

Rapidly search topic

linear regression	1
multiple linear regression	14
gradient descent	21
types of gradient descent	27
polynomial regression	<u>33</u>
ridge regression	34
Lasso regression	42
elasticnet regression	46
logistic regression	47
classification matrix	63
softmax regression	70
polynomial logistic	75
Decision tree	93
ensemble learning	97
voting ensemble	99
bagging	100
random forest	103
adaboost	108
k-mean clustering	111
gradient boosting	122
stacking	
hierarchical clustering	127
kNN	130
SVM	132
naive bayes	138
Xgboost	146
DBSCAN	167
Thank you note :)	171

Machine Learning

D Linear Regression

Supervised learning

linear regression

simple linear regression

multiple linear regression

polynomial linear regression
regularization.

Let's say, we have data of CGPA, and placement package.

CGPA	Package
-	-
-	-
-	-
-	-

simple linear regression

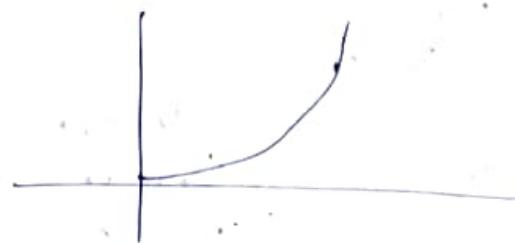
CGPA	Gender	Race	State	Package
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-

multiple linear regression

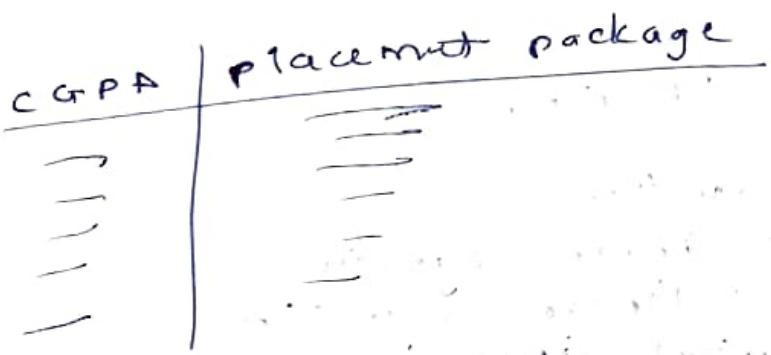
Linear



if we don't have linear data

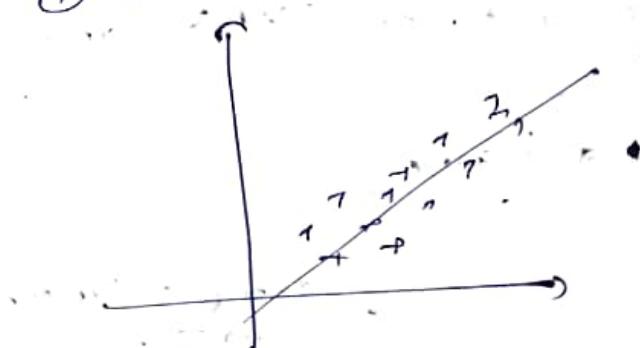


①

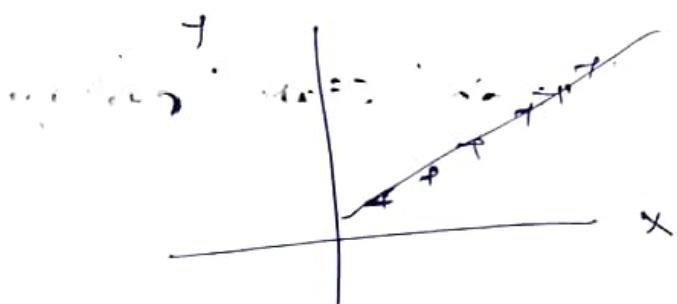


our task to predict package on basis
of CGPA.

① We need to plot the data



bit fit line
why not linear
it mostly doesn't
exist in real world.



$$y = mx + c$$

What it does?

