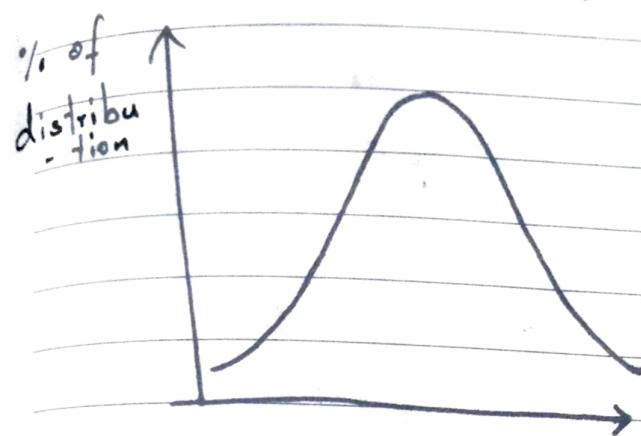
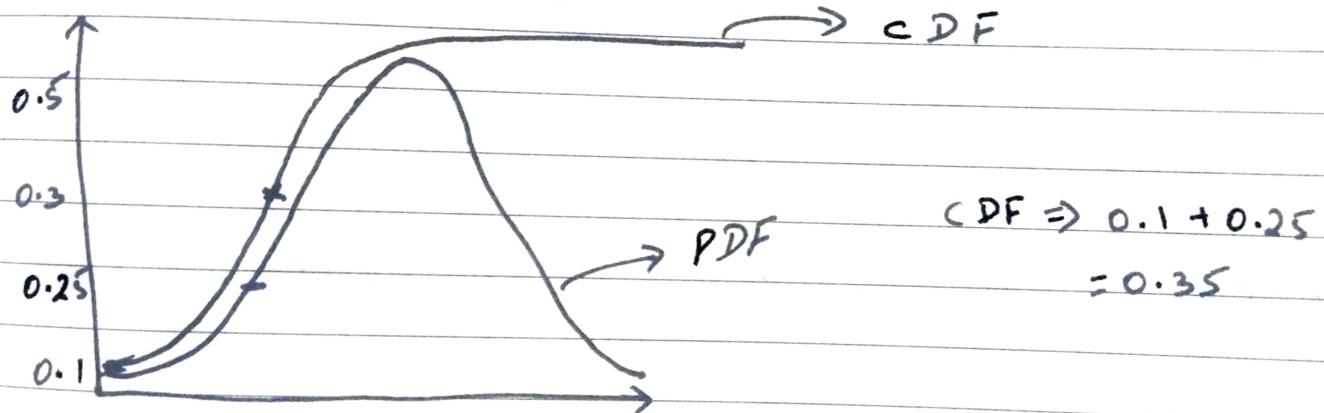


Probability Density function (PDF)



