

LINEAR REGRESSION / LEAST SQUARES

(b) the 'mean' cuts through the average value of 'y'
 ↳ gives us a good starting point to compare.

how well this line 'b' }
 fits the data ↓

$$= (b - y_1)^2 + (b - y_2)^2 + (b - y_3)^2 + \dots$$

is call sum of Squared Residuals = 25

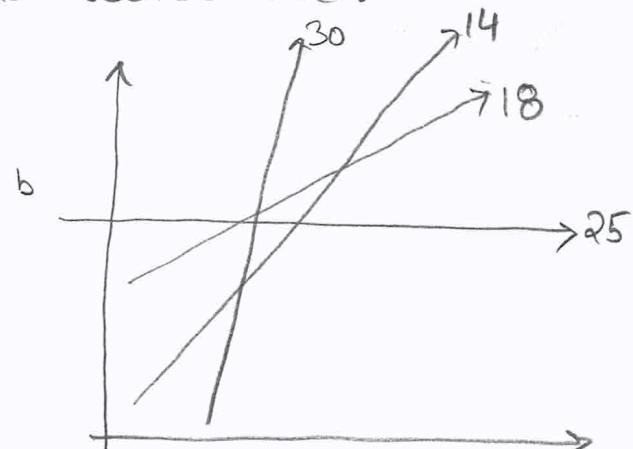
Residuals = difference b/w real data & the fitted line 'b'.

Rotating line 'b' gives better sum of sq. residual.

upto a certain slope & gets worse later.

$$y = b + ax$$

↓
slope
intercept

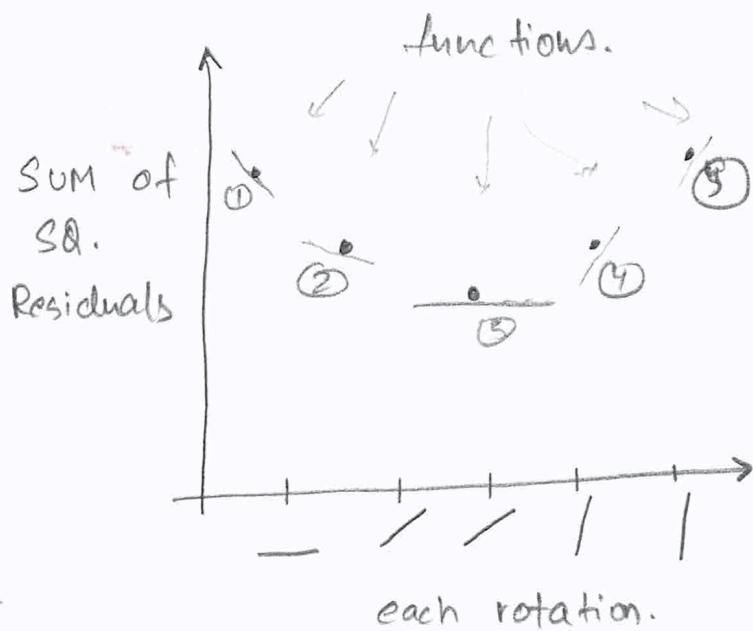


To Find the values of 'a' & 'b' so that we

MINIMIZE the sum. of Sq. Residuals.

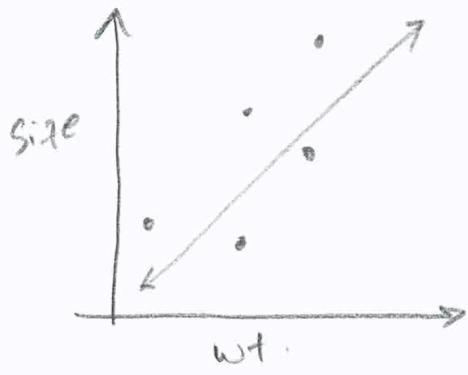
$$= \underbrace{((ax_1 + b) - y_1)^2}_{\text{value of line at } x_1} + \underbrace{((ax_2 + b) - y_2)^2}_{\text{observed value of } x_2} + \dots$$

To find the optimal rotation of the line \Rightarrow take derivative of the functions 'o' in graph



Derivatives give us the slope of the function at every point.

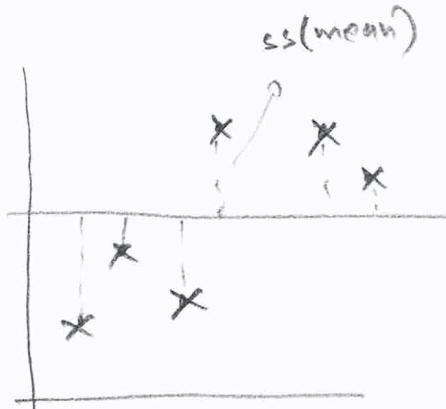
Function ① slope is steepest. & Function ③ has slope = 0 (least sq.)



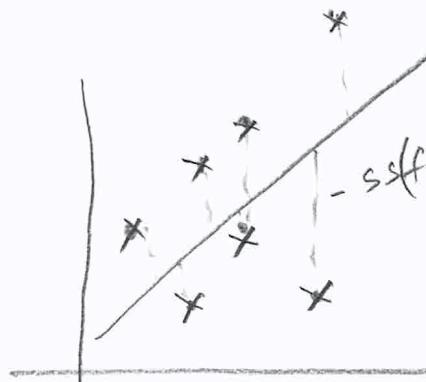
with this line

1) calc. R^2 to find relation b/w size & wt. large value imply large effect.

