

Statistical Modeling

40%	F. exam
20%	Role
20%	Project
20%	Mid course

Rajan Sir :- The professor is interested in building "Statistical Models" (Data Science).

→ Teaches beyond Stats, and uses the statistical ideas to build models (50% of data science).

Diff Equation

Integration

1st & Second order D.E

VECTORS

Matrix & determinants (Important)

Basic topics for this module.

Matrices & Determinants
 $(O = \text{diag})$ - Inverse of P

Matrix multiplication with transpose of P
if $A = P^{-1}B$ then $B = A^T P^T$
Transpose of transpose is original matrix
If P is orthogonal then $P^T = P^{-1}$
Orthogonal matrix with transpose of P

Co-relation Coefficient

To know the relation between two variable we use "CO-Variance"

$$\text{Variance} (\sigma^2) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (\text{for 1 variable})$$

$$\text{Co-variance} = S_{xy}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- Covariance is There will be a positive relationship if both x & y increases or decreases together ($S_{xy}^2 = +ve$)
- There will be negative or opposite relationship if both x & y changes oppositely ($S_{xy}^2 = -ve$)
- There is no relationship if x changes, then y does not change. ($S_{xy}^2 = 0$)

(*) The problem is the covariance value b/w two Variable 'x' & 'y' Changes of the units of 'x' & 'y' are changed, so we cannot predict how strong is the relationship b/w the two variables with diff values of Co-variance

Therefore,

we use standard score for each variable and calculate the co-variance.

This co-variance in terms of z-score is called as Co-relation Coefficient (γ).

$$\gamma = \frac{1}{n-1} \sum_{i=1}^n (Z_x)(Z_y)$$

$$-1 \leq \gamma \leq 1$$

$\gamma \approx 1 \rightarrow$ Strong positive relationship

$\gamma \approx -1 \rightarrow$ Strong negative relationship

$\gamma \approx 0 \rightarrow$ no relationship

Simple Linear Regression :-

- ① The Variable which we want to predict on, is called Predicted Variable / Outcome Variable / Target Variable.
- ② The Variables which are used to predict the Predicted variable is called Predictor Variable / Regressor Variable.

Since we use only one variable to predict the Predicted variable. \therefore Simple.

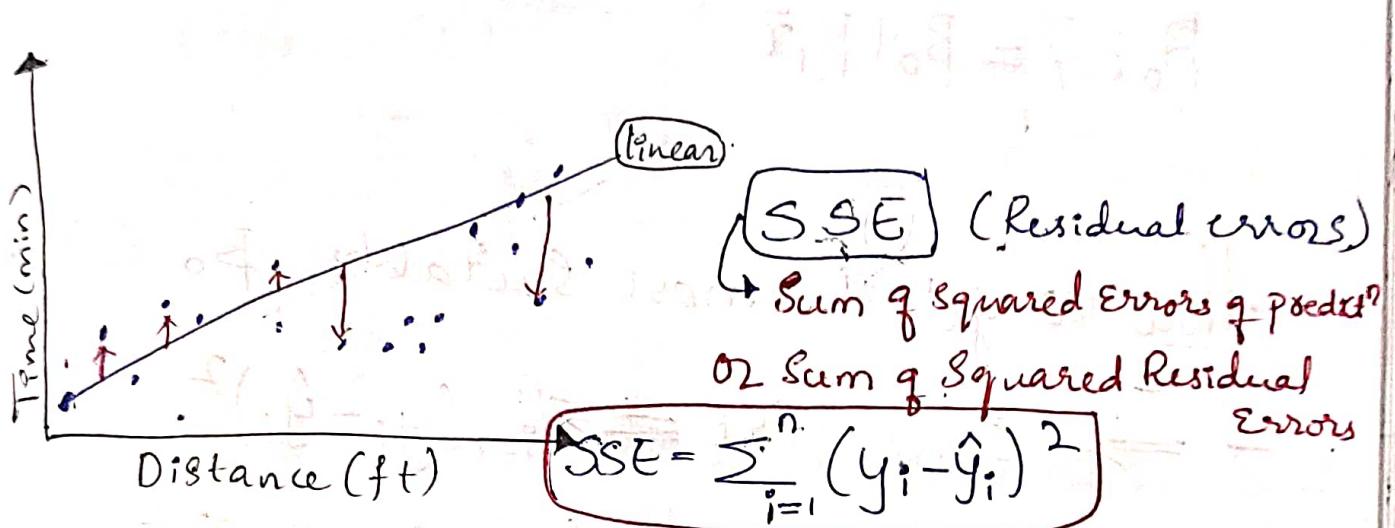
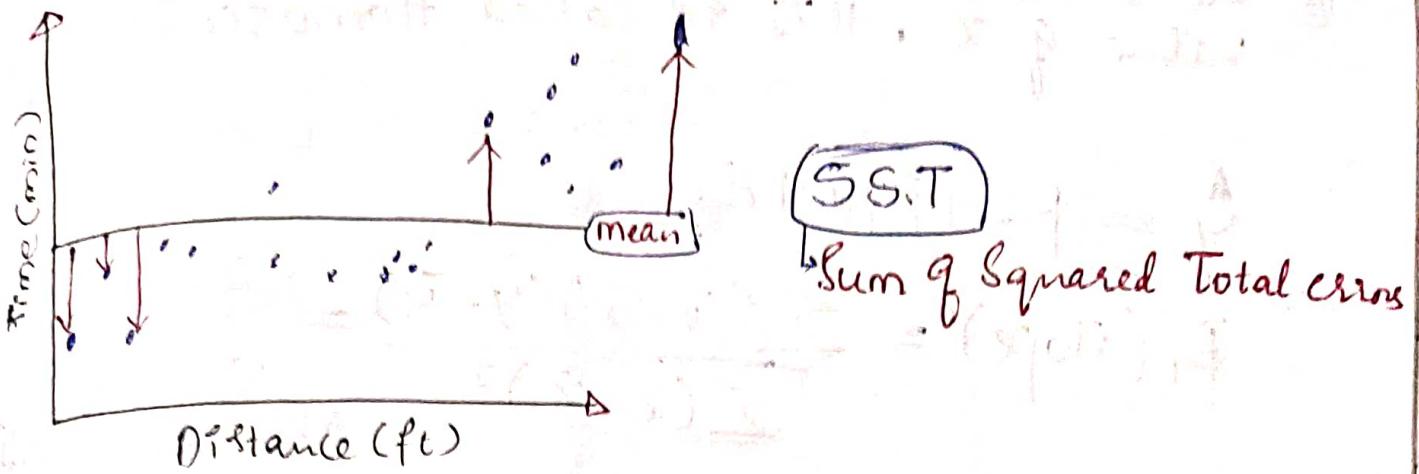
The value of m The performance of the model line will only results in the linear change \therefore linear

$$Y = mx + c$$

\hat{Y} \rightarrow predicted Value

\bar{Y} \rightarrow mean value

R^2 \rightarrow Says how better is the model.



① The best line would be the line with the least SSE
this can be achieved by systematically changing m & c.

In Data Science field we ~~can't~~ use β_0 & β_1
in place of m & c.

$$\Rightarrow \hat{y} = \beta_0 + \beta_1 x \rightarrow \text{prediction}$$

$$y = \beta_0 + \beta_1 x + \epsilon \rightarrow \begin{cases} \text{Data model} \\ \text{measured value} \end{cases}$$

\downarrow Errors.