Cumulative Density Function

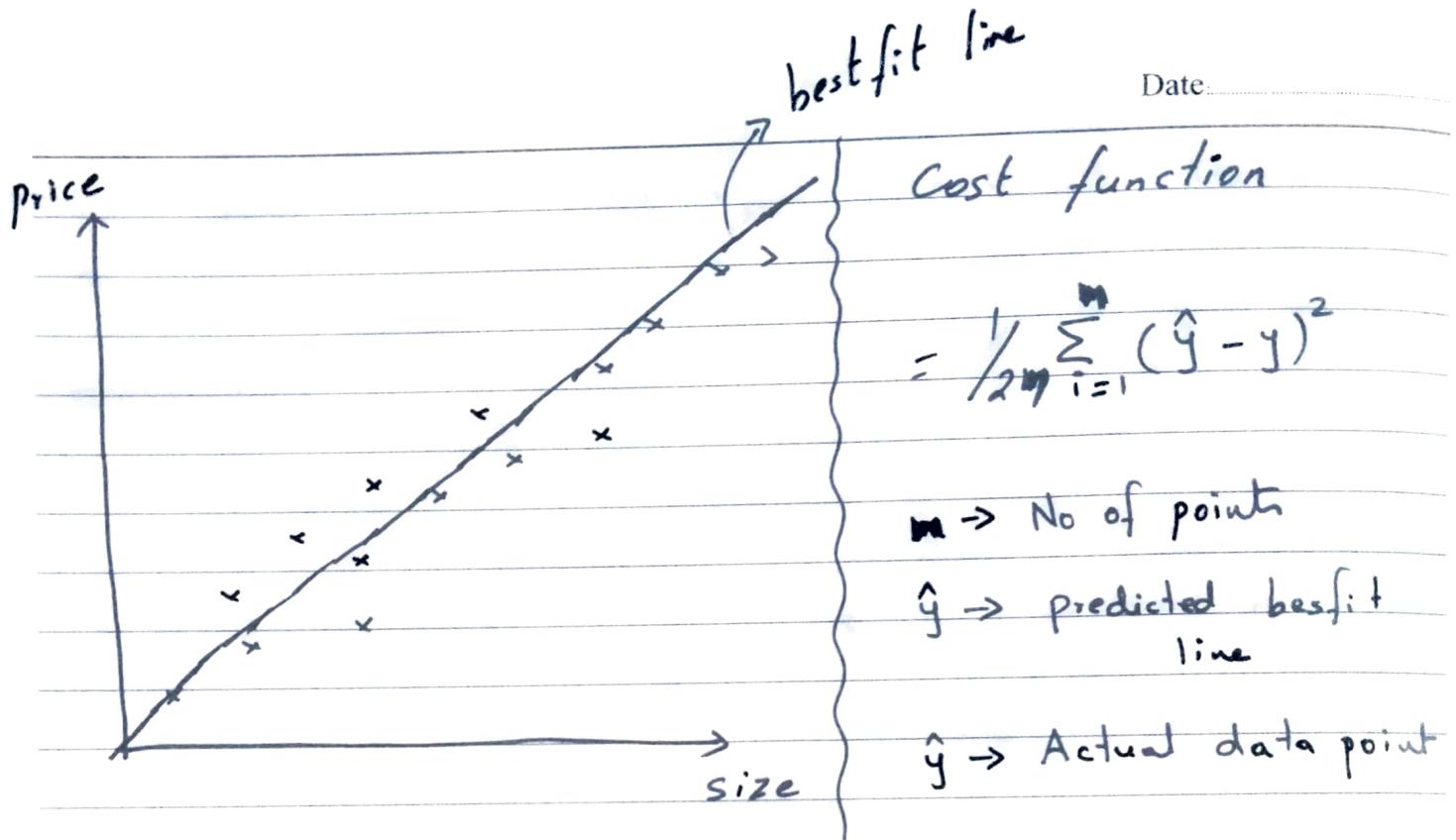


Linear Regression indepth Math

$$y = mx + c \rightarrow \text{Best fit line}$$

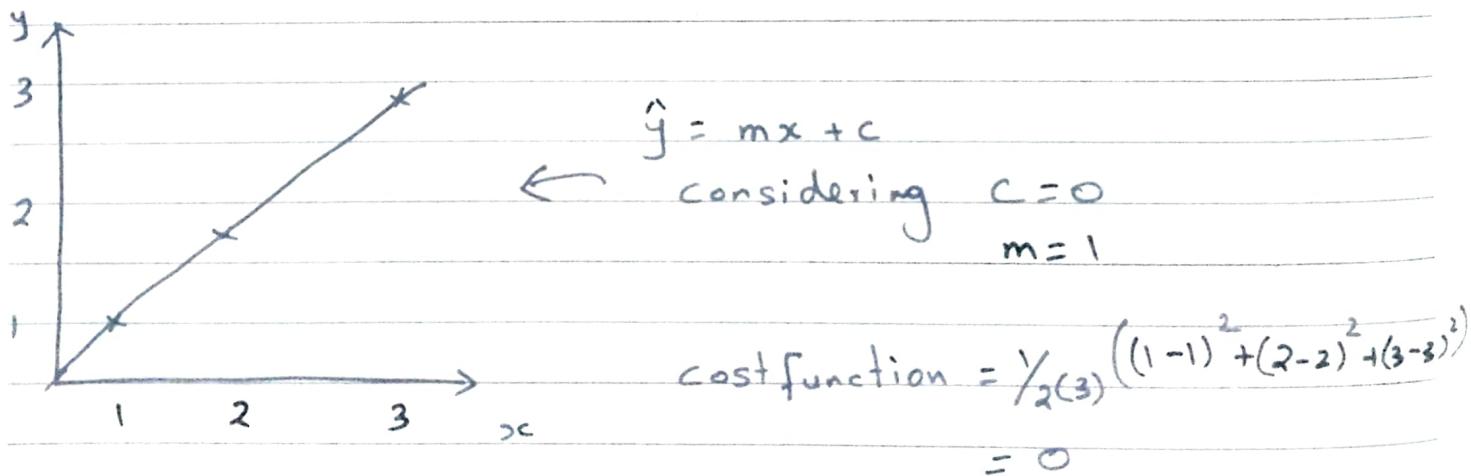
$m \rightarrow$ slope

$c \rightarrow$ intercept

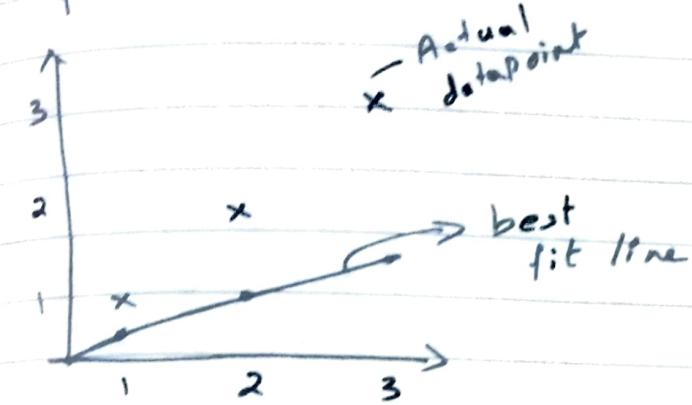


best fit line is actually a \hat{y}