- 2) calc. p-value if R^2 is statistically significant
- 3) find size given wt.



$$R^2 = \frac{SS(\text{mean}) - SS(\text{fit})}{SS(\text{mean})}$$



LOGISTIC REGRESSION: MAXIMUM LIKELIHOOD, R^2 , p-value

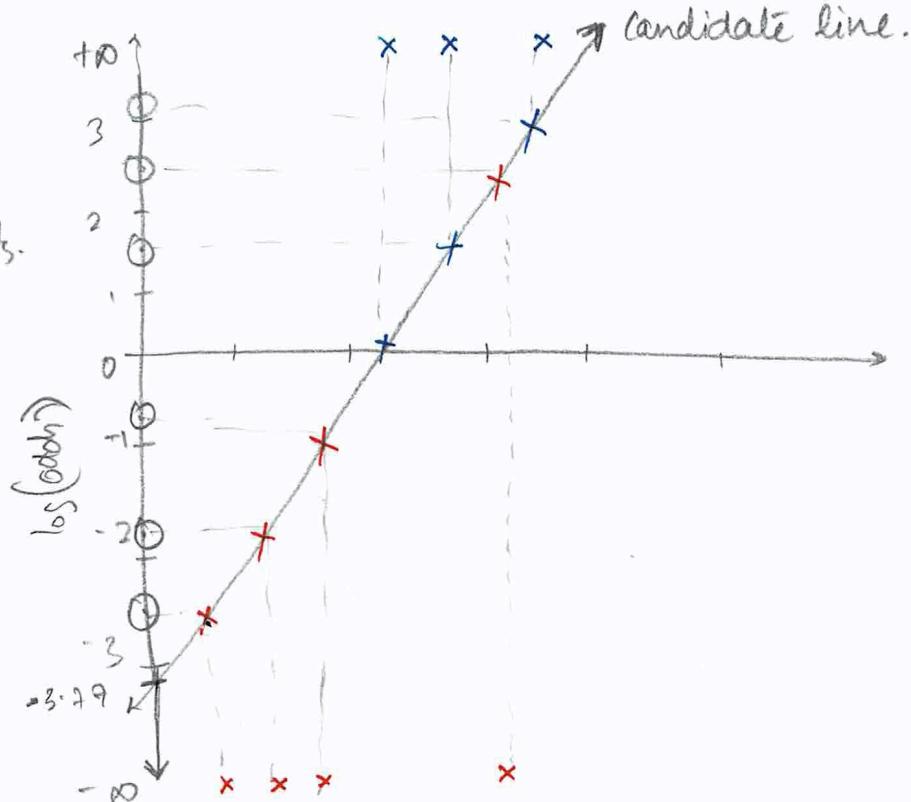
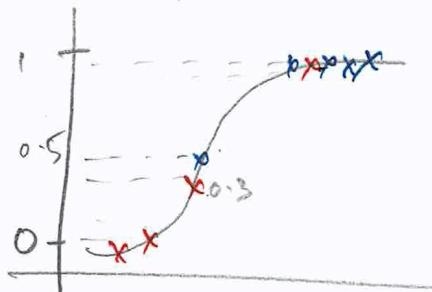
Project the original data on candidate line. These points are log (odds) (O) pts.

Transform

log (odds) to



$$P = \frac{e^{\text{log}(\text{odds})}}{1 + e^{\text{log}(\text{odds})}}$$



x - likelihood that these are obese $\Rightarrow 0.45, 0.9, 0.91, 0.92, \dots$

x - likelihood that these are Not obese:
 $(1 - \text{Pr}(\text{being obese}))$

$$= (1 - 0.9), (1 - 0.3), (1 - 0.01), \dots$$

likelihood of data
 points on squiggle line) $= (0.49 \times 0.9 \times 0.91 \times 0.92) \times$
 $\quad [(1 - 0.9) \times (1 - 0.3) \times (1 - 0.01)]$

We take log & find MAX value for this

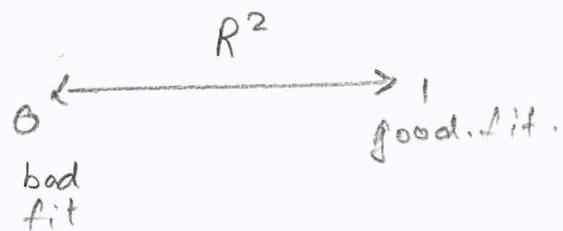
$$\begin{aligned} \text{Max } \log (\text{likelihood of data}) \\ \text{on squiggle line}) &= \log(0.49) + \log(0.9) + \log(0.91) + \log(0.92) + \\ &\quad \log(1 - 0.9) + \log(1 - 0.3) + \dots \\ &= -3.77 \end{aligned}$$

$$LL(\text{fit}) = \log (\text{likelihood of data on squiggle line}) = -3.77$$

probability of obesity = $\frac{\text{no. of obese}}{\text{total}} = \frac{\text{xxxxx}}{\text{xxxxx} \times \text{xxx}} = \frac{5}{9} = 0.55$
 $= LL(\text{probability})$.