We assume that the errors are coming from Normal distribution with mean zero, with some S.dev

④ The error distribution is independent of the value of x , this is called Homoscedasticity.

$$y = \beta_0 + \beta_1 x + \epsilon$$

term = 8

$$\beta_1 \text{ (slope)} = \frac{\sum ((x - \bar{x})(y - \bar{y}))}{\sum (x - \bar{x})^2} = \frac{s_{xy}^2}{s_x^2} = \frac{S_{xy}}{S_x}$$

Intercept
Covariance

Intercept
Correlation

$$\beta_0 + \bar{y} = \beta_0 + \beta_1 \bar{x}$$

How to get the most suitable β_0 & β_1 values

$$\Rightarrow \text{Since } SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\Rightarrow SSE = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x))^2$$

$$\frac{d(SSE)}{d\beta_0} = 0$$

$$\frac{d(SSE)}{d\beta_1} = 0$$

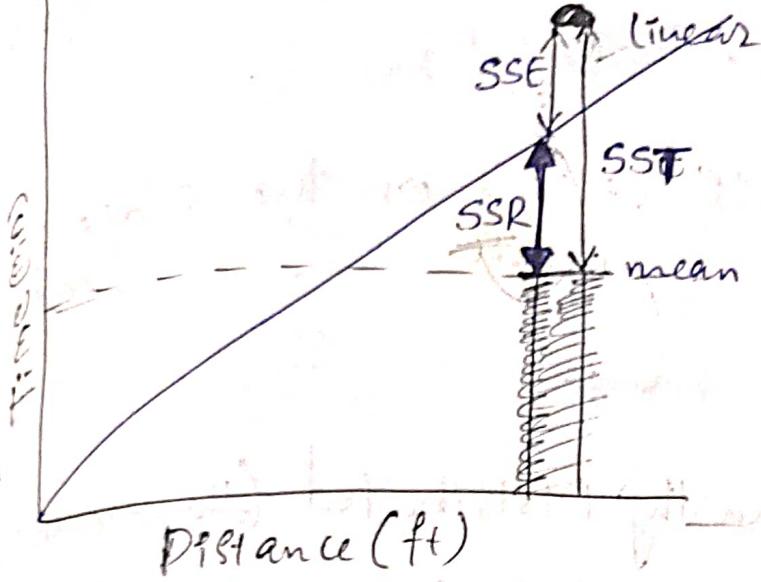
diff SSE w.r.t β_0

diff SSE w.r.t β_1

⑤ β_0 has no importance if the x does not have 0 as the data in its + 0.8

(make y intercept = 0 during this situation while drawing trend line)

SSR :- Sum of Squared Regression Errors



$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\underline{SSR} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$SST = SSE + SSR$$

$$I = \frac{SSR}{SST} + \frac{SSE}{SST}$$

$$\boxed{\frac{SSR}{SST} = I - \frac{SSE}{SST} = R^2}$$

R^2 Close to 1 = Good model

$R^2 \rightarrow$ How much of the total error is explained by the model

↳ what proportion of errors that has been rectified by the model

SSR \rightarrow Explained errors

SSE \rightarrow Unexplained errors.

$$\therefore \underline{SST = Ex. err + Unex. err}$$

To quantify the model is good or not

① R^2 value.

② Hypothesis testing on the slope of the model

Assumptions

- * Errors are Normally Distributed. (Homoscedasticity)
- * Every measurement is independent of every other measurement
- * Homoscedasticity
- * No Influential outliers.

R Studio

$$Q = \frac{322 + 322}{122} = 1$$

dt LR \leftarrow lm(dt\$time ~ dt\$distance.)

Linear model.

Built time regressed upon distance.

Now,

dt LR Shows the Co-efficients:-

i.e. B_0 & B_1 Values

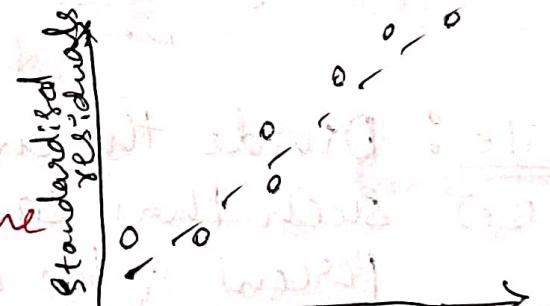
- ④ Summary(dtLR) :-
- Residuals :- Test whether it is Skewed or not
 - Co-efficients :-

Estimate	std err	t-value	P-value
			=
 - Residual std error is 6.909 with this we accept or reject null hypothesis
 - Multiple R-squared :-
 - F-statistic :-

- ⑤ plot(dtLR) :-
- Residuals VS Fitted (This shows if there is a relationship in the diff measurement. if relationship exists, then the ② assumptions will fails)

→ Normal Q-Q

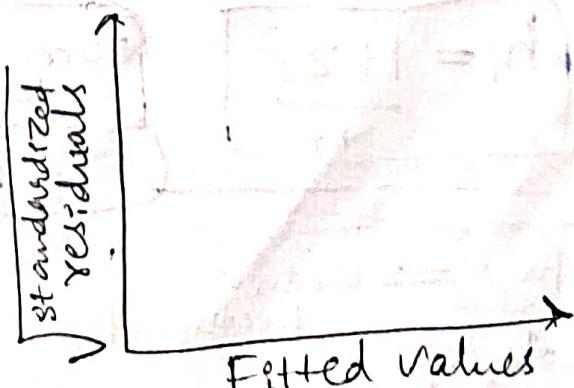
The residues must be close to Normal Q-Q line



→ Scale - Location

(checks for Homoscedasticity)

If H.sce exists then value of errors must be same.

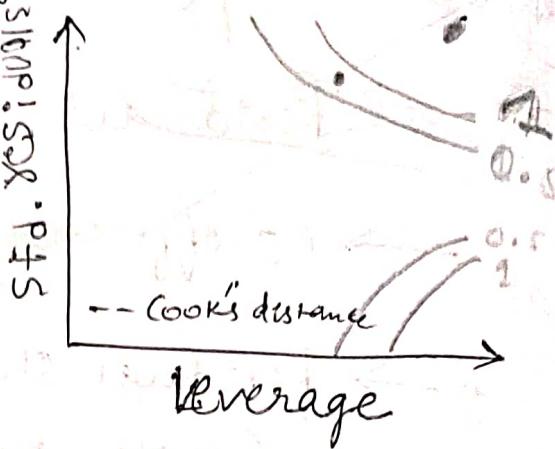


→ Residuals Vs Leverage

(*) Any data point that has an Cook's-d value > 1 is considered as influential outliers.

↳ we must remove that data point and rebuilt the model.

(*) Any data point that has an Cook's-d value b/w 0.5 & 1, those are called suspected outliers.



Quartile: Divide the entire data into 4-sub-range such that each sub-range contains 25% of the data.

Quantile: Divide the entire data into k sub-range (k -quantile) such that each subrange contains $100/k$ percent of the data.

Leverage: How far away is the predictor variable from the mean of predictor variables.

$$h_i = \frac{1 + z_i^2}{n}$$

$$\text{Sum of leverage} = \text{mean of leverages} \times \text{Sample size}$$

Sum of leverage \leftrightarrow Number of parameters

$$\frac{\text{No of parameters}}{\text{Sample size}} = \text{mean of leverages}$$

Example

- $h > 2$; considered as outliers
 $\text{mean}(h)$ ~~generalized about bounded with P.D.~~
 Therefore leverage tells how far the x value far from mean of x .
 \therefore Farther the x from $\text{mean}(x)$, higher the leverage.
 \therefore Therefore \uparrow the leverage values or h .
 Values for each x value is greater than 2 , then we consider that point as influential outliers.

