

# Regression models

1) Simple Regression

$$y = b_0 + b_1 x$$

↑      ↑      ↗

Dependent var      Independent var

Y intercept (constant)

Slope (slope)

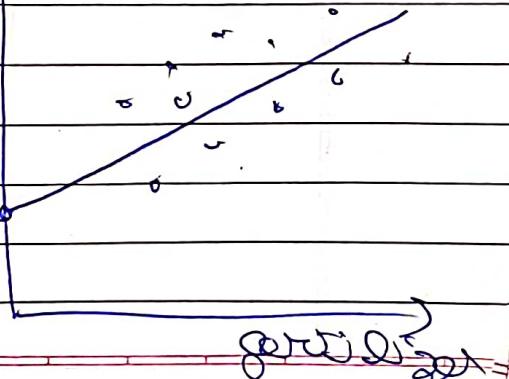
Suppose we want to predict  
height with Based on fertilizer  
used

$$\text{Predict} = b_0 + b_1 \times \text{fertilizer}$$

$$\text{Let } b_0 = 8 \text{ tonnes}$$

$$b_1 = 3 \text{ tonnes/ kg}$$

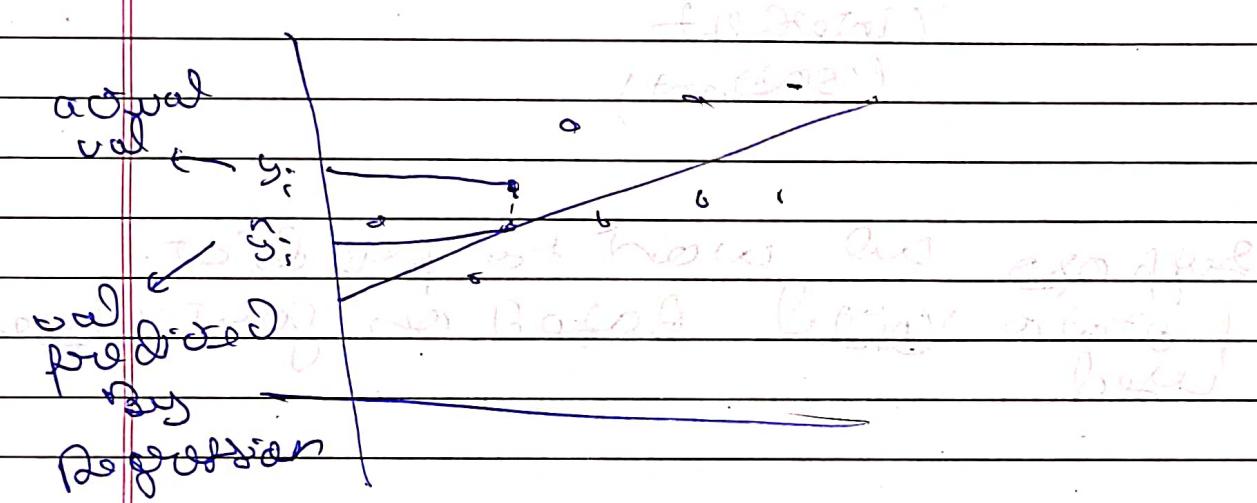
Graph will be like of



We can say that  
for every 1 t  
of fertilizer  
predicted is by 3 t

Now how to find best  
Regression Line

we do this by ordinary  
least square



standardized & fit to  $y_i - \bar{y}$  (residual)

∴ we need to find

$$\hat{y} = b_0 + b_1 x_1$$

$b_0, b_1$  such that

sum  $(y_i - \hat{y}_i)^2$  is minimized

→ multiple linear Regression

Let's say we have a  
big company we have  
there R & O Spend, Administration,  
Marketing and Sales, Production,  
and Profit

We want to build a model that will predict the profit of the company.

## Dawson

$$y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

↳ independent variables ↳ dependent variable

...and we stop along the road  
to take a look at the rock  
samples I have collected.