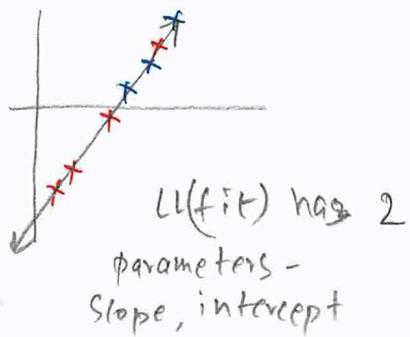
$$R^2 = \frac{LL(\text{prob.}) - LL(\text{fit})}{LL(\text{prob.})} \leftarrow R^2 \text{ from log likelihoods.}$$

$$LL(\text{prob.}) = LL(\text{overall prob.})$$



A Chi-Square value $= 2 [LL(\text{fit}) - LL(\text{prob.})]$

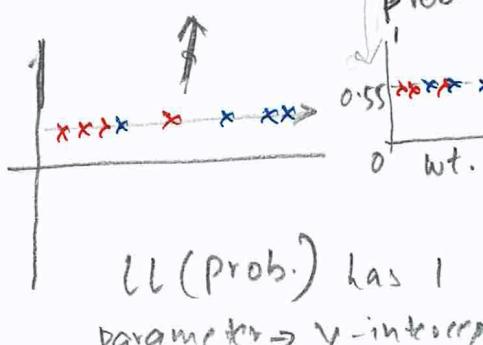
Deg. of freedom:



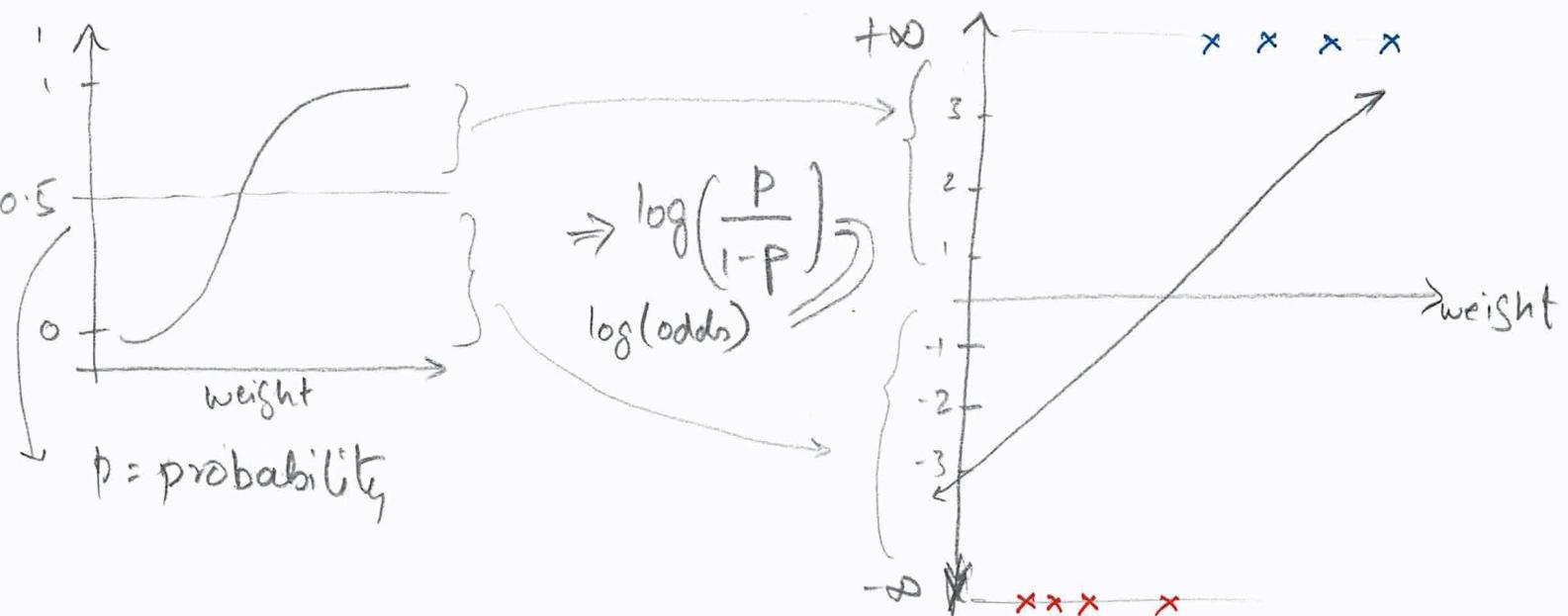
$$\therefore df = 2 - 1$$

$$\phi\text{-value} = (\text{Chi Sq. val. } \& \text{ df})$$

$$\log \left(\frac{\text{obese}}{\text{not obese}} \right) = \text{translate back to prob.}$$



LOGISTIC REGRESSION - COEFFICIENTS.