Studentized Residual :-

$$\text{std res.}_i = \frac{e_i}{\sqrt{\text{MSE}(1-h_i)}}, \quad 2 < \text{std res.}_i < -2$$

Std residual tells, how the errors of predictions (y -value) influence.

Combining both Leverage and std residual we get **Cook's D**

Measures overall influence of an observation by seeing the regression co-efficients.

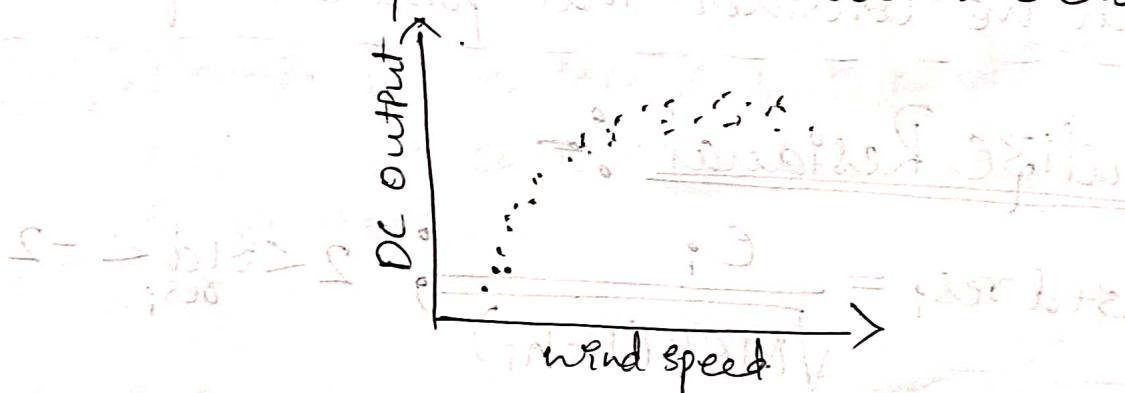
$$D_i = \frac{1}{P} (\text{std res.}_i)^2 \left(\frac{h_i}{1-h_i} \right)$$

Tukey's four Quadrant approach

- *) If the assumed linear relationship is actually non-linear for the model, then we will do non-linear transformation i.e. $x^2, x^3, \log(x), \log(y)$, etc, to get a linear relation.

→ How to know that there is a non-linearity?

→ The Scatter plot will be having a pattern as shown below.



Tukey's ladder

Ladder for x.		
Up ladder	Neutral	Down ladder
\dots, x^4, x^3, x^2, x	$\sqrt{x}, x, \log x$	$-\frac{1}{\sqrt{x}}, -\frac{1}{x}, -\frac{1}{x^2}, -\frac{1}{x^3}, -\frac{1}{x^4}, \dots, \log x$

Ladder for y.		
Up ladder	Neutral	Down ladder
\dots, y^4, y^3, y^2, y	$\sqrt{y}, y, \log y$	$-\frac{1}{\sqrt{y}}, -\frac{1}{y}, -\frac{1}{y^2}, -\frac{1}{y^3}, -\frac{1}{y^4}, \dots, \log y$

- *) First do transformation of x.

Move toward y^2, y^3
(upladder on y) ---

or
(downladder on x^2)
towards $\log x, -\frac{1}{\sqrt{x}}$

more toward y^2, y^3
(upladder on y) ---

or

toward x^2, x^3 ---
(upladder on x)

Move toward
 $\log x, -\frac{1}{\sqrt{x}}$
(downladder on x)
or

Move toward $y, -\frac{1}{\sqrt{y}}$
(downladder on y)

Move toward x^2, x^3
(upladder on x) ---

or

Move toward $y, -\frac{1}{\sqrt{y}}$
(downladder on y)

Multiple linear regression

* When there is more than one X-variable.

↳ Some or all variables could be Categorical

↳ The variables could be interacting among themselves

↳ The variable could depend on each other

$$\rightarrow Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

$$\rightarrow Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + (\beta_3 x_1 x_2) + \dots + \beta_n x_n + \epsilon$$

If only ~~two~~ x_1, x_2 variable we have only one Interaction term

$$\rightarrow Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + (\beta_3 x_1 x_2) + \epsilon$$

Always start with model putting all Variables and only pairwise interaction terms

To address that the variables are related to each other, we must look at the co-relation b/w them.

if some of the variables be related to each other
for eg:- x_1 is related to $x_2, x_3, x_4, x_5, \dots, x_n$

These variables are multicollinear

(*) Multicollinearity: -

$$x_1 = \gamma_1 + \gamma_{12}x_2 + \gamma_{13}x_3 + \dots + \gamma_{1n}x_n + \epsilon \quad R^2_1$$

$$VIF = \frac{1}{1 - R^2_1}$$

$$x_k = \gamma_k + \gamma_{k1}x_1 + \gamma_{k2}x_2 + \dots + \gamma_{k(k-1)}x_{k-1} + \gamma_{kk+1}x_{k+1}$$

$$+ \dots + \gamma_{kn}x_n + \epsilon.$$

$$R^2_k$$

$$VIF = \frac{1}{1 - R^2_k}$$

(*) Variance Inflation Factor:

$$VIF = \frac{1}{1 - R^2_k}$$

or highest correlation coefficient of any pair of all the independent variables

R^2

$$VIF = \frac{1}{1-R^2}$$

if R^2 is small

R^2	VIF
0.5	2
0.75	4
0.8	5
0.9	10
0.95	20
0.99	100

The Standard VIF is greater than 4.

$$\boxed{VIF > 4}$$

If these are the equations, i.e.:-

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n + \epsilon$$

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$SSE = \sum_{i=0}^n (y_i - \hat{y}_i)^2 \Rightarrow \sum_{i=0}^n (y_i - (\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n))^2$$

$$\frac{d(SSE)}{d\beta_k} = 0$$

We will get k equations :-

This is too much of Calculus.
Since we have large number
of variables.

∴ we move to Numerical Techniques

Numerical technique :-

- ④ It involves the assumption of values of $\beta_0, \beta_1, \beta_2, \dots$ to calculate SSE.
- ④ Now, Change all the β -values by sum amount and see that SSE changes or not
- ④ If it decreases, Keep Changing the value until it reach to minimum value and starts to increase from there
- ④ ∴ The SSE with min value is considered as our model.
- ④ The changing of these values is done by using Gradient Decent method.

R-Studio :-

$$ft\text{LM} \leftarrow lm(ft\$A \sim ft\$B + ft\$C + ft\$D + ft\$E)$$

summary(ftLM) :-

→ Residuals :-

Test whether it is skewed or not.

→ Co-efficients :-

	Estimate	sd. Errn	t value	Pvalue
Intercept	31.6084	7.1051	4.449	0.006
B	7.0676	1.0031	7.046	0.067
C	-0.4201	0.2413	-1.741	0.14216
D	-0.2375	0.5086	-0.467	0.66018
E	-0.7936	0.2977	-2.666	0.04658

→ Multiple R-Squared value : 0.9186

Adjusted R-squared : - 0.8535

F-statistic : 14.11 on 4 and 5 DF

P-value : 0.006232

B-Values (Co-eff)
are not significant

Since P value is
0.14 & 0.60, that
means we cannot
reject null hypothesis

This could be happening because C & D might be related to each other.
we must check the relation using VIF.

(*) Vif(ftLM) ← Library(Car)

	B	C	D	E
	1.031045	1.616	1.415	1.209.