## Assumptions of Linear Regression

Exhibit No. 30 (Continued)

~~1. Random Error term~~

~~2. Nonlinear relationship~~

~~3. Autocorrelation~~

~~4. Heteroscedasticity~~

(1) ✓ (2) ✗

5. No multicollinearity

6. Normal distribution of error terms

7. Homoscedasticity

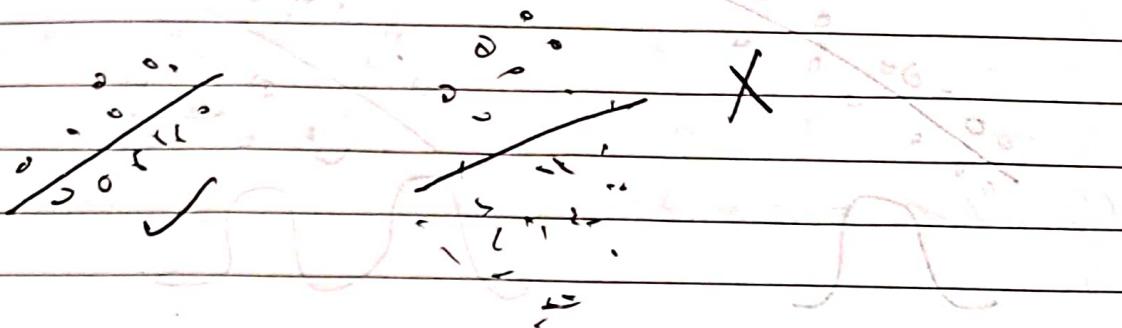
8. No autocorrelation

we should use linear regression in 2/3 and if

to make sure we can use linear regression and have some checks

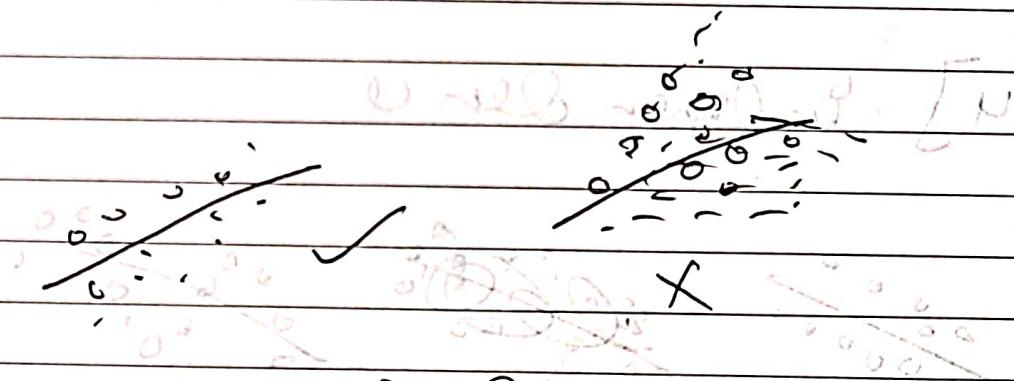
## 1 Linearity

(Linear Relationships between X and Y)



## 2 Homoscedasticity

(Equal variance)

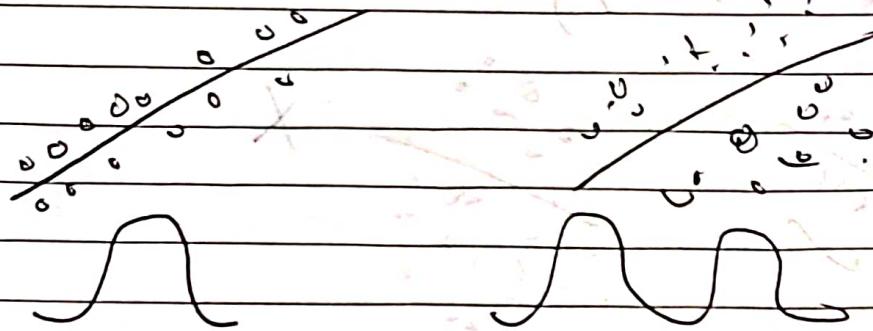


and don't want  
an inc or dec conc  
slope

voz sinquer denendon  
independent var  
Hence same res

### 3] multivariate normality

standardized residuals

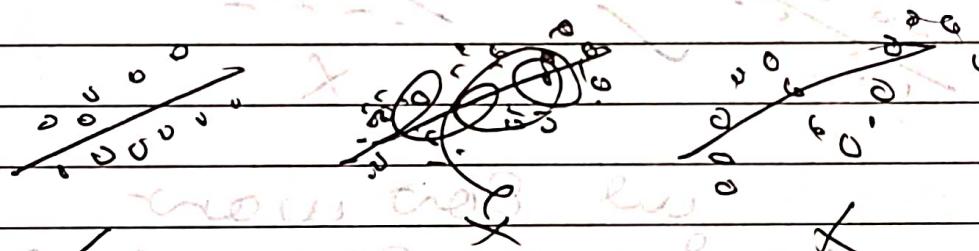


(we want normal)