The coefficients are represented in terms of log (odds) graph

$$y = -3.48 + 1.83 \times (\text{weight})$$

↓ ↓
 intercept slope

Python Result:

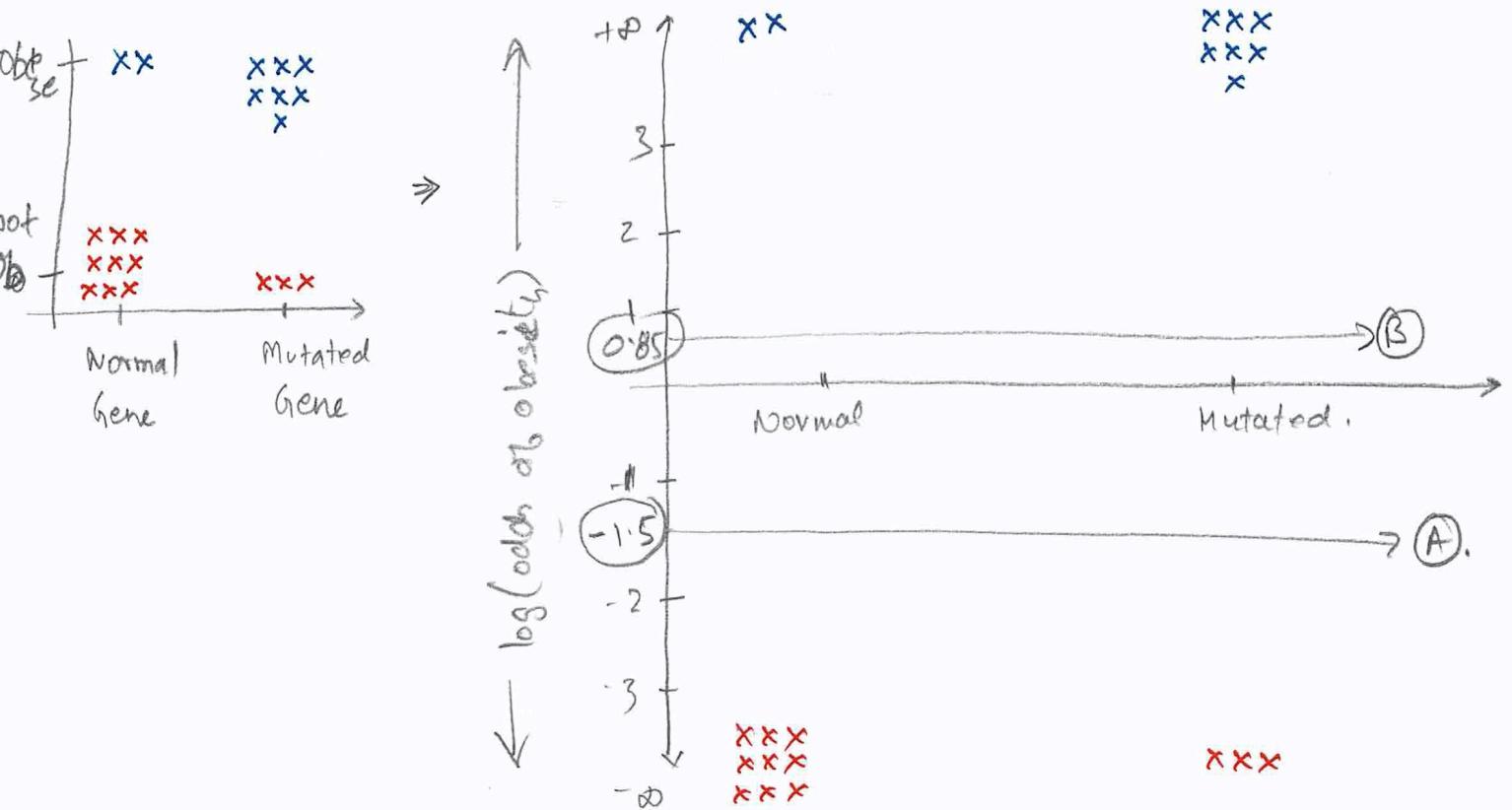
Coefficients:

	Estimate.	Std.Error	Z-value	p-value (> 1z)
(Intercept)	-3.48	2.36	-1.47	0.14
weight	1.83	1.08	1.67	0.09

$Z\text{-value} = \frac{\text{estimate}}{\text{std. error}}$ \Rightarrow (no. of Std.Dev. the estimated intercept is away from ϕ) \Rightarrow Wald's test.

Since its less than 2 Std.Dev. away from ϕ , its not statistically significant. & this is confirmed by p-value. (large)

For every unit of wtg. gained, log(odd of y/n) increases by 1.83.



$$\log(\text{odds})_{\text{Normal}} = \log\left(\frac{2}{9}\right) = -1.5 \rightarrow \textcircled{A}$$

$$\log(\text{odds})_{\text{mutated}} = \log\left(\frac{7}{3}\right) = 0.85 \rightarrow \textcircled{B}$$

} t-test

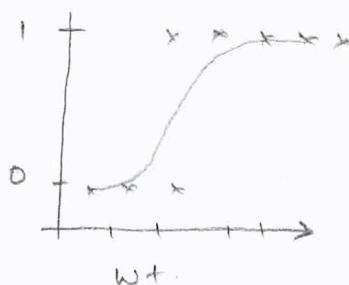
$$\text{Size} = \textcircled{A} \times B_1 + (\textcircled{B} - \textcircled{A}) B_2 = -1.5 B_1 + 2.35 B_2$$

Coeff:	Estimate	Std. Error	Z-val.	Pr(z)
(intercept)	-1.5	0.78	-1.92	0.054
Mutated	2.35	1.04	2.25	0.024

$\log(\text{odds ratio})$ in log scale that tells how much having the 'mutated' ↑/↓ the odds of being obese. ↓ greater than 2 i.e. statistically significant.