Now remove the variable with highest VIF and rebuild the model till you get ~~relevant~~ values.

ftLM2 ← lm(A ~ B + D + E).

(*) Now check the assumptions are true!

→ plot(ftLM2).

(+) By doing these steps, still we are not getting good model!, then.

we must check the variable interaction terms.

Including interaction terms in R-studio is as follows:-

ftLMint ← lm(A ~ B + D + B:D)

Colon

↳ Interaction term

Categorical Variable

- * If there are categorical variables such as gender, geographic region, occupation, marital status, etc. in the regression analysis,

we need to insert ~~n~~ n-1 dummy variable into the regression analysis replacing the category

↳ In R-studio

import:

library(fastDummies)

Now select the column which is categorical
↳ and write the code to remove it as follows

→ $\text{SalGen} \leftarrow \text{dummy_cols}(\text{Data}, \text{Select_Columns} = \text{c("Gender")}, \text{remove_selected_columns} = \text{TRUE})$

→ It will create 'n' number of Gender Column (i.e. Gender_male & Gender_Female).

Since we need only 'n-1' number of Column (i.e either Gen-male / Gen-Female)

∴ we remove any only column by Subset

$\Rightarrow \text{SalGen} \leftarrow \text{Subset}(\text{Data}, \text{Select} = -\text{c(Gen-male)})$.

- * Now we must consider these categorical Variables as continuous variable.

Linear Regression - Summary

- The aim is to build a model of a Numerical Variable (Continuous data) using one or more other related variables.
- Relationship b/w Variables is established by calculating the Co-variance and/or Co-efficient Co-efficient
- The variables that are used to build the model are called predictor / regressors variables
- The variable that is modelled is called the predicted variable / outcome variable.
- The best possible prediction for the outcome variable is only its mean value, the performance of this "no-regressor" model is quantified by SST or Mean Squared Total (MST).
- When there is only one regressor variable, the model is called SLR model.
- The reg var may appear in either linear or non-linear form in SLR.
- Non-linearities can be identified using Tukey's four quadrant approach and built into the model using Tukey's ladder
- An SLR is written in the form $y = mx + c$
↳ 'x' variable may appear in any non-linear form.

- The model building exercise involves finding the value of 'm' and 'c' that minimise the "Sum of Squared Errors" (SSE)
- The performance of the SLR is quantified by using
 - R-squared
 - The hypothesis testing on the slope(m)
- $R^2 = \frac{SSR}{SST}$ (R^2 denotes the proportion of total error explained by model)
- When the number of regressor variables is more than 1, we have "multiple linear regression" (MLR)
- In MLR, the regressor variables might be correlated among themselves.
 - Correlation among regressor variables may be studied using "Variance Inflation Factor" (VIF)
 - The variables with VIF values greater than 4 are considered to be multicollinear with other variables and may be dropped from the model building exercise.
- Each categorical regressor variable is replaced by $n-1$ dummy variables, where n is the number of categories in that variable.
- Interaction among the regressor variables are modelled as product of the variables.

→ R-Squared value in MLR is defined as

$$1 - \frac{MSE}{MST}$$

↳ Mean errors are determined using appropriate DOF.

↳ MSE is obtained by dividing SSE by $n-k-1$ (DOF) where k is the number of regressor variables used in the model. $MSE = \frac{SSE}{n-k-1}$

↳ MST is obtained by dividing SST by $n-1$. $MST = \frac{SST}{n-1}$

→ The decision to reject or not reject the model is done by calculating the F-statistic, which is defined as the ratio of MSR to MSE.

↳ MSR is obtained by dividing SSR by k .

$$MSR = \frac{SSR}{k}$$

→ When the model is not rejected, the hypothesis testing on each regressor variable is done using t-test.

→ When one or more regressor is found to be not statistically significant from zero, the model is rebuilt by dropping that/those variables.

→ In MLR, it is better to start with the most complex model that has all the variables and all "pairs" of interactions.

↳ The nonsignificant variables from this are dropped and a new model is built.

↳ The process is repeated until the remaining regressors are significant.

Classification Models II

Next.

The outcome could be either Continuous or Categorical.

What if Continuous \rightarrow linear regression

Categorical \rightarrow (next topic): logistic regress
(Binary)

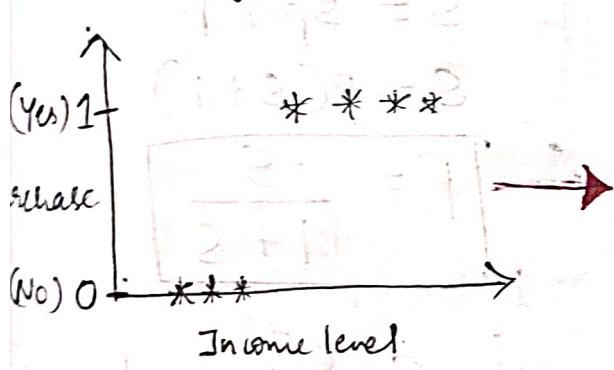
Logistic Regression

↳ Binary classification technique

Binary Classification:

↳ Outcome variable is either 0, or 1

↳ Regressor variable are numerical or categorical

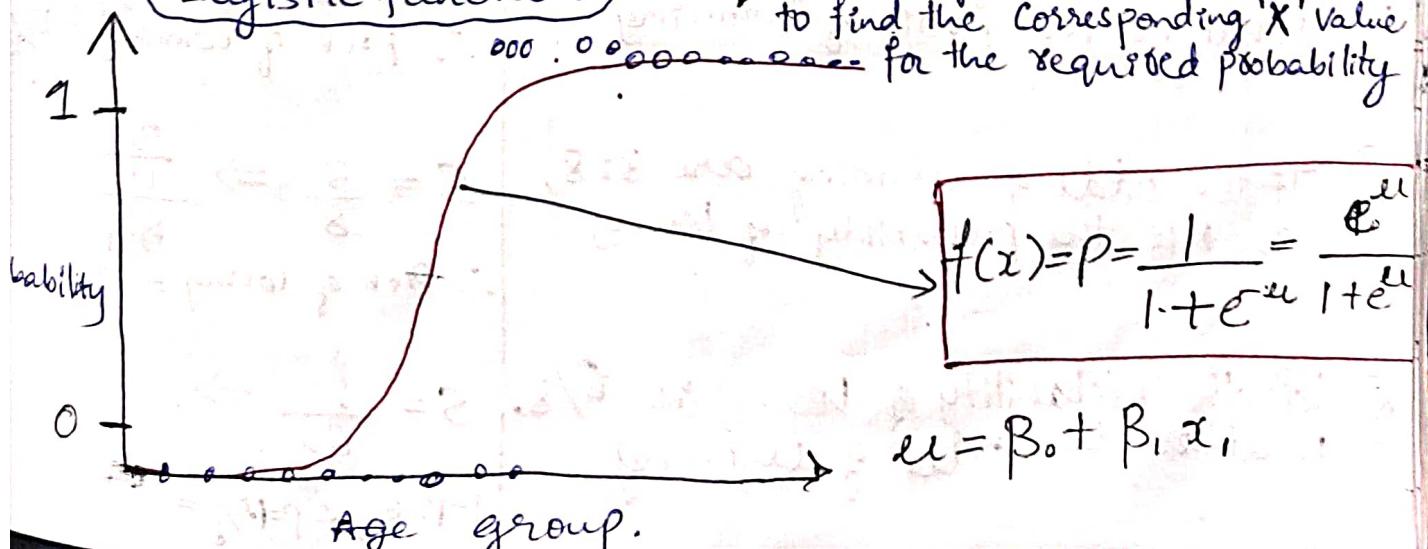


For this kind of data (Binary)
if treat the outputs as continuous
and build a linear regression
model, then that would
not give exact prediction.