It tries to cover all datapoints such
as to minimise error and closeness.

from sklearn.model_selection import
train_test_split

from sklearn.linear_model import
LinearRegression

lr = LinearRegression()

lr.fit(x-train)

lr.predict(x-test)

$m = LR.\text{coef}$ - weightage

$c = LR.\text{intercept}$ - offset in case if $y \neq mx$

techniques to find m and c :

To find m & b , there are two solutions

1) closed form solution

2) Non-closed form solution

Closed form solution?

If it is a mathematical form of expression, it can contain number, variable, operation functions, but no limits, differentiation and integration.

From example

$$\text{To solve } \frac{ax^2 + bx + c}{\text{roots}} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

In Non-closed form, we need to use approximation techniques to reach there.

In linear regression, there are two methods of calculating m and b .

1) OLS (ordinary least square)

and by sklearn in internal

2) Approximation (gradient descent)

Most of people explains it using gradient descent.

But why?

→ When we work with higher dimensions, to get answers from direct formula becomes difficult. There, so gradient descent becomes efficient technique

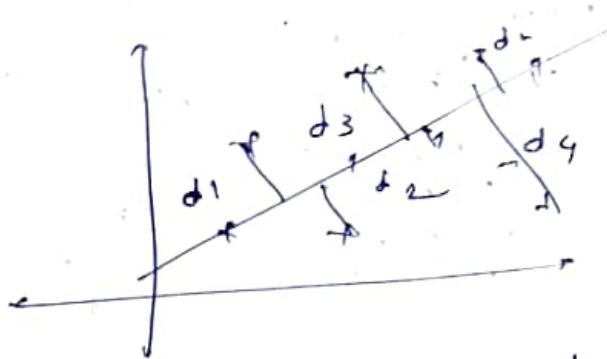
OLS - sklearn.LinearRegression()
gradient-descent - SGDRegressor

OLS direct formulae

$$m = \frac{\sum_{i=1}^n (w_i - \bar{w})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$c = \bar{y} - m\bar{x}$$

\bar{x} and \bar{y}
are means
 $i = \text{current row}$



$$\text{Error} = d_1 + d_2 + d_3 + d_4 + d_5 - \dots - d_n$$

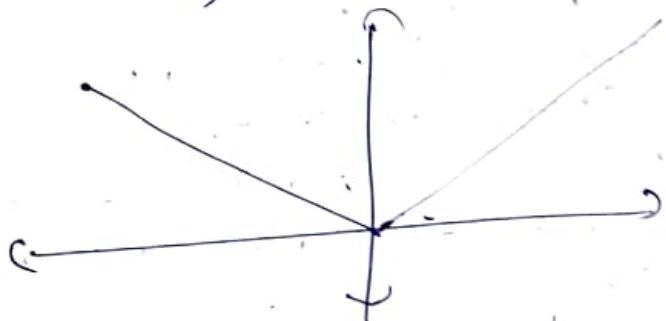
↓ To handle +ve & -ve both

$$\text{errors} = d_1^2 + d_2^2 + \dots + d_n^2$$

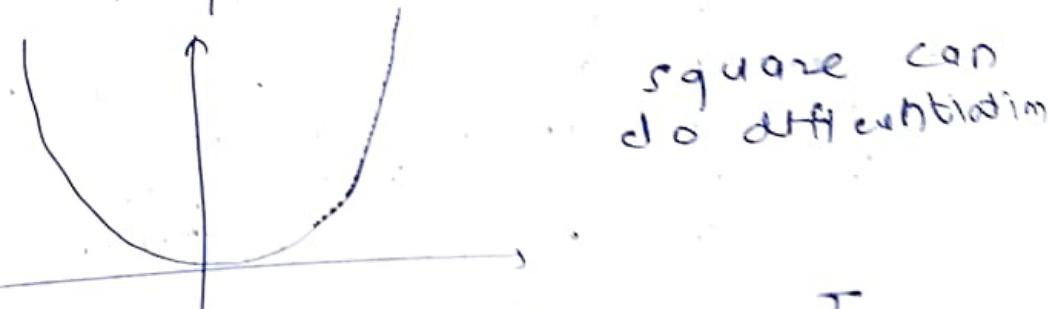
why not $|d_i|$?