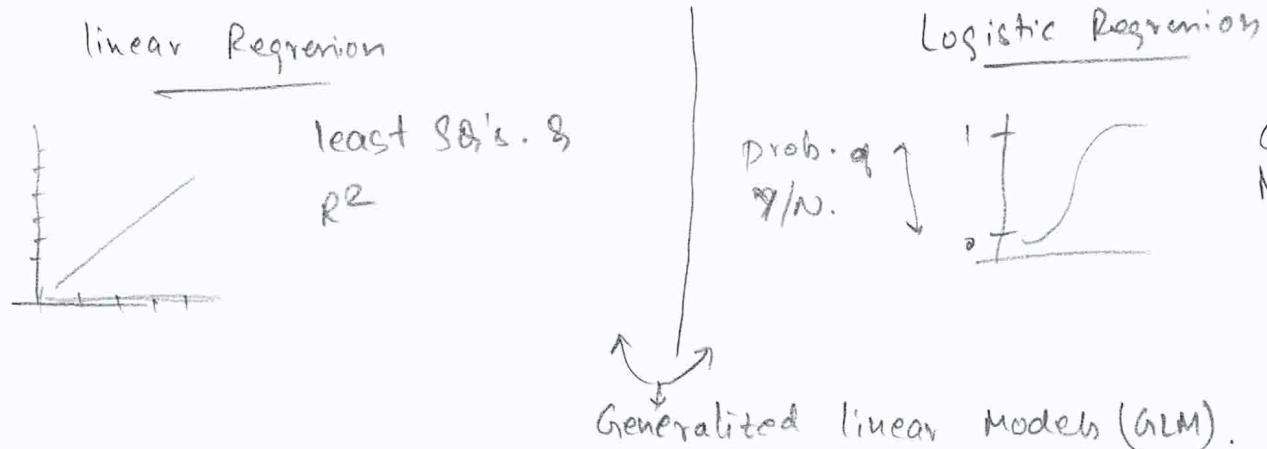
LOGISTIC REGRESSION.

Logistic Reg. predicts if something is T/F, Y/N.



This curve tells the probability of a 0/1 given wt. Although this tells us the probility, it's usually used for classification.

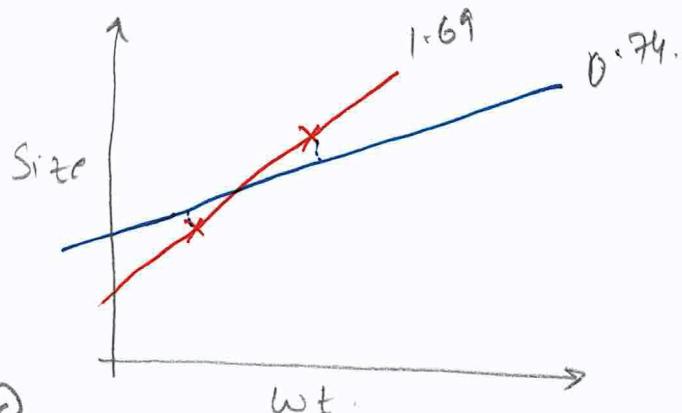
Eg. when prob > 50% \Rightarrow '1' else '0'.



LASSO & RIDGE REGRESSION; ELASTIC RIDGE

By introducing small amount of Bias, we can significantly drop Variance.

least $SQ = y\text{-intercept} + \text{Slope} \times wt.$



$$\begin{aligned} \text{Ridge} &= \frac{\text{Sum of } SQ \text{ residuals/likelihood}}{+ \lambda \times \text{slope}^2} \\ &\quad (\text{linear}) \quad (\text{logistic}) \end{aligned}$$

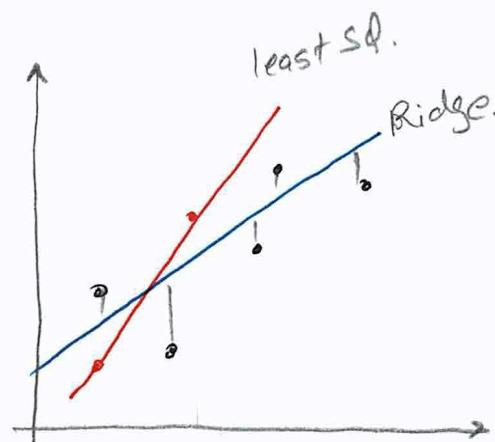
Ridge regression contains all parameters except, y-intercept.

How to determine the best ' λ ' value? ?

$$\begin{aligned} \text{Size} &= y\text{-intercept} + \text{Slope} \times wt + \text{diet difference} \times \text{flight} + \text{speed} \times \text{AirSpeed} \\ &+ \lambda (\text{Slope}^2 + \text{diet diff}^2 + \text{speed}^2 + \dots) \end{aligned}$$

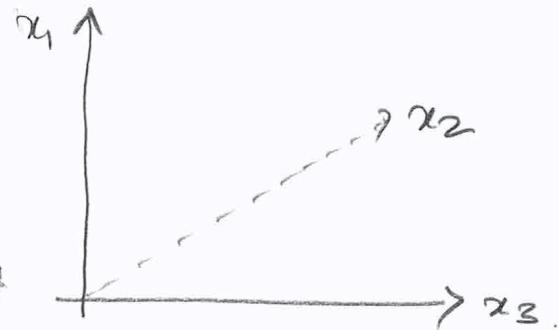
Advantage:

When the sample size are relatively small, the Ridge Regression can improve prediction (made from new data) (reduce variance) by making predictions less sensitive to train data by adding Penalty to the thing that must be minimized. \rightarrow



Even when there isn't enough data to find least sq. parameter estimate, the Ridge Regression can still find Solution with Cross Validation & penalty.