① We find the Probability values and get the Sigmoid / psychometric curve as shown below.

Logistic function

Once we get this curve, we have
to find the corresponding 'X' value
for the required probability



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

$$x = B_0 + B_1 x_i$$

Logistic - Sigmoid - Ogive - Psychometric function

Mathematically

$$f(x) = p = \frac{1}{1 + e^{-x}} = \frac{e^{-x}}{1 + e^{-x}} = \frac{e^{-\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}}$$

(1)

Since this is extremely non-linear function.

∴ we define the probability in terms of Odds.

→ we introduce a term called Odds (S)

$$S = \text{Odds} = \frac{P}{1-P}$$

$$\Rightarrow S(1-P) = P$$

$$S - SP = P$$

$$S = Sp + P$$

$$S = P(S+1)$$

$$P = \frac{S}{1+S}$$

* If the prob of winning is $6/12$.
what are the odds of winning

$$S = \frac{6/12}{1-6/12} = \frac{1}{1} = 1:1$$

* If the odds of winning are $13:2$
what is the prob of winning

$$S = \frac{13}{2} \Rightarrow \frac{13}{15}$$

∴ Prob of winning = $13/15$

* If the odds of winning are $3:8$,
what is the probability of losing

$$S = \frac{3}{8} \Rightarrow \frac{3}{11}$$

∴ Prob of losing = $8/11$

* If the probability of losing is $6/8$,
what are the odds of winning

$$S = ? \Rightarrow$$

$$1-P = 6/8 \therefore P = 6/8 =$$

Since from ①

$$S = \text{odds} = \frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}}$$

$$= \frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}}$$

$$\therefore S = e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n}$$

Since the above equation is non linear we convert into linear

apply log on b.e

$$\log(S) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n \rightarrow \text{Solve for Betas}$$

\rightarrow Logit function.

Probability (P)	Odds (S)	$\log(S)$
0	0	undefined
0.10	1/9	-2.20
0.25	1/3	-1.39
0.5	1	0
0.75	3	1.39
0.9	9	2.20
0.95	19	2.94
1	∞	inf

Q) Now we must set the boundary.

$$\therefore \log(s) = 0 = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

→ Boundary

In 1-dimensional;

$$\log(s) = \beta_0 + \beta_1 x_1 = 0$$

$$\Rightarrow x_1 = \frac{\log(s) - \beta_0}{\beta_1}$$

This will be the boundary

In 2-dimensional;

$$\log(s) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$$

$$\Rightarrow x_2 = \log\left(-\frac{\beta_0}{\beta_2} - \frac{\beta_1}{\beta_2} x_1\right)$$

This is the boundary

This is of the form $y = mx + c$

$$m = -\frac{\beta_1}{\beta_2}$$

$$c = -\frac{\beta_0}{\beta_2}$$

* In linear regression for finding the values of β_0 & β_1 we used minimized SSE.

i.e., for least value of SSE we take/find the corresponding β_0 & β_1 value.

→ But SSE does not work for logistic regression problems. Since we are solving for probability which is highly non-linear and has lot of minima and maxima.

∴ we introduce a function which deals with Probability → log likelihood function

Log Likelihood

$$\log(\text{likelihood}) = \sum_{i=1}^n [y_i \ln \hat{y}_i + (1-y_i) \ln (1-\hat{y}_i)]$$

Cost Function

here:

$y_i \rightarrow p_i \rightarrow \text{Probability of the } i\text{th point to be in the Category 1}$

$\hat{y}_i \rightarrow \hat{p}_i \rightarrow \text{Predicted probability of the } i\text{th point to be in the Category 1.}$

- ✳ The accuracy will be high when the $y_i \approx \hat{y}_i$ (The prediction matches with actuality)

we get \hat{y}_i from:

$$\log(s) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

i) take anti log of this, we get odds

ii) get the probability from odds, which is predicted probability.

- ✳ The Cost function is used to find the β_0, β_1, \dots values by using gradient descent method.

R-Studio

`floesLgR <- glm(floes$A ~ floes$B,
family = 'binomial')`

`summary(floesLgR)`

↳ Deviance Residuals:

	Min.	Q1	Median	Q3	Max
	-1.95015	-0.321	-0.05335	0.265	1.72

(Coefficients:

Model	Estimate	Std. Error	Z Value	P-Value
Intercept	-20.40782	4.52332	-4.512	6.43e-06
floes\$B	0.4259	0.09482	4.492	7.09e-06

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 123.156 on 91 DOF

Residual deviance: 49.937 on 90 DOF.

Number of Fisher Scoring iterations: 7

④ Since the p-values are lesser than 0.05
 \therefore cannot reject the model (Since we reject the null hypothesis if $\text{intcept } \& B_1 = 0$).

* we do not get F-value in logistic Reg

* Since we know that F-statistic is.

$$F\text{-stat} = \frac{M.E.E}{M.U.E.E}$$

∴ we donot have any explained/unexplained errors in log-reg therefore there is no F-stat.

* Therefore we look at Null deviance and Residual deviance.

* Since we got the β_0 & β_1 values from the summary of model.

i.e

$$\text{log}(S) = -20.40782 + 0.4259 * \text{Age.}$$

and loglikelihood can be found by.

dividing the Residual deviance by -2

$$\text{i.e } \text{log}(likelihoood) = \frac{49.937}{-2} = -24.97.$$

$$\boxed{\text{Residual deviance}} = -2 \times \text{log likelihood}_{\text{without model}}$$

$$\boxed{\text{Null deviance}} = -2 \times \text{log likelihood}_{\text{without model}}$$

* ∵ we treat Residual deviance as explained errors and Null dev as total errors.

$$\therefore \text{Pseudo } R^2 = 1 - \frac{\text{Residual Deviance}}{\text{Null Deviance}} \approx 1 - \frac{\text{SSR}}{\text{SST}}$$

Most often we do not use Pseudo R^2 since it tells us that how close is P_i to \hat{P}_i but we are not worried about that (But interested in categorising the data into yes or no).

We can find the 'x' value (Age) for the corresponding probability \rightarrow Boundary

$$\log(S) = -20.41 + 0.43 * \text{Age}$$

if $P=0.5$ boundary of Age = ?

\Leftrightarrow for $P=0.5$ $S=1$.

$$\log(S) = 0$$

$$\therefore -20.41 + 0.43 * \text{Age} = 0$$

$$\Rightarrow \boxed{\text{Age} = 47.9}$$

i.e above 47.9 age will say yes is the model prediction for $P=0.5$

Also - From ①

For age of 'x', which category does he falls.

$$\log(S) = -20.41 + 0.43 * 50$$

$$\log(S) = 0.89$$

$$S = e^{0.89} = 2.43$$

$$S = 2.43$$

$$P = \frac{S}{1+S} = 0.71$$

For this model what kind of errors are possible is given below by Confusion Matrix or Truth Table.