cost function should be minimum



if $m = 0.5, c = 0$

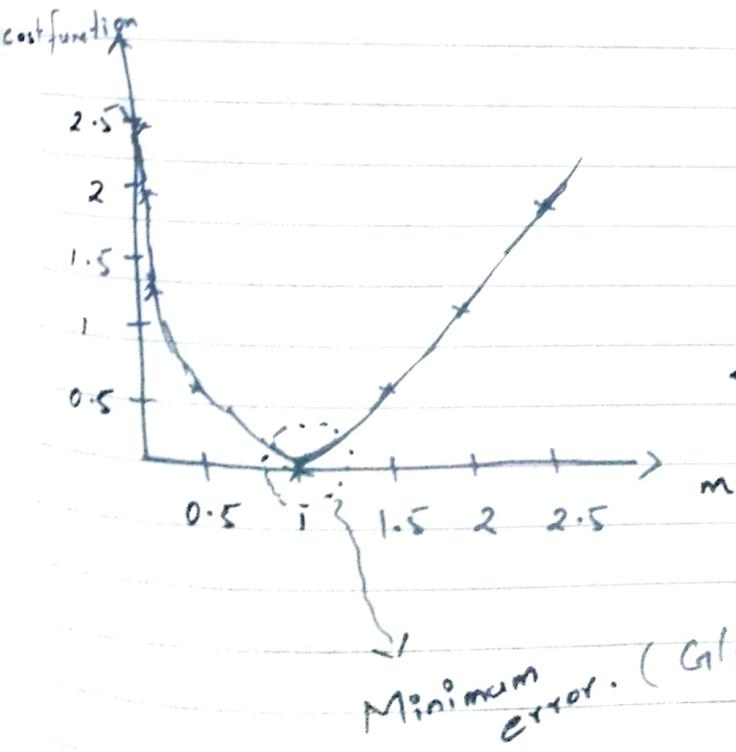


$$\hat{y} = (0.5)(1) + 0 \\ = 0.5$$

cost function

$$= \frac{1}{2}(3) \times [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2] \\ = 0.58$$

Gradient Descent



for different slope (m) selection the cost function changes.