For least sq.- since there are 3 params, ideally we should have 3 data pts.



But this is not a requirement

for Ridge.

$$\text{LASSO} = \text{Sum of sq. residuals} + \lambda \times |\text{slope}|$$

The big difference b/w lasso & Ridge \Rightarrow Ridge can only asymptotically shrink the slope to zero but in lasso, slope=0.

$$\text{Eq: Size} = C + m_1 x_1 + m_2 x_2 + m_3 x_3 + m_4 x_4 + \dots$$

when x_3 & x_4 are not significant, Ridge cannot make m_3 & $m_4 = 0$ but lasso can exclude useless variables from Eq.

Lasso works best when there are lot of useless variables.

Ridge works best when most of the variables are useful

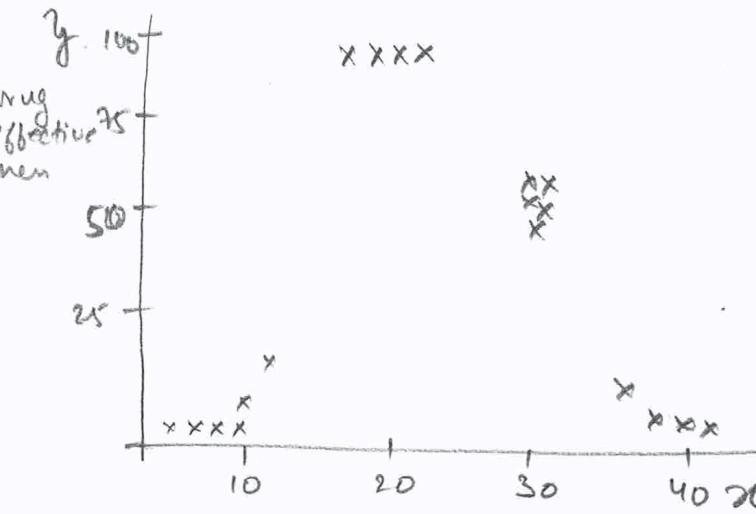
	lasso	ridge
--	-------	-------

$$\text{Elastic Net} = \text{Sum of Sq. residuals} + \lambda_1 \times |\text{slope}| + \lambda_2 \times \text{slope}^2$$

Elastic Net-groups & shrinks the parameters associated with correlated variables & leaves them eq. or removes them all at one time. when the no. of parameters are in millions.

REGRESSION TREE

DTs difficult to apply linear Regression line on this data set. ∴ regression tree.



We determine how to divide the x_i by trying different thresholds of x_i & calculating the sum of sq. residuals. (SSR) & pick the one with MIN. SSR . This becomes the Root Node

③ When there multiple x_i , we pick x_i with least (SSR) as Root Node.

PRUNE: prone to over fit the regression tree.

① Calculate the tree score

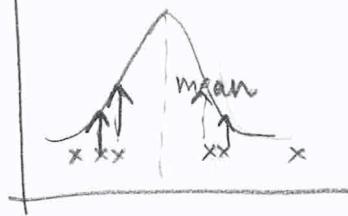
no. of terminal nodes (leaf)

$$\text{Tree Score} = SSR + \alpha T$$

↓ tree complexity penalty
 lowest SSR
 using cross validation.

SATURATED MODELS & DEVIANCE

LOGISTIC REGRESSION

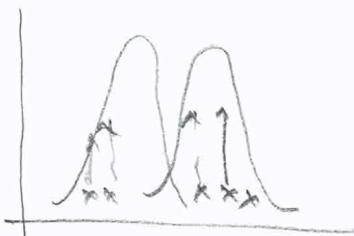


the normal curve is a model of the data.

We use this instead of the original data to
calc. the ^{estimated} probabilities & statistical tests.

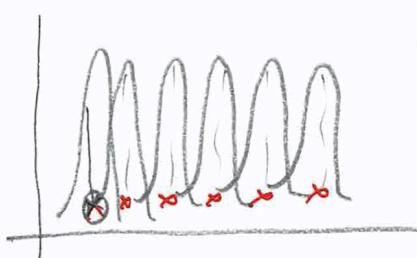
likelihood of data Null Model: $0.2 \times 0.5 \times 0.8 \times 1.5 \times \dots = 0.03$.

$$LL(\text{Null Model}) = \log(0.03) = -3.49$$



No of parameters = 2

$$LL(\text{Proposed Model}) = 3.75$$



Saturated Model = makes out no. of parameters we can use. i.e., one parameter per data pt.

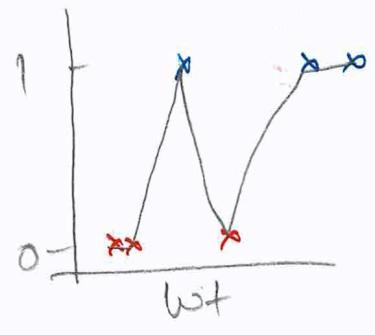
$$LL(\text{Saturated Model}) = 1291.5$$

$$\text{Residual Deviance} = 2 [LL(\text{Saturated model}) - LL(\text{proposed})]$$

↳ 2 times makes these LL difference to have Chi-sq. distribution. with df = no of parameters.

$$\text{Null deviance} = 2 [LL(\text{saturated}) - LL(\text{null})]$$

$$R^2 = \text{Null deviance} - \text{Residual Deviance}$$



In logistic Regression, the Saturated Model perfectly fits the Sigmoid.

$$\therefore R^2 = \frac{LL(\text{null}) - LL(\text{proposed})}{LL(\text{null})}$$

↳

$$\text{Residual deviance} = 2 \times -LL(\text{Proposed})$$

$$\text{Null deviance} = 2 \times -LL(\text{null}).$$

K-NEAREST NEIGHBORS.