		Prediction	
		Class 1(S)	Class 0(F)
		Predicted to be true which is true.	Predicted to be -ve which is false.
Actual	Class 1 (Success) +	Predicted to be positive which is false.	Predicted to be -ve which is true.
	Class 0. (Failure) -		

		Prediction	
		Positive	Negative
		True Positive (n ₁)	False negative (n ₂)
Actual	Positive	False Positive (n ₃)	True Negative (n ₄)
	Negative		

$$n_1 + n_2 + n_3 + n_4 = N \text{ (Total persons)}$$

$$\text{Accuracy} = \frac{n_1 + n_4}{N}$$

$$\text{Sensitivity} = \frac{\text{Correctly identified to be true}}{\text{Actually Positive}} = \frac{n_1}{n_1 + n_2}$$

$$\text{Specificity} = \frac{\text{Correctly identified to be -ve}}{\text{Actually Negative}} = \frac{n_4}{n_3 + n_4}$$

$$\text{Precision} = \frac{\text{Correctly identified to be true}}{\text{Predicted Positive}} = \frac{n_1}{n_1 + n_3}$$

According to the which type of question we are solving, we must consider ~~the~~ Accuracy, Precision/specificity/Sensitivity accordingly.

In case of unbalanced data, accuracy is not good quantity to consider, ∴ we consider a new metric called Kappa Metric.

Kappa Metric :-

→ This says that "How well the model is better than Random Classification."

$$\text{Kappa} = \frac{\text{total Accuracy} - \text{Random Accuracy}}{1 - \text{Random Accuracy}}$$

$$\text{total Accuracy} = \frac{\text{Correct Predictions}}{\text{Total}}$$

$$\text{Random Accuracy} = \frac{\text{Actual False} * \text{Pred. Fal}}{\text{Total}} + \frac{\text{Actual True} * \text{Pred. True}}{\text{Total}}$$

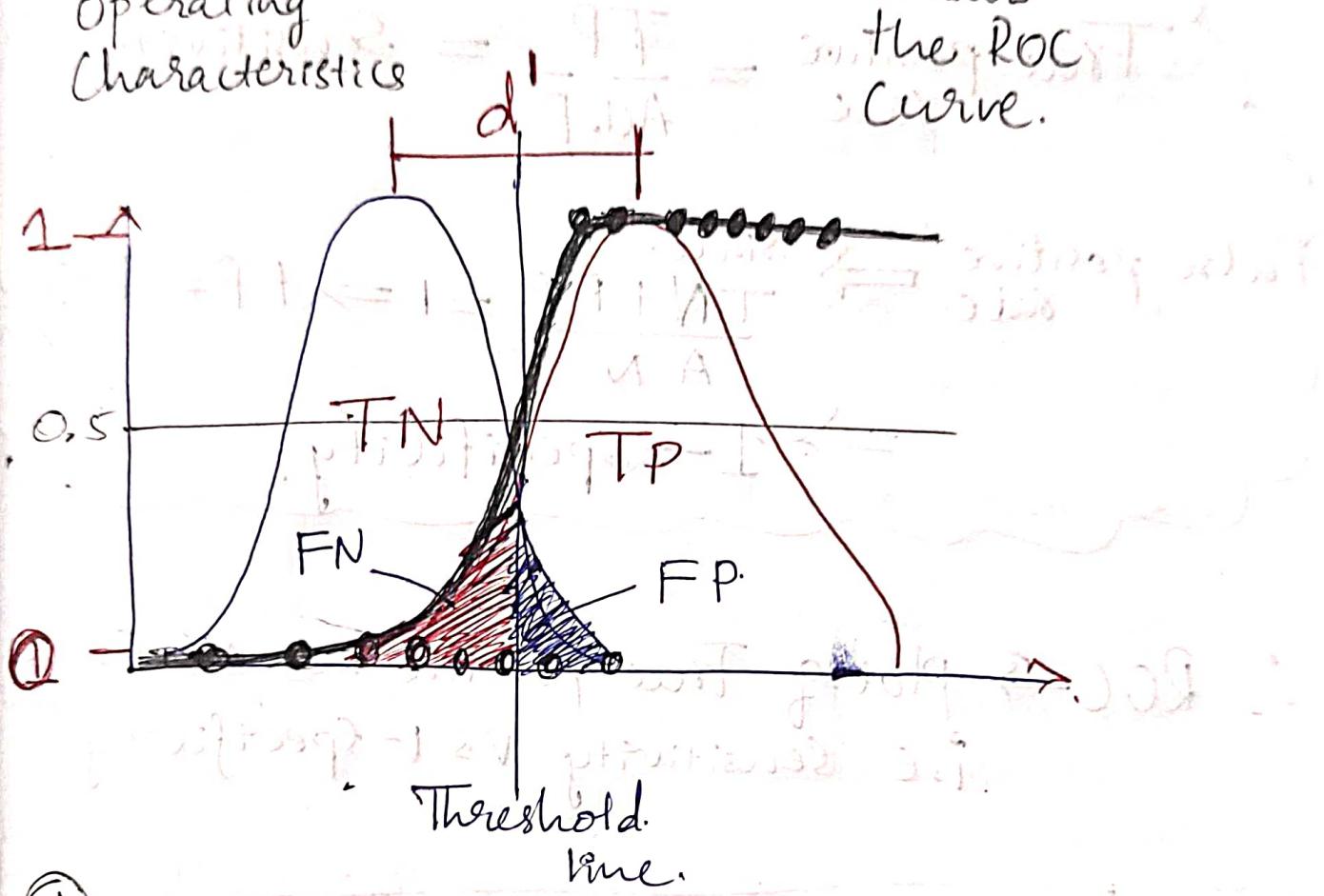
Kappa Value.

< 0	NO Agreement
0 - 0.2	Slight
0.21 to 0.4	Fair
0.4 to 0.6	Moderate
0.6 to 0.8	Substantial
0.8 to 1	Almost Perfect.

ROC Curves and AUC

Receiver operating characteristics

Area under the ROC curve.