← this will help us to selecting (m) value slope value

Minimum error. (Global minima)

in order to move down to reach global minima

we need to write

Convergence theorem

$$m = m - \underbrace{\left(\frac{\partial m}{\partial m} \right)}_{\hookrightarrow \text{slope}} \times \alpha \quad \begin{array}{l} \xrightarrow{\text{learning rate}} \text{leaning} \\ \xrightarrow{\text{small value}} \text{rate should be small value} \end{array}$$

if the slope is negative

$$\begin{aligned} m &= m - (-\text{ve}) \times \alpha \\ &= m + \alpha \\ &= \text{greater value} \\ &\quad \text{less value} \end{aligned} \quad \begin{array}{l} \xrightarrow{\text{Negative slope}} \\ \text{Left} \rightarrow \text{greater} \\ \text{Right} \rightarrow \text{less} \end{array}$$

if the slope is positive

$$\begin{aligned} m &= m - (+\text{ve}) \times \alpha \\ &= m - \alpha \\ &= \text{less value} \\ &\quad \text{greater value} \end{aligned} \quad \begin{array}{l} \xrightarrow{\text{positive slope}} \\ \text{Left} \rightarrow \text{less} \\ \text{Right} \rightarrow \text{greater} \end{array}$$

in each iteration the cost function will reduce and try to attain global minima.

Ridge by Lasso Regression

Ridge if the slope is huge (steep), we try to penalizing with the help of ridge regression with the below formula

$$\text{cost function} + \lambda \frac{(\text{slope})^2}{m^2} \quad \left[\begin{array}{l} \lambda \text{ is btw} \\ 0 \text{ and } \infty \end{array} \right]$$

It will help to reduce the steep of the slope.

By this we can reduce overfit condition

Lasso

$$\text{Cost function} + \lambda |\text{slope}|$$

$$\downarrow$$

$$\text{Cost function} + \lambda \times |m_1 + m_2 + m_3|$$

Lasso performs well when we have many features.

$$R^2 \text{ score} = 1 - \frac{SS_{\text{res}}}{SS_{\text{mean}}} \quad (\text{sum of residual})$$