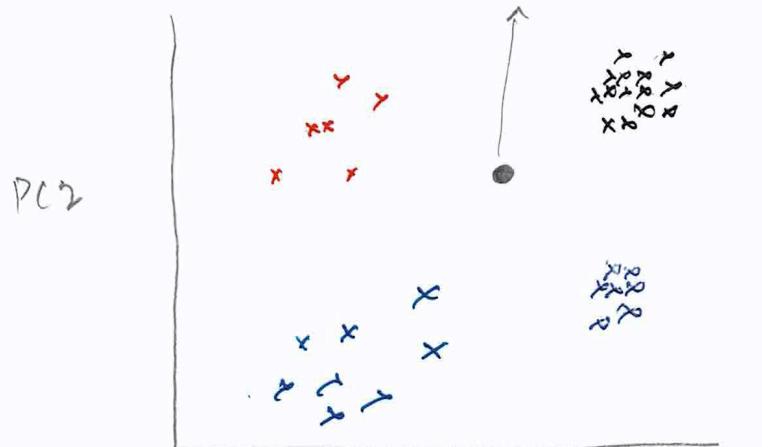
A simple way to classify data.

When $K=5$, we assign the point to the cluster category that gets "MAX. votes".

Red - 6, Blue - 2
Black - 14, Blue II - 7.

Since Black got most votes, it's assigned to Black cluster category.

When the point is b/w 2 categories -
if we can avoid. when $k = \text{odd number}$.



K-NEAREST NEIGHBORS.

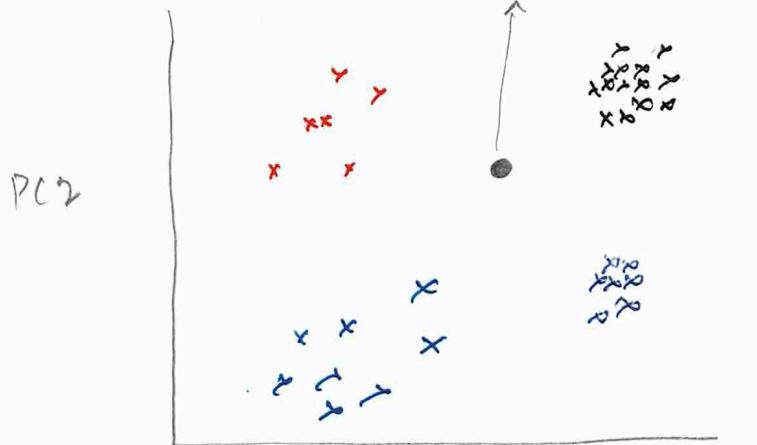
A simple way to classify data

when $K=5$, we assign the point to the cluster category that gets "MAX. votes".

Red - 6, Blue - 2
Black - 14, Blue II - 7.

Since Black got most votes, its assigned to Black cluster category.

When the point is b/w 2 categories -
if we can avoid. when $k = \text{odd number}$.



RANDOM FORESTS

- ① From the original 'df', randomly select samples \Rightarrow
Bootstrapped dataset.

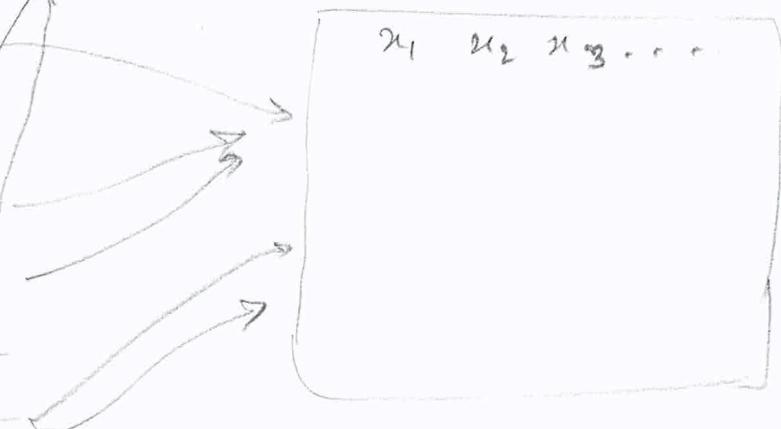
* A sample can be picked more than once.