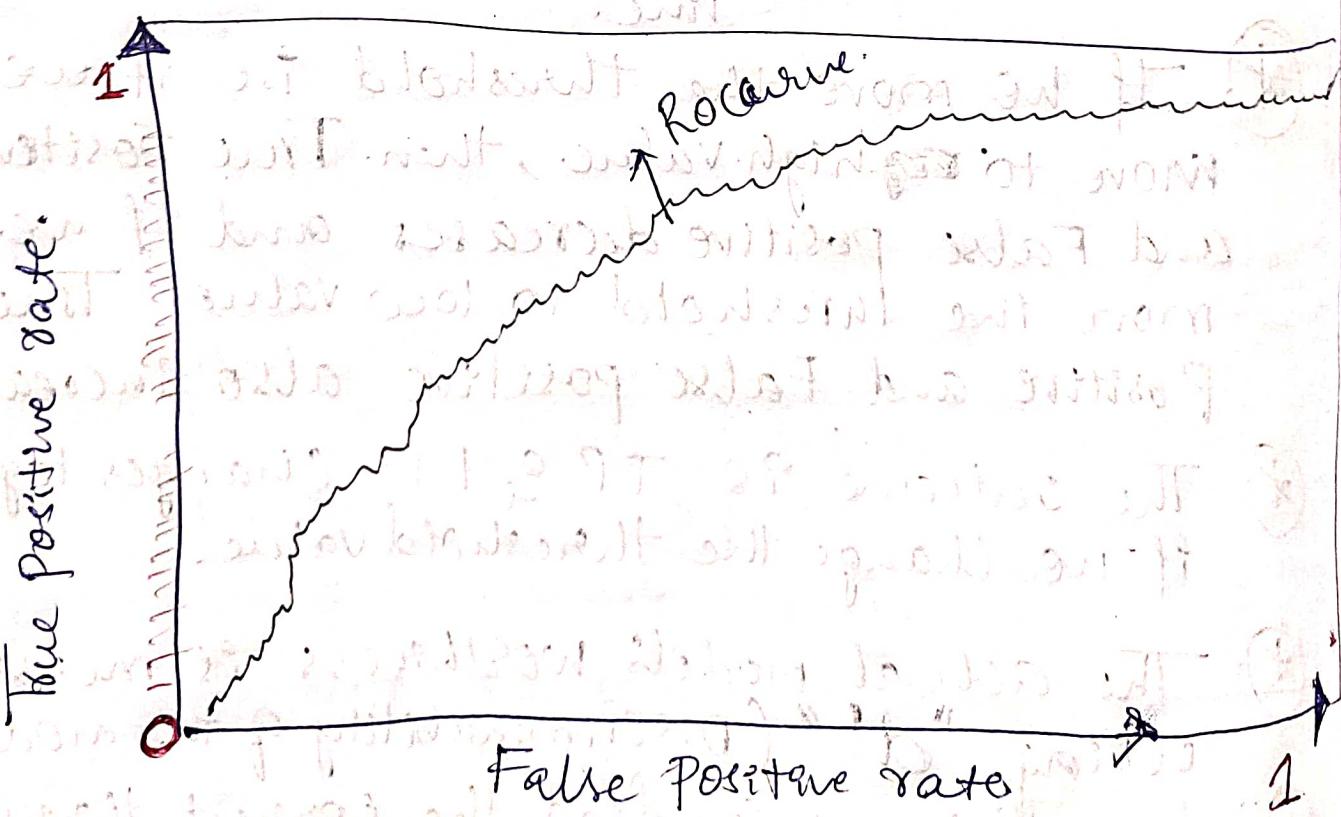
- (*) If we move the threshold i.e if we move to high value, then True positive and False positive decreases and if we move the threshold to low value. True Positive and False positive also increases.
- (*) The outcome is TP & FP changes together if we change the threshold value.
- (*) The actual model's worthiness is measured using "d" (Discriminability of the model)
- (*) From this we can find the correct threshold value.

④ The graph of Sensitivity vs. 1-Specificity. E.g. is called as ROC curve.

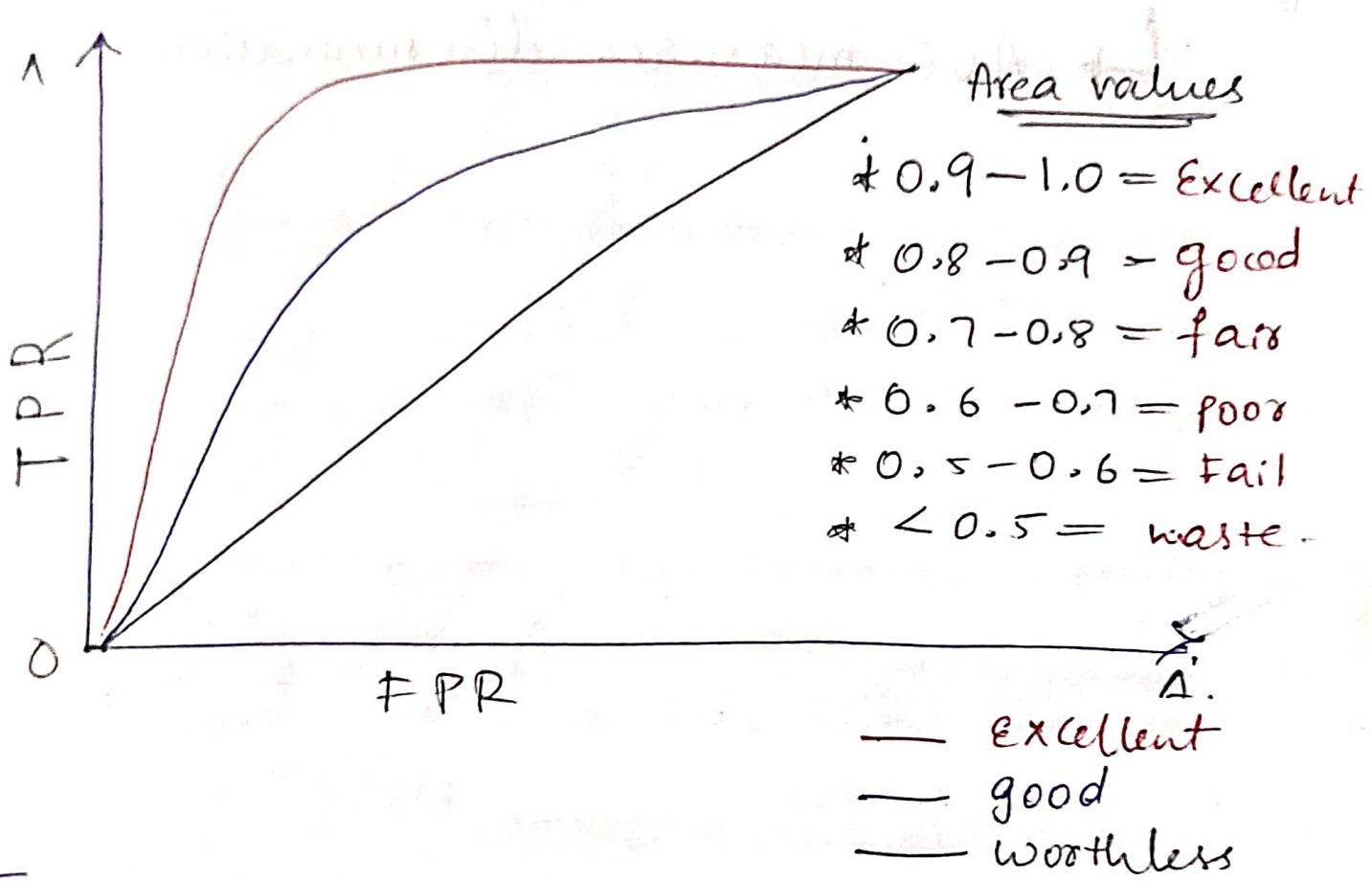
$$\text{True positive rate} = \frac{\text{TP}}{\text{Act. P}} = \text{Sensitivity}$$

False positive rate. \Rightarrow Since $\frac{\text{TP} + \text{FP}}{\text{AN}} = 1 \Rightarrow \text{FP} = 1 - \frac{\text{TN}}{\text{AN}}$
 $\Rightarrow 1 - \text{Specificity}$

\therefore ROC \rightarrow Plot of True positive rate vs False Positive rate, i.e. Sensitivity vs 1-Specificity.