outliers stabilization

Next step involve differentiation



mod graph
→ continuous
but not diff

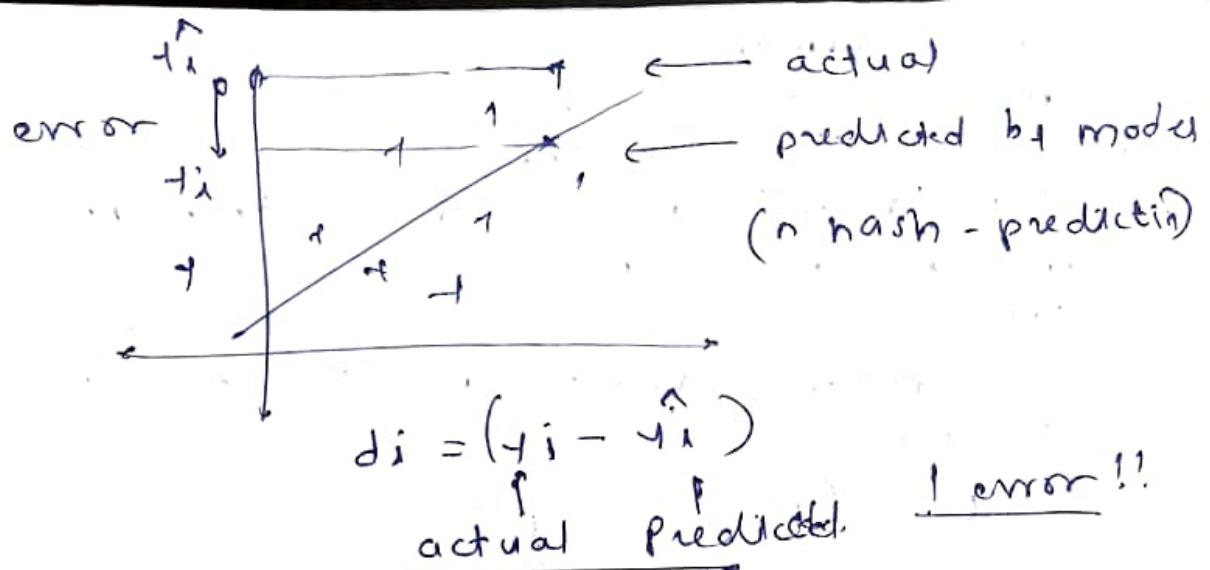


square can
do differentiation

$$E = \sum_{i=1}^n d_i^2$$

E
loss
function
error

I would say,
it needs closest pass to every point.
In fact I can say,
I need to find m and c such
that the value of eqn (E) should be
minimum.



total $E = \sum_{i=1}^n (d_i)^2$

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

Lines can be any. But we need only one with minimum E.

$\hat{y}_i = mx_i + b$: using line eqn

$E(m, b) = \sum_{i=1}^n (y_i - mx_i - b)^2$

so, we need m and b.

by changing m & b, we can increase & decrease error.

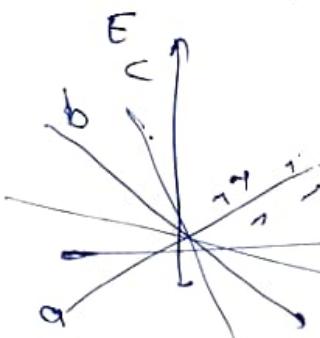
$y = f(x)$ → m means, after changing $y = f(x)$ we get y changed.

so, we can modify alone m; alone c, or alone both m and c.

we start with m,

we assume $c = 0$

$$E(m) = \sum_{i=1}^n (y_i - mx_i)^2$$



possible lines.