Original df'

	x_1	x_2	x_3	\dots
y_1				
y_2				
y_3				
y_4				
\vdots				

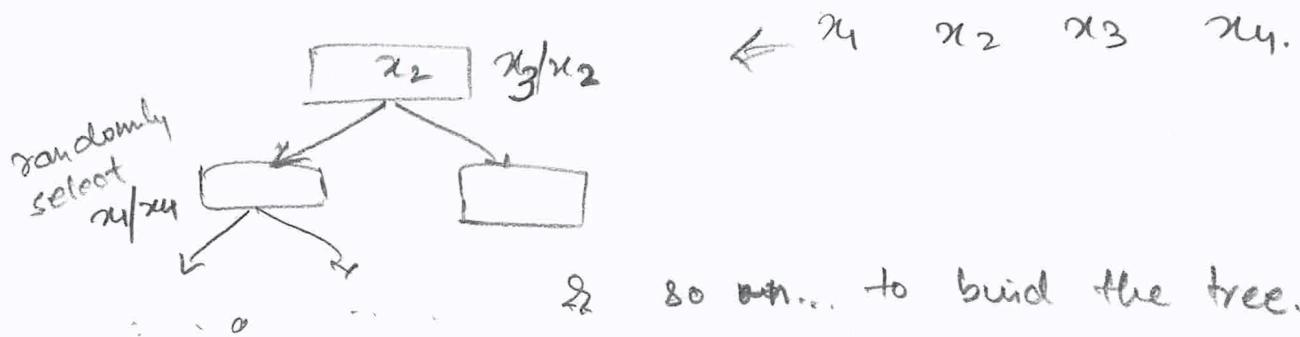
Out-of-Bag
dataset

Bootstrapped 'df'



typically
 y_3 of
data do
not end
up in
bootstrapp
df.

- ② Create a decision tree on bootstrapped df. but use only random variables (x_1, x_2, \dots) i.e. select only 'n' variables instead of all x_i .



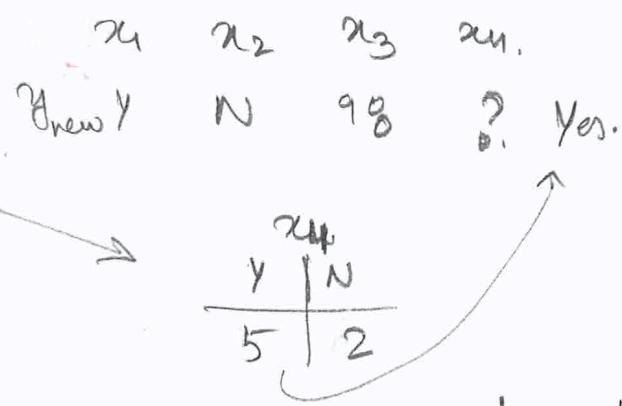
& so on... to build the tree. by
considering only random variables
at each node.



Using a bootstrapped sample & considering only a subset of variables at each node/step gives wide variety of trees.

The Variety make Random Forest more effective when compared to Decision trees.

We take this data & run it down all the trees in RF to & choose one with max. votes.



[Bootstrapping the data AND using aggregate to make decision is called BAGGING.]

The out-of-Bag was not used to create the trees ∵ this df is run across all trees that were built.

Accuracy of RF - proportion of out-of-Bag samples that were correctly classified.

The samples that were incorrectly classified \Rightarrow Out-of-Bag Error

How to choose no. of x initially? -

Compare Out-of-Bag Error for a RF built using 2 variables per step vs 3 variables per step... & choose the most accurate RF.

Summary:

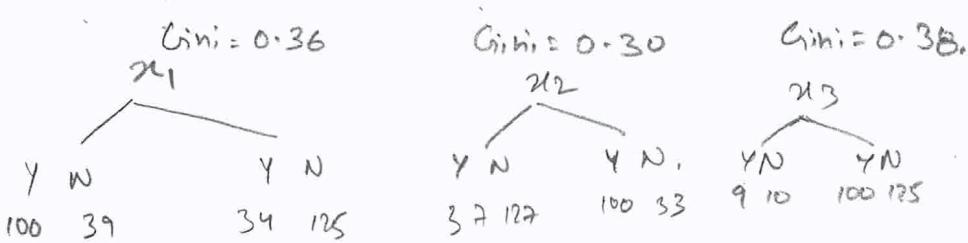
change no. of variables →

- 1) Build a RF (typically start with \sqrt{x} of no. of variables)
- 2) Estimate Accuracy of RF

DECISION TREE

A) Root nodes; Internal nodes; Leaf nodes;

Chest pain	BP	Blocked arteries	Heart Disease (Y)
x_1	x_2	x_3	
N	N	N	N
Y	Y	Y	N
Y	N	Y	N
.	.	.	.



1) How to determine the Root Node
 x_1, x_2, x_3 .

2) Take each x_i & calculate the "Gini Impurity" & pick the x_i with least Gini impurity value.

$\Rightarrow x_2 \Rightarrow$ Root Node.

- 3) To pick Internal node - repeat the same process & pick the x_i with least impurity.
- 4) If the Node itself has lowest impurity, it becomes Leaf node.

- B) Wt. (0.4)
- | | | |
|---|---|---|
| 5 | 4 | 1) Sort x_1 - low \rightarrow high. |
| 2 | 2) Cal. the avg. b/w x_i | |
| 4 | 3) Cal. the Gini impurity for each Avg. x_i | |
| 7 | | |

- C) for Ranked x_i like, 1-4, multiple choice \rightarrow Cal. the Gini impurity at each value of x_i & pick the lowest

Decision trees have downside of often being over fit.

- Missing data-
a) find the co-relation b/w the missing column & another column in df. Use linear Regression to predict the value of missing value; or choose mean/median
b) for categorical - choose the category with max. value; i.e in ($Y=8, N=10$), substitute with 'N'; or find the column with highest correlation & predict the value of missing data