$$SS_{\text{res}} = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$$

$$SS_{\text{mean}} = \frac{1}{n} \sum_{i=1}^n (y - \bar{y})^2$$

To check Multi collinearity
 Statsmodel OLS summary interpret

1. We need to check $\beta_0, \beta_1, \beta_2, \beta_n$

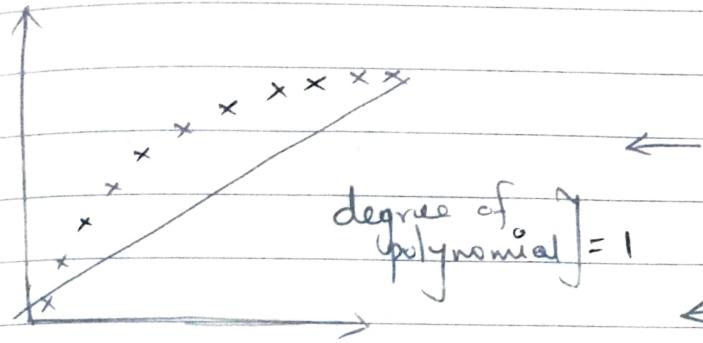
2. R^2 & Adj R^2 should near 1

3. Std error should be very very less

4. P value should less than 0.05

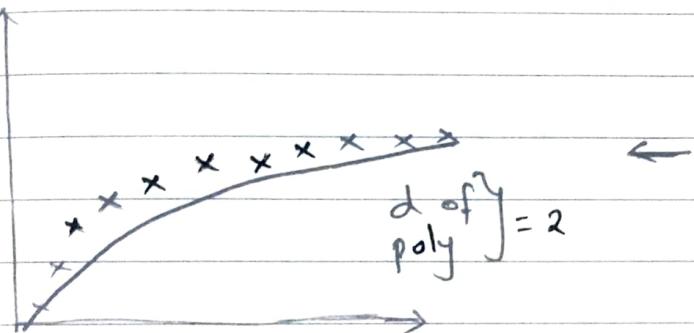
Bias by Variance

In terms of
Train data

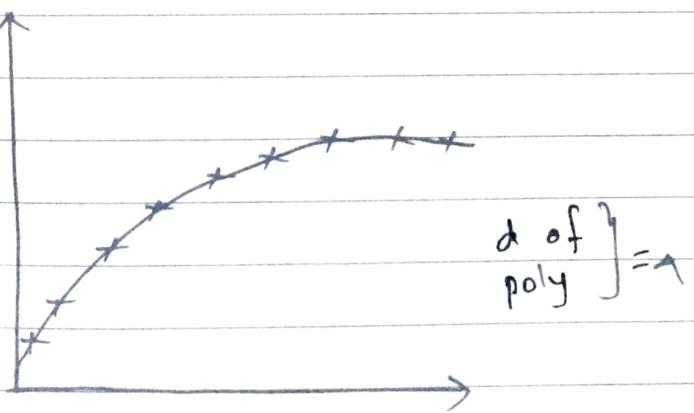


When the error is
high it is underfitted
model

high bias by high
variance



low bias by low Variance
good model to choose

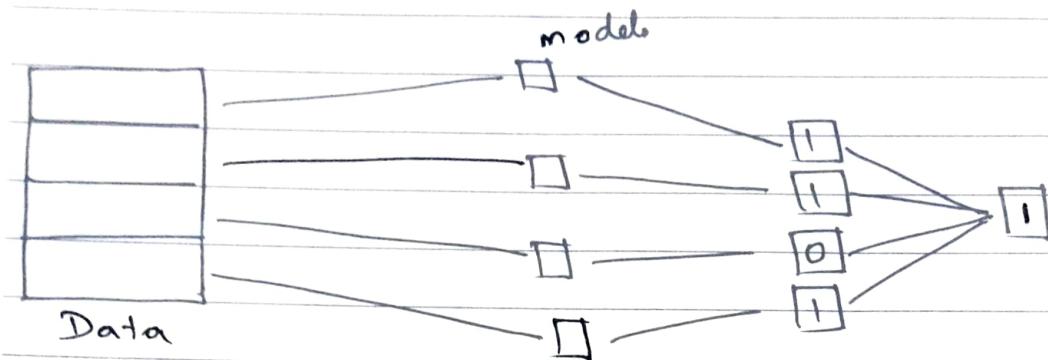


When the error is
very very low, it is
Overfitted model

low bias by high
variance.

Consider bias \rightarrow error of train data
Variance \rightarrow " " test data.

Bootstrap Aggregation



splitting data and passing into various models and getting the highly occurred result

e.g. Random Forest

Decision tree

low Bias, high variance

Random forest

low bias, low variance.

R^2 by Adj R^2

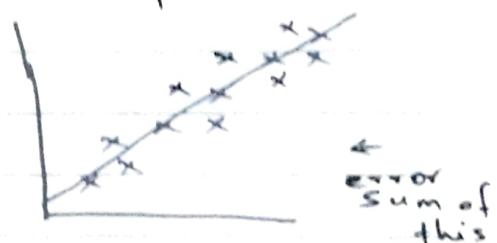
ranges (0 to 1)

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{mean}}} \quad (\text{sum of residual})$$

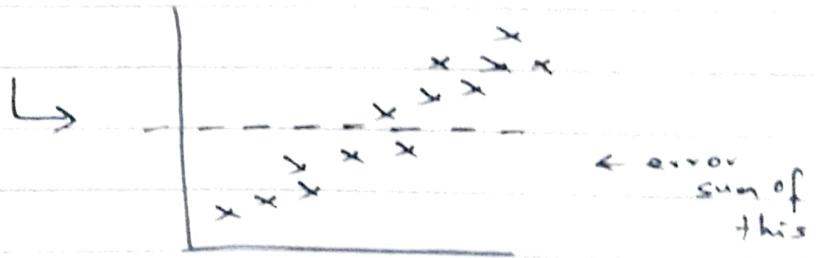
(sum of mean)

Sum of residual (or) error

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \hat{y} \rightarrow \text{Predicted points.}$$

Sum of mean (or) total average

$$\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad \bar{y} \rightarrow \text{mean}$$

 R^2 can become below 0

When the predicted best fit line is worse than Average line it will become zero. It means that is the worst model.

