

3.) K-fold cross validation.

K=5

n = 500

validation data size $\frac{500}{5} = 100$



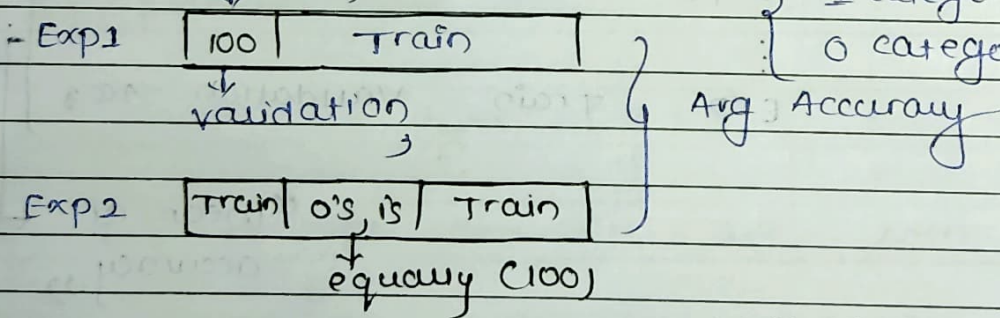
4.) Stratified K Fold cv \rightarrow { Imbalanced dataset }

0's, 1's K=5

\downarrow equally

n = 500

{ 1 category = 350
0 category = 150 }



HYPERMETER TUNNING WITH CROSS VALIDATION

> Finding the best parameter while training the models.

1.) Grid search cv [Grid search + cv]

> Technique of finding optimal or best parameter value from given set of parameter.

> Grid search

penalty { L_1 , L_2 , elastic, None }

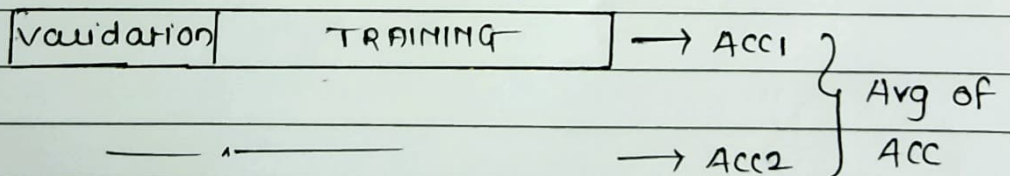
Combination of parameter,

solver { lbfgs, liblinear, newton-cg, newton-cholesky, saga, saga, }

"Grid search cv making sure that go ahead and checking every combination try to find out best combination with high or best accuracy,"

> let's assume that L_1 & lbfg is best parameters
 logistic regression

↓
 cross validation [k fold cv]



> Disadvantage

① Time complexity increases with huge dataset for training the model.



2.) Randomized search cv

> Randomly select the best set of parameters.

> eg: $n_iter = 10$, $cv = 5$
 10 different combination + $cv = 5 \Rightarrow 50$

↓

select best parameter

> Advantage

① Time complexity decreases ↓