distribution

(normality test)

### 4] independence



crossed lines

not correlated

independence

all over must be independent

they should not depend on each other

on each other

Ques 1 5) check for multicollinearity

$\rightarrow$  Jobs of multicollinearity

can we work with 192 variables

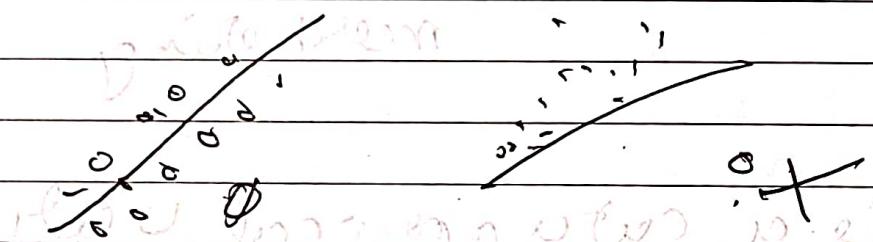
independent var not correlated with each other

$$x_1 \neq x_2$$

$$x_1 \sim x_2$$

↳ ~~↳ check for multicollinearity~~  $\times$  Diagnosis of

b) outlier check (optional)



(S) decide whether removal of outliers before or not (depends)

standard way handling outliers  
outlier point to produce  
regression line & diagnosis

## Dear Dummy variables

~~These are also called indicator variables~~

We have ~~Maximin to New York City~~ <sup>not</sup>

Project	R&D	admin	marketing	Finance
1	1	0	0	0
2	0	1	0	0
3	0	0	1	0

to build linear regression

$$y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots$$

~~constant & other required~~

R&D

admin

marketing

Stat is a categorical var  
So to take this we create  
dummy variable

Suppose we have  
2 categories in state  
newyork and california

So instead of strict rules  
a new column

P	Q & D	admrx	markets	NY	CA
1	1	1	1	1	1
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0

So we only include in NY  
NY and ignore CA and means  
not in NY

5: Regionen der

$$y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + b_4 D_1$$

this may seem biased but when many columns fix 0 the counts can be will adjust values of CA (Default val)

## Dummy variable trap

if we include both the dummy variables then there is no default constraint and also

$$D_1 = 1 - D_2$$

So by having

$$P_1 + P_2 \text{ we get } 1$$

$\rightarrow P$  value

consider a coin toss

$H_0$ : fair coin

null hypothesis

fair coin means it has not a guranteed

so it can be either head or tail

at 95% confidence level

it is 1% chance

assuming coin is fair

now we toss a coin and get tails

1 time - 50% significance

2 time - 25% confidence

3 time - 12% value

4 time - 6%  $\alpha = 5\% (0.05)$

5 time - 3% reject null

6 time - 1% if  $p < \alpha$  [hypothesis]

$H_0$  is wrong

very low p-value

these are the p values

if we assume not fair coins

1 time 100%

2 time 100% } p-value

6 time 100% } p-value

so if we want to reject  $H_0$  at 5% level then we need 20% p-value

→ Building a model

↳ Regression is done  
consider all the independent  
variables

(because

i) too many variables will  
not give optimal result

ii) model becomes too  
complex

→ Some methods

1] All-in

2] Backward elimination

3] Forward selection

4] Bidirectional elimination

Step-wise

Regression

### 5] Score comparison

→ All in

take all variables

if we have good prior  
knowledge then only  
use

2] it is compulsory

3] preparing for Bayesian  
elimination

→ Bayesian Elimination  
STEP - 1 :- Select significance  
level ( $\alpha = 0.05$ )