Adj R^2

When we are adding more features, the R^2 value will automatically increase, coz it will try to assign some coefficient value, It happens always.

So we need to use Adj R^2

$$\text{Adj } R^2 = 1 - \frac{(1 - R^2)(N-1)}{N-p-1}$$

(No of independent variables)

$p \rightarrow$ No of predictors

$N \rightarrow$ No of sample size

It will penalize the attributes that are not correlated with target.

Adj R^2 always less than or equal to R^2

Hypothesis testing

We use this to evaluate 2 mutually exclusive statement on Population using sample data.

steps

① Make initial assumption (H_0)

mutually exclusive
2 events can't occur at the same time eg tossing a coin

② Start Collecting data

③ decide to Reject or Not Reject (H_0)

$H_0 \rightarrow$ Null hypothesis

	H_0	H_1
Not Reject	OK TP FP	Type 2 Error FN TN
Reject	Type 1 Error	OK

Type I Error

Actually True by test result false
Not Reject Reject

Type II Error

Actually False by test result true
Reject Not Reject

T-Test, chisquare test, ANOVA Test

if we are considering only one categorical column at that time we must use one sample proportion test

eg:-	Gender	AGE	Weight
	M	45	..
	F	32	...
	M	17	...
	M	29	..
	F	19	..
	F	7	--

if we take this to test

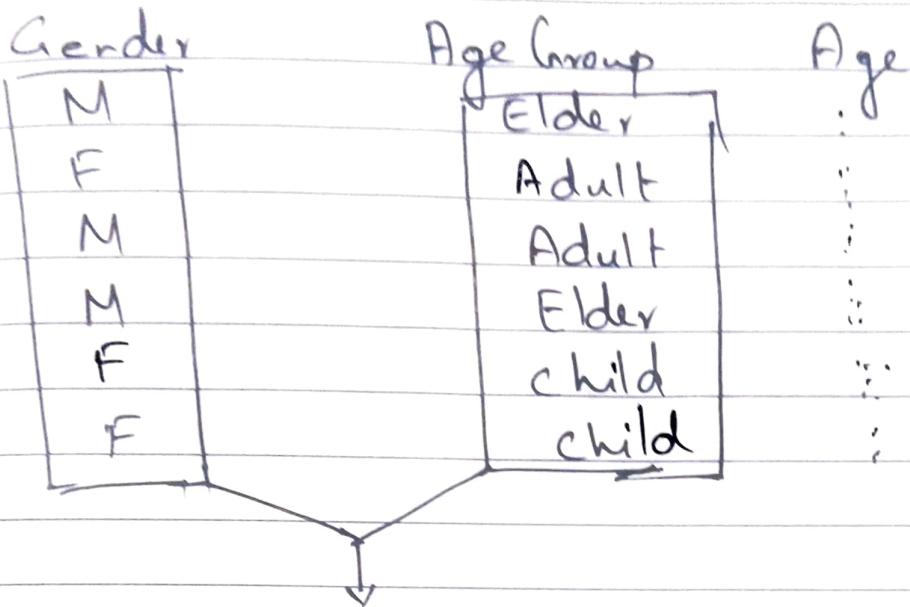
$H_0 \rightarrow$ No difference in proportion

$H_1 \rightarrow$ There is a proportion

$P_{val} \rightarrow \leq 0.05$

We will if P_{val} less than 0.05 or 5% reject H_0 by accept H_1 .

if we consider 2 categorical features for test then we should use Chi-square test.



$H_0 \rightarrow$ based on age group No diff in proportion

$H_1 \rightarrow$ difference is there

for One continuous variable

T-Test

for 2 continuous variable

T-Test / Correlation Test

ANOVA Test

One^{or} more Categorical by Numerical

features but the Categorical feature have more categories then we should use ANOVA test.