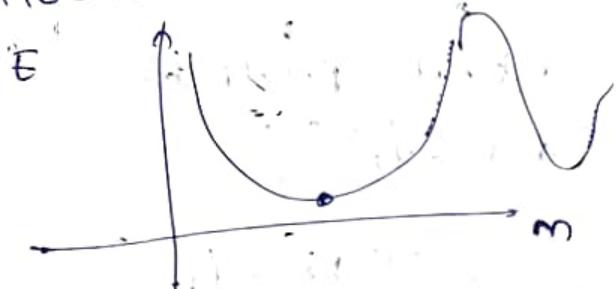
In around origin, we

can rotate 360°.

(5)

we say Reg Line.

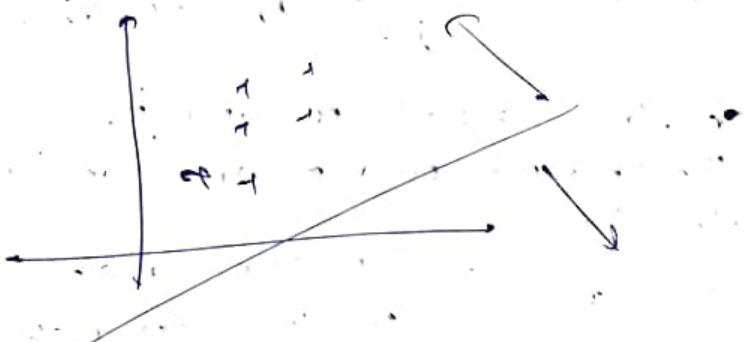
at assume
 $m = 10, E = 50$
 Now, wt I rotate line
 Line b will do less error than line a.
 Line a more more.
 wt a's more more.
 Line c will do less error than line b
 so, after rotating. m, E changes as
 follows.



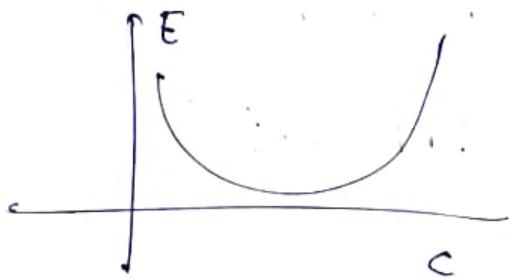
after a point, again errors will be increasing.

Now, wt go with c.

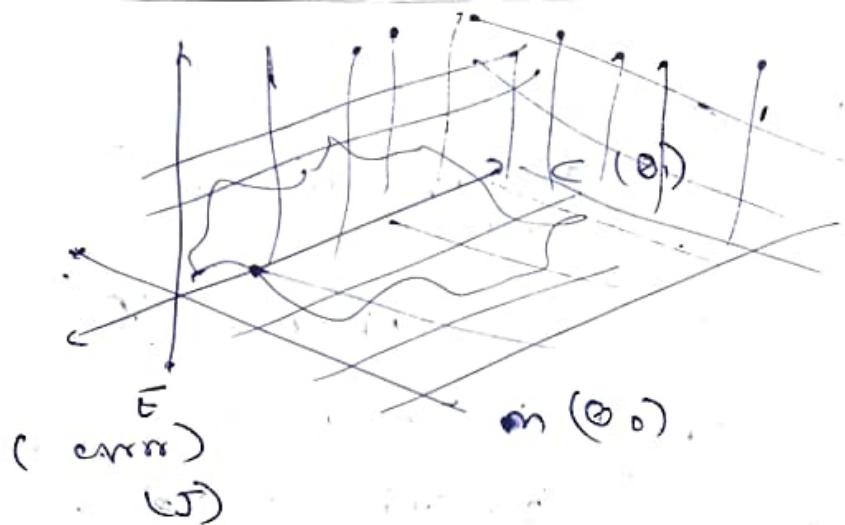
$$\therefore E(b) \doteq \sum_{i=1}^n (y_i - u_i - b)^2, (m = 1)$$



Now, we can either go up or down
 so start there calculate errors, more
 and minimize errors.



This is for c and E.



for certain values of m and b , error will be maximum. Similarly, for certain it will be minimum.

But how we can get maximums and minimums (heights and depths) for errors there?

-by maxima and minima in calculus.