→ STEP - 2 :- go all in

STEP 3 :- consider  $H_0$  &  $H_1$   
p value

→ If p value is less than

then accept  $H_0$  else accept  $H_1$

if  $P > SL(\alpha)$  then STEP 4  
or Finish

STEP 4:- Remove var,  
without var with  
Highest p-value

STEP 5:- ~~Egg model~~  
without var

(go to Step 3)

Forward Selection

STEP 1:- Select ~~all~~  $\leq \alpha$  (e.g.  
var ( $\alpha = 0.05$ )

STEP 2:- Fit all  
possible simple regression  
and select one with  
highest p-value

STEP 3:- Keep in var  
and fit all models  
with initial predictor

Step 4:- Consider others  
P-value is less then  
go to Step 3 otherwise FIN  
(keep the previous model)

## Bi-Diretional Elimination

Step 1 :- Select a Design preference level to enter and to stay in the model

Design preference :  
SLENTER = 0.1, SLSTAY = 0.05

Step 2 :- Perform steps of primary selection

( P < S L entry + to enter )

Step 3 :- Perform steps of backward elimination

( P < S L entry + to entry )

If no new var can enter  
or limit then → FIN

all possible models

Step 7: Select the model  
of goodness of fit

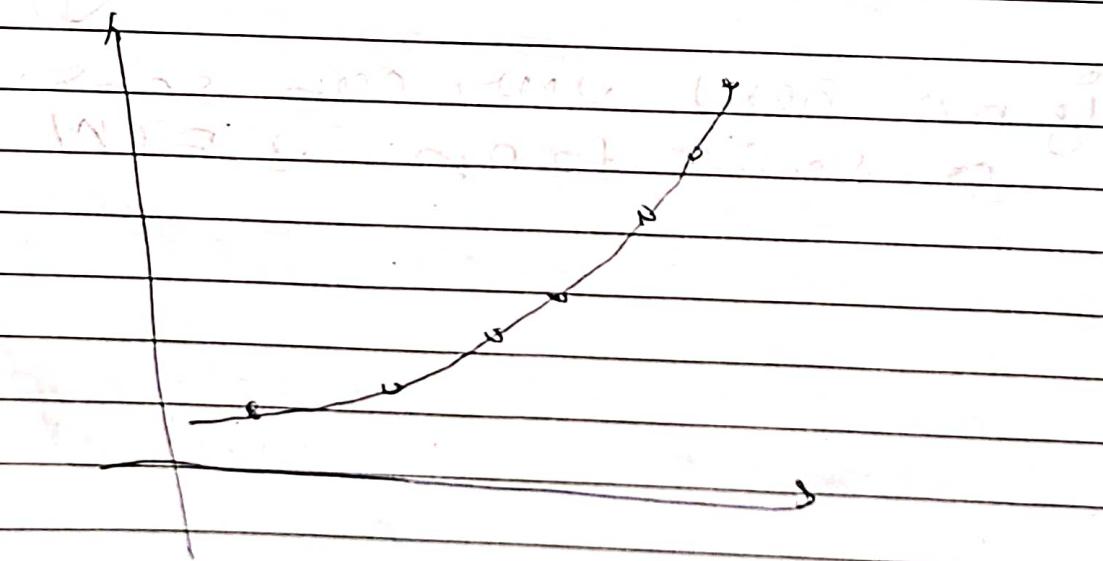
for linear regression all  
models

Step 8: Select the model  
of goodness of fit

→ Polynomial linear regression

→ Model selection based on  
all available Data set

→ Best fit model



a curve would lie between a line

operation

$$y = bx^2 + \cancel{bx} + \cancel{c}$$

why is it called Diagonal

it is called Diagonal because we always found on the graph below the x-axis

if coefficient of  $x^2$  is positive  
we call it Diagonal Regressor

What is Diagonal Regression?