Metrics in classification problem

Confusion Matrix

FPR (Type I Error)

FNR (Type II Error)

Recall (TPR, Sensitivity)

Precision (Positive Pred val, Specificity)

Accuracy

F Beta

Cohen kappa

ROC AUC

PR Curve

If we have balanced target feature we can consider "Accuracy"

If we don't have balanced target then we can consider Recall, Precision, F_{beta}.

Confusion Matrix

		Actual	
		0	1
Predicted	0	TP	FP
	1	FN	TN

→ Type I Error

↓

Type II Error

Always we should reduce Type I, II Errors.

$$FPR = \frac{FP}{FP+TN}$$

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

Recall / TPR / Sensitivity

$$\frac{TP}{TP+FN}$$

Precision / Specificity /
positive Prod value

$$\frac{TP}{TP+FP}$$

Out of all +^{ve} values
how many +^{ve} actually
predicted correctly

Out of all predicted +^{ve}
values, how many are actually
positive values.

Whenever FP is much more important
go with Precision (eg:- email spam
detection Problem)

Whenever FN is much more important
go with Recall (eg:- Cancer detection)

If Both FP by FN is more important
then we go with F_{Beta} score.

$$F_{\beta\text{eta}} = \frac{(1+\beta^2) \cdot \text{Precision} \cdot \text{Recall}}{\beta^2 \cdot \text{Precision} + \text{Recall}}$$

$$\beta = 1 \rightarrow F_1 \text{ score}$$

$$\beta = 0.5 \rightarrow F_{0.5} \text{ Score}$$

$$\beta = 2 \rightarrow F_2 \text{ Score}$$

When Both FP by FN are important $\rightarrow F_1$ Score

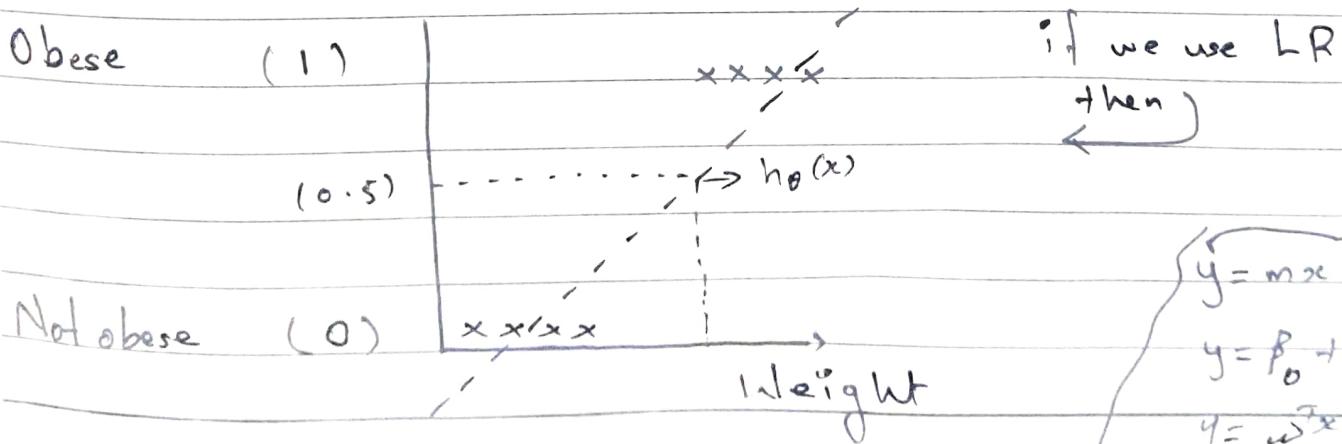
FP have more impor than FN $\rightarrow F_{0.5}$ Score

FN have more impor than FP $\rightarrow F_0$ greater than F_1

Logistic Regression

* Binary Classification

* Multiclass Classification



Reasons for not using Linear Regression for Classification

Whenever we hv lot of outliers the best fit line completely gets deviated

if we get bestfitline greater than 1 len than 0 wrt datapoint at that time we have no option.

We have to use logistic regression to overcome these problems.