On minimum, our slope becomes 0. for any equation of a line, we need to find its derivative, and put it to 0.

so, we will differentiate our ~~m~~ and error function $E(f)$, and (with respect to both m and c) & put it 0.

If it was an equation, which was only depended on x , then I, after taking,

$$E(x) \quad \frac{dE}{dx} = 0 \quad \text{if it was only dep. on } x.$$

but if $f(u, v)$,

so I use partial diff

$$\frac{\partial E}{\partial m} = 0 \quad \text{&} \quad \frac{\partial E}{\partial c} = 0$$

by this, I will get two eqns, & from them eqns, I will calculate m and c .

solve the eqns,

$$\frac{\partial E}{\partial c} = \frac{\partial \sum_{i=1}^n (y_i - mx_i - c)^2}{\partial c}$$
$$= \sum_{i=1}^n \frac{\partial}{\partial c} (y_i - mx_i - c)^2$$

apply chain rule

$$= \sum_{i=1}^n (2 \cdot (y_i - mx_i - c) \cdot \frac{\partial}{\partial c} (y_i - mx_i - c))$$

$$= \sum_{i=1}^n (2(y_i - mx_i - c)) \cdot (0 - 0 - 1)$$

$$= \sum_{i=1}^n (-2)(y_i - mx_i - c) = 0$$

$$\sum_{i=1}^n (y_i - mx_i - c) = 0$$

$$\therefore \frac{\sum_{i=1}^n y_i - \sum_{i=1}^n mx_i - \sum_{i=1}^n c}{n} = 0$$

$$= \frac{\sum_{i=1}^n y_i}{n} - \frac{\sum_{i=1}^n mx_i}{n} - \frac{\sum_{i=1}^n c}{n} = 0$$

$$= \bar{y} - m\bar{x} - \frac{n c}{n} = 0$$

$$\boxed{c = \bar{y} - m\bar{x}}$$

$$E_n = \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2$$

$$\frac{\partial E}{\partial \bar{m}} = \sum_{i=1}^n \frac{\partial}{\partial m} \left(y_i - mx_i - \bar{y} + m\bar{x} \right)^2$$

$$\begin{aligned} \frac{\partial E}{\partial \bar{m}} &= 2 \sum_{i=1}^n \left(y_i - mx_i - \bar{y} + m\bar{x} \right) \cdot \\ &\quad \frac{\partial}{\partial m} \left(y_i - mx_i - \bar{y} + m\bar{x} \right) \\ 0 &= 2 \sum_{i=1}^n \left(y_i - mx_i - \bar{y} + m\bar{x} \right) \\ &\quad (0 - x_i - 0 + \bar{x}) \end{aligned}$$

$$0 = (-2) \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})(x_i - \bar{x})$$

$$0 = \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})(x_i - \bar{x})$$

$$0 = \sum_{i=1}^n \left[(y_i - \bar{y}) - m(x_i - \bar{x}) \right] (x_i - \bar{x})$$

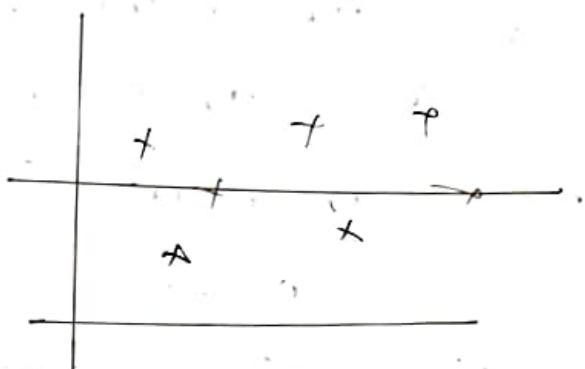
$$0 = \sum_{i=1}^n \left[(y_i - \bar{y})(x_i - \bar{x}) - m(x_i - \bar{x})^2 \right] = 0$$

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

$$\boxed{m = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Regression metrics

② MAE - mean absolute error



Simply the distance between actual and predicted.

$$MAE = \frac{|y_1 - \hat{y}_1| + |y_2 - \hat{y}_2| + \dots + |y_n - \hat{y}_n|}{n}$$