$\rightarrow$  Polynomial Diagonal Regression  
is special case of multiple linear regression

$$y = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots$$

$$y = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots$$

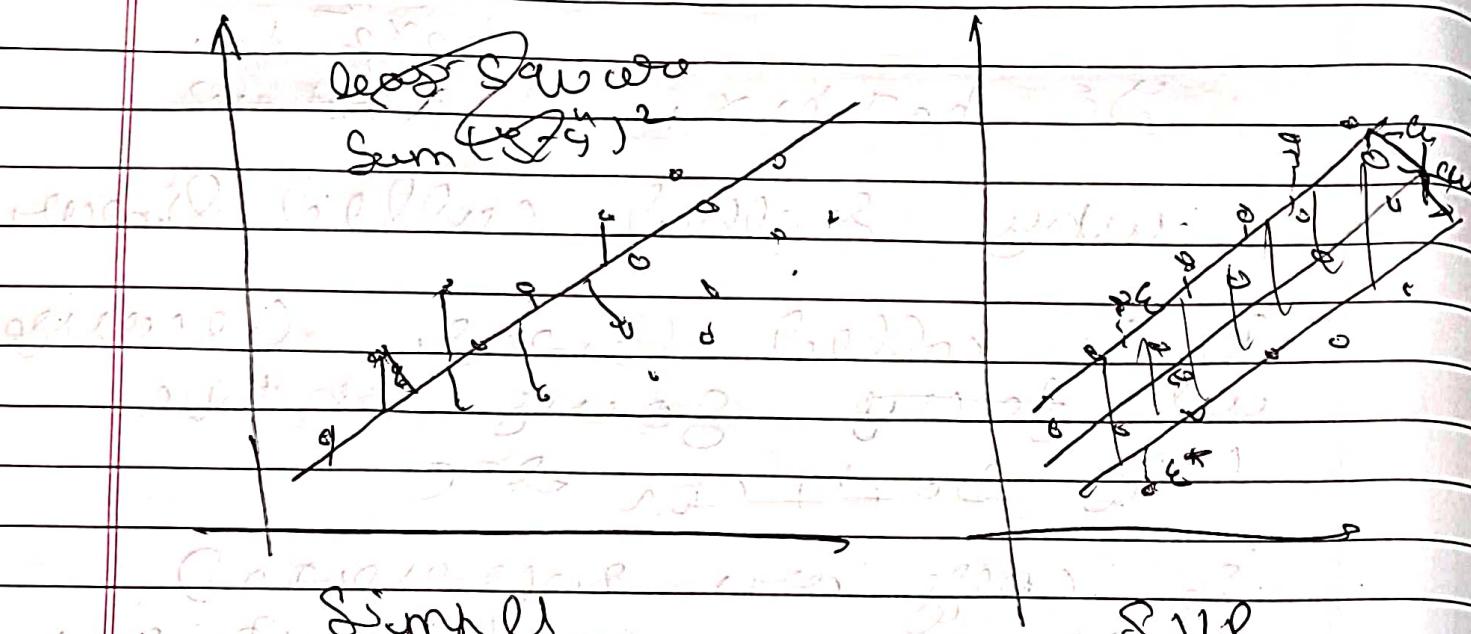
$$y = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots$$

$$y = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots$$

$$y = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots$$

$$y = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots$$

$\rightarrow$  Support us over progress  
(Easier [linear])  
 in all time



Best fit line of data are  
 same

Simple line we have a  
 tube called a  $\sigma$ -intercept  
 tube. Here  $\sigma$  we ignore  
 errors of data points in the tube (loss  
 error of least square method)

for linear reg we have

$$\text{Sum } \{(y - \hat{y})^2\}$$

for sum we have now

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^m (y_i - \hat{y}_i)^2 \rightarrow \min$$