Logistic Regression

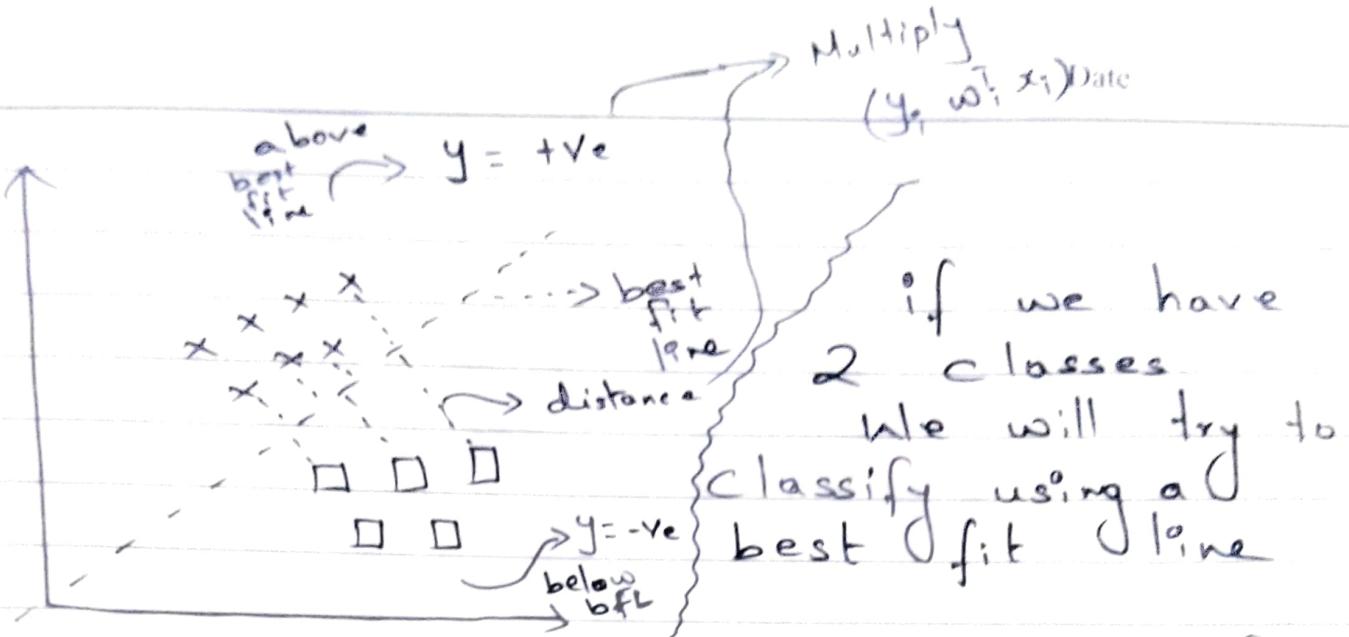
$$\text{Cost function} = \sum_{i=1}^n y_i w_i^T x_i \rightarrow \text{it should be Max}$$

Optimizer

Sigmoid Function

$$z = \sum_{i=1}^n f(y_i, w_i^T x_i)$$

$$= \frac{1}{1+e^{-z}} \rightarrow 0 \text{ or } 1$$



then we will find the distance between the plane (or) best fit line and all the data points

Intercept is 0 in this case bcoz it passes through origin

$$\begin{aligned} y &= mx + c \\ y &= w^T x + b \rightarrow \text{any} \\ y &= \beta_0 + \beta_1 x \end{aligned}$$

$$w^T x + b \rightarrow \text{Intercept}$$

$$w^T x + 0$$

\hookrightarrow it will become distance b/w plane & datapoint

if it is above the plane +ve distance
below " " -ve "

multiply with y_i

$$(+ve) \cdot (+ve) = +ve$$

$$\sum_{i=1}^n (y_i w^T x_i)$$

$$(-ve) \cdot (-ve) = +ve$$

\hookrightarrow cost function

Cost function is Z

sigmoid $\frac{1}{1+e^{-Z}}$ \rightarrow it will range between
only 0 & 1

For Multi class classification

Change the parameter

multiclass : OVR