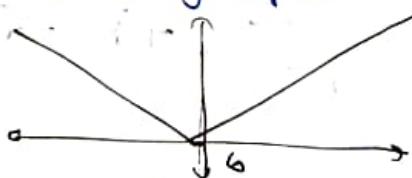
$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n}$$

Advantage

- ① Same unit
- ② Robust to outliers

Disadvantage

Absolute mod graph



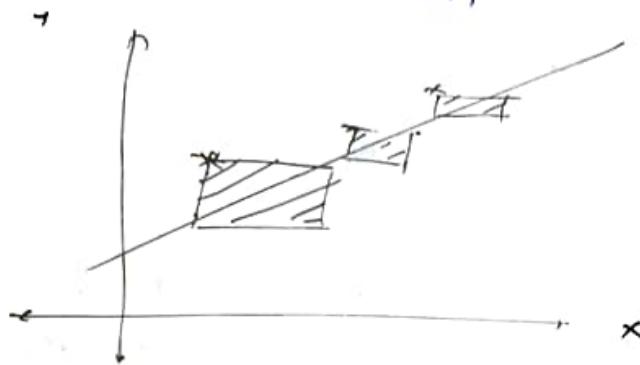
cannot differentiable at 0.

We apply optimization, so we need differentiation techniques

So that's why Mean squared error was introduced.

2) mean squared error (MSE)

$$MSE = \frac{\sum_{i=1}^n (y_i - f_i)^2}{n}$$



Advantage

can use as loss function, as it is differentiable.

Disadvantage

1) unit squares-

$y^2 \rightarrow y$ interpretation.

2) Not robust to outliers to square

3) RMSE (root mean squared error)

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{N}}$$

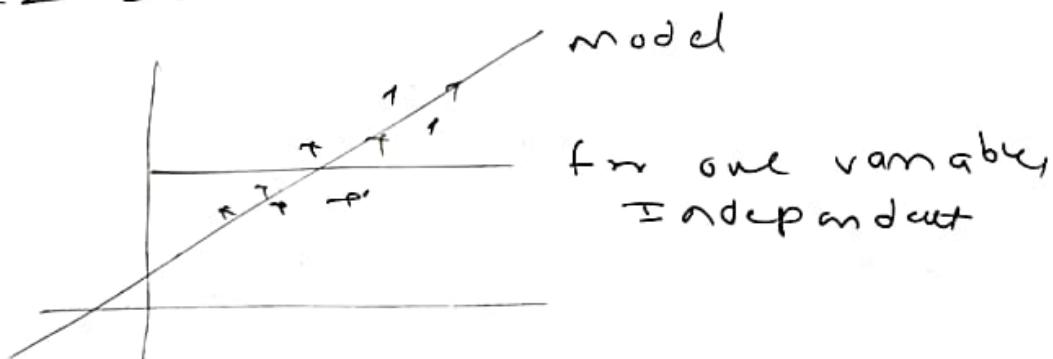
Advantage

same unit

Disadvantage

Not much robust to outliers.

4) R² score



when we calculate R^2 score, we just compare how much better our LR line, then mean line.

coefficient of determination /
Goodness of fit

$$R^2 \text{ score} = \frac{SS_R}{SS_m}$$

SS_R = sum of squared error in LR

SS_m = sum of squared error in mean.

$$R^2 \text{ score} = 1 - \frac{\left[\sum_{i=1}^n (y_i - \hat{y})^2 \right] \text{reg}}{\left[\sum_{i=1}^n (y_i - \bar{y})^2 \right] \text{mean}}$$

comparison of worst to the best

special case

1) case 1) R^2 score = 0

$$\frac{\left[\sum_{i=1}^n (y_i - \hat{y})^2 \right] \text{reg}}{\left[\sum_{i=1}^n (y_i - \bar{y})^2 \right] \text{mean}} = 1$$

both errors are same. No advantage of regression, As model is unable to pick up ~~diff~~ advantage of model features.