Points  
above

Points  
below

Points

below

more linear fit

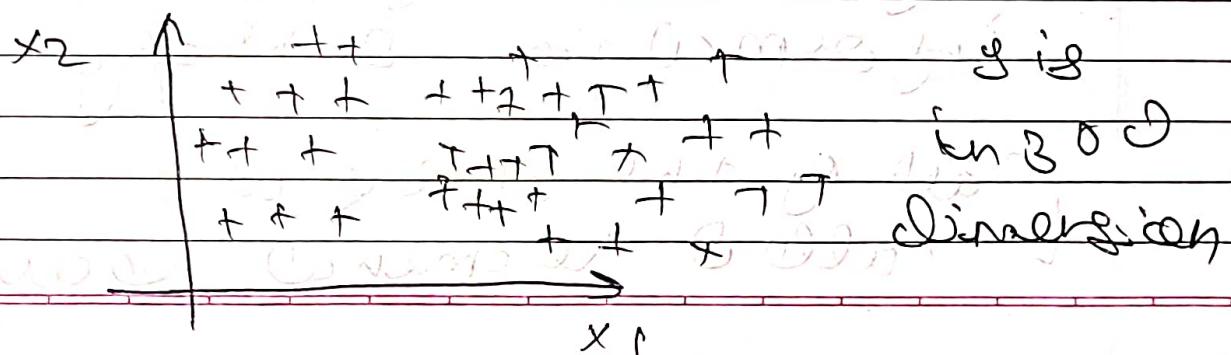
→ Decision Tree Regression

Suppose we have

independent variables

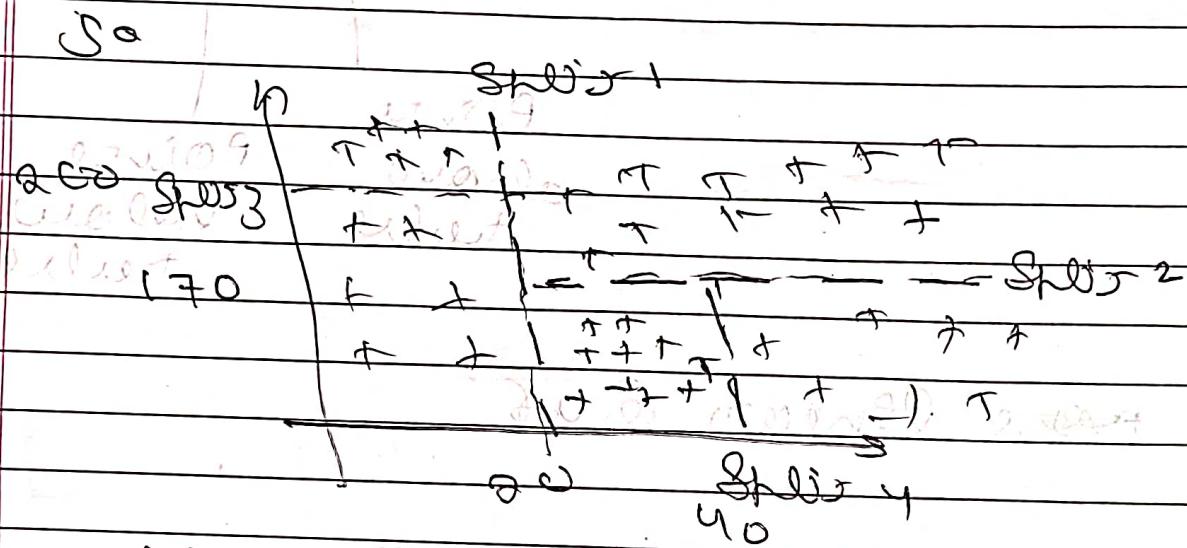
$x_1$  and  $x_2$  and we

need to predict  $y$



now a decision tree

algorithm will split the data points such that some information is added and it stops when information added is less than  $\epsilon$ .



(this is based on our)