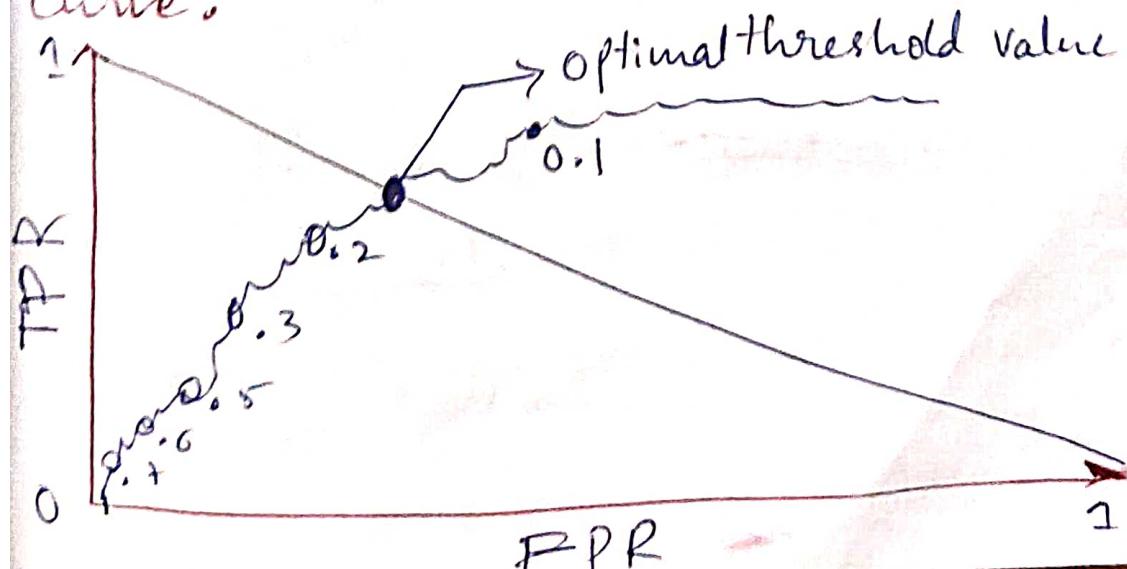


- (*) The area in the plot is one. Since 1×1
- (*) Area under ROC curve tells how good is the classification (i.e. Discriminability)
- (*) Larger the area better is Discriminability



Final Question is how do we pick the best threshold for the model?

By drawing a diagonal from left to right and mark the point where diagonal cuts the ROC Curve.



④ ROC will help to get the best threshold for model.

⑤ AUC will help to determine LRM is good or not; higher the AUC ; better the model

↳ AUC measures discrimination

Naïve Bayes

- Naïve Bayes is used for multi-class classification.
- It is the application of Bayes theorem.

Naïve Bayes Algorithm

- Naïve Bayes:- Computes $P(y_i | x)$ by using Bayes theorem and instead computes the inverse conditional probability $P(x | y_i)$.
- A simple classifier that performs surprisingly well on a large class of problems.
- This type of methods are called Generative learning Models.

find $\rightarrow P(A|B) = ?$

Given $\rightarrow P(B|A)$

O.K.t. $P(A|B) = \frac{P(A)}{P(B)} \times \underbrace{P(B|A)}_{\text{Cond prob}}$

\downarrow marginal probability

This is Bayes theorem.

But In Naïve Bayes classification method we proceed with the above step along with the "Naïve Assumption". For eg \rightarrow to find $P(A|B, C)$

We assume that B & C are independent of each other and we rewrite as $P(A|C)$

Example :-

What is the probability of winning the match when the Indian cricket team wins the toss & chooses to field.

Event A : Toss : $A = \{W, L\}$

Event B : Choice : $B = \{\text{Field, Bat}\}$

Event C : Match outcome : $C = \{W, L, D\}$

To find

$$P(C=W | A=W, B=\text{Field}) = ?$$

We make an assumption that Event B and Event A are independent. This assumption is called Naive Assumption.



$$P(C=W | A=W, B=\text{Field}) = P(C=W | A=W) \times P(C=W | B=\text{Field})$$

$$\therefore P(C=W | A=W, B=\text{Field}) = \frac{P(A=W, B=\text{Field} | C=W)}{P(A=W) \times P(B=F)}$$

$$\Rightarrow \frac{P(A=W | C=W) \times P(B=\text{Field} | C=W) \times P(C=W)}{P(A=W) \times P(B=F)}$$

⇒ The above eqⁿ is final eqⁿ and solve for it.

R-Studio:-

```
library(e1071)
```

```
nb_model <- naiveBayes(class ~ ., data =  
trainHouseVotes84)
```

It gives all the prob values which must be used in the equation.

- The Naive assumption we makes is never true.
- Still Naive Bayes does surprisingly well in a lot of situations
- It works best if predictor variables are Categorical
- It is a parameter-free model

Another important idea is "When we have Continuous data we Create some bins with some conditions and make it Categorical, then we solve using Naive Bayes algorithm".

Gradient Descent

$$y = f(x_1, x_2)$$

Let, $y = 3x_1^2 + 4x_2 - 3$



Gradient is given as

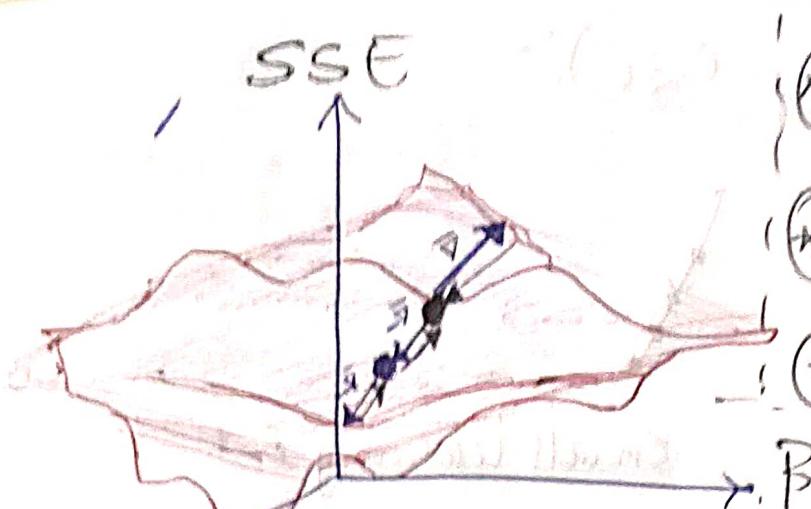
$$\vec{\nabla}y = \frac{\partial y}{\partial x_1} \hat{x}_1 + \frac{\partial y}{\partial x_2} \hat{x}_2$$

→ Gradient is always point in the direction of steepest increase in y .
 * Since we are interested in the steepest decrease in y ∴ Steepest decrease is $-\vec{\nabla}y$

* In regression problem our aim was to minimise SSE w.r.t. β values.

∴ In the above figure instead of x_1, x_2 we have β_1, β_2 and instead of y we have SSE





- (*) Start with a random β -value.
- (*) find the gradient (∇) and $-\bar{\nabla}$
- (*) In opp direction take a small step $\rightarrow \beta_2$, i.e. change β_1 or β_2 by small amount
This change is called as Step Size.

- (*) Now we reach to a new point and find the $-\bar{\nabla}$ once again.
- (*) Every time we reach the new point calculate the sum of squared errors and observe that the SSE is lower than the previous then. Continue doing this process; by the time SSE starts to increase, stop going beyond and consider the lowest SSE and corresponding β_1 & β_2 values.

$$\theta_0^{(t+1)} = \theta_0^{(t)} - (\alpha) \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}$$

(New point) (Original point) Learning Rate Slope of steepest descent (direction of fall)
 (Step size)

- (*) Setting the step size is an art and it depends on the problem given.

Learning Rate (Step size)



Big learning rate

Small learning rate

- (1) Choosing a big learning rate cannot help to get exact minimum point in less time. Value also makes lot of steps.
- (2) Choosing very small learning rate also makes lot of steps.

∴ Choosing a step size is an art and comes by practice.

Optimization Algorithms

① Derivative Based Algorithms (Gradient Based)

- ↳ Steepest descent method.
- ↳ Newton-Raphson method.
- ↳ Conjugate gradient etc.

② Derivative-free methods

- ↳ Random Search method
- ↳ Genetic algorithm
- ↳ Simulated annealing etc.