2) case 2) R^2 score = 1

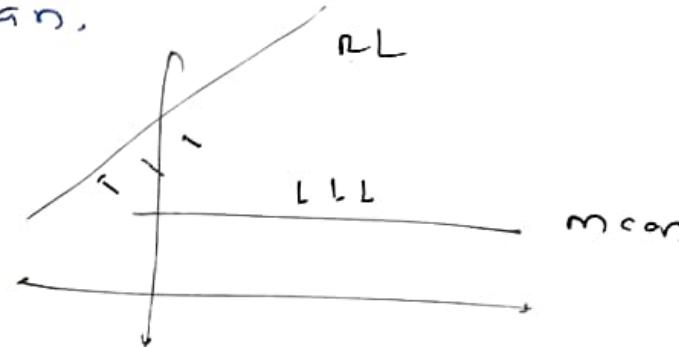
$$\frac{\left[\sum_{i=1}^n (y_i - \hat{y})^2 \right] \text{reg}}{\left[\sum_{i=1}^n (y_i - \bar{y})^2 \right] \text{mean}} = 0$$

means no error by regression. Perfect

case 3) Negative R₂ score

$$\text{if } SSR > SSM$$

means RL is very bad then mean.



Need to look at data & algorithm

Interpretation:

if say R₂ score is 0.8

it means our input columns are able to explain 80% of variance in target column. Rest 20% we don't have actually.

Advantage apply universally on all datasets

Disadvantage

if we add more columns in input, our R₂ score starts increasing.

Adjusted R₂ Score

$$R_{2\text{adj}}^2 = 1 - \left[\frac{(1-R_2)(n-1)}{(n-k)} \right]$$

R₂ = Normal R₂ score

n = total number

k = total no of input columns

Let us explain this by two cases:

case 1)

adding irrelevant column,

case 2)

adding relevant column

case - I

our R₂ score will either somewhat decrease or it will constant.

$$R_{2\text{adj}} = \frac{(1-R_2)(n-1)}{(n-1-k)} \downarrow$$

$$(1-i) \downarrow$$

so overall R_{2adj} ↓.

case - II relevant

$$R_{2\text{adj}} = \frac{(1-R_2)(n-1)}{(n-1-k)} \downarrow \curvearrowleft$$

$$1-i \downarrow \Rightarrow \uparrow$$

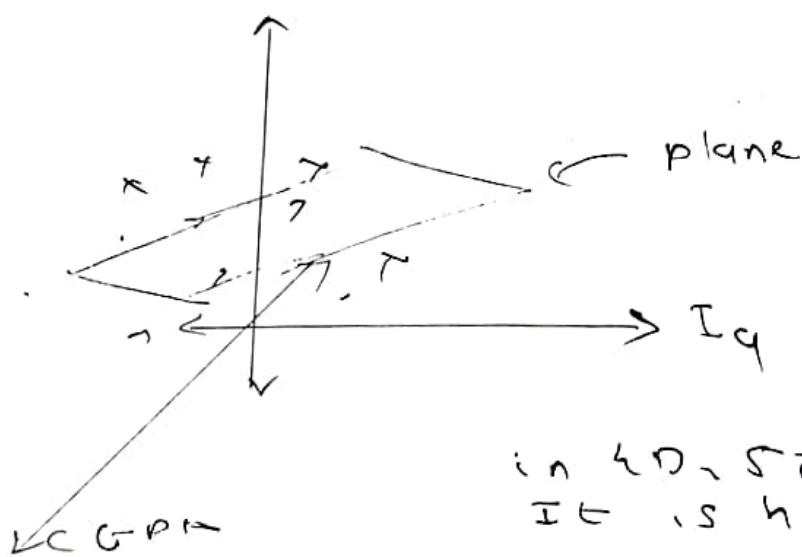
so overall increased .

Multiple Linear regression.

comes when there are many inputs.

Iq	subject	CRYPH package
-	-	-
-	-	-
-	-	-

let assume it takes 2. ignore sub
for now
package



in 4D, 5D ... nD
it is hyperplane

$\Rightarrow 