is happening as we can see it's a tree

seems to make some

nodes general

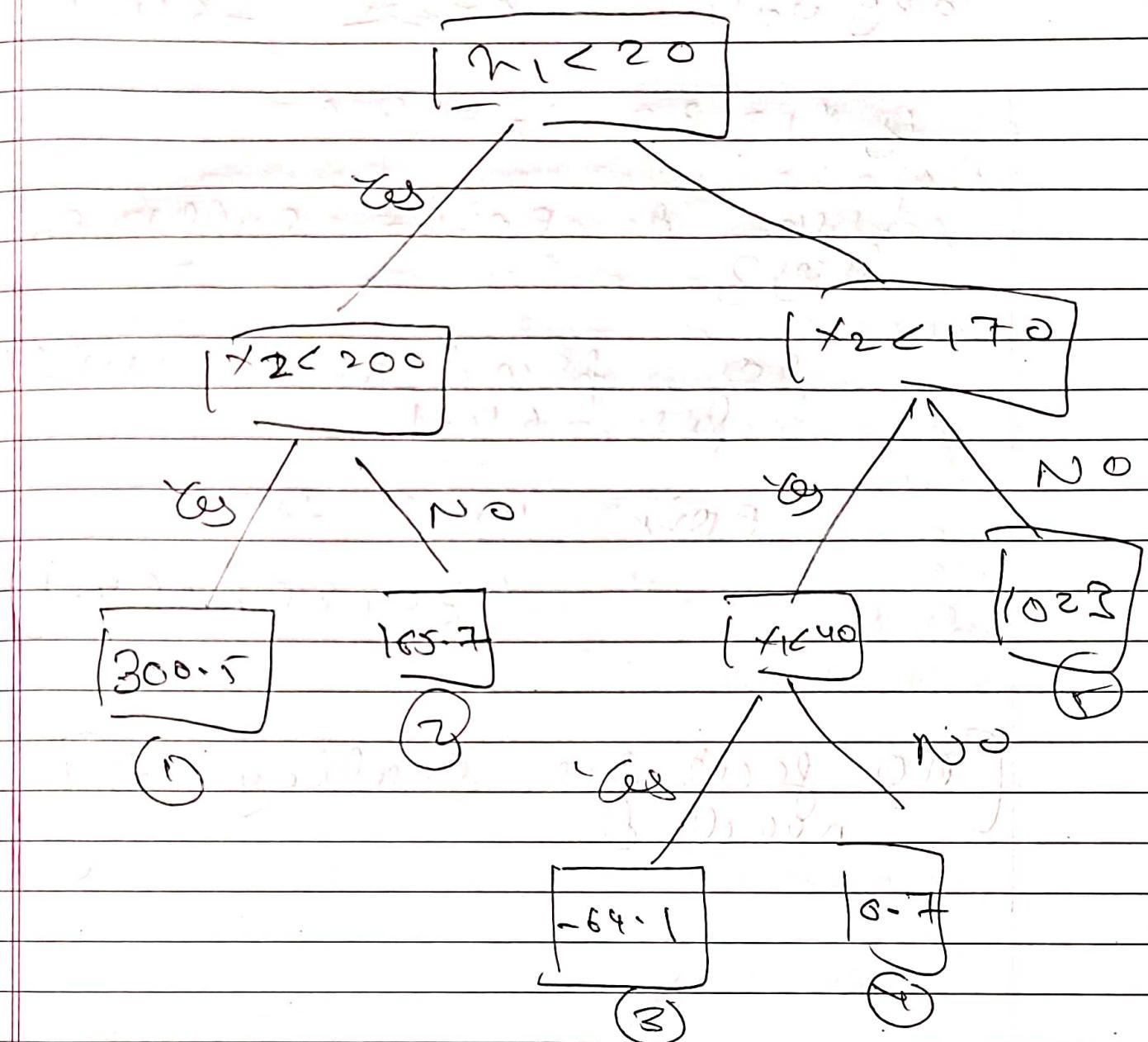
algorithm uses

information on the off

[all of the splitting are]

called technical jargons

## Decision tree



when we go to choose  
from any of the paths  
from root to leaf of  
we will get a float

now every leaf has  
a value = avg of all  
the data points

Let's say we have 3 vertices

avg value  $\rightarrow$  ~~65.5~~ -64.1

$$\text{if } x_1 = 30, x_2 = 20$$

then the point falls on  
line 3

so coordinates will  
be -64.1

point will be

$$(30, 20, -64.1)$$

[no feature scaling is  
needed]

→ Random Forest Regression

worse on Ensemble Learning

intuition behind

1] Pick Random K Data

Pairing from Training  
Data

2] Then Build Decision  
tree from these Random  
Point

3] we choose number  
Ntree & we want to  
build a build and  
repeat steps 1 and 2

4) for a new test point  
out of Ntree for it  
the value of Y then we  
take avg of all predicted  
values and assign that  
value to new test  
point

By decimal we can

500 + 500

Decimal is better than  
division tree because  
we don't rely on one  
tree for any chance  
in decimal we can  
use all trees but it is  
very difficult to impact  
on root of tree. And  
since we take any of  
multiple trees similarly  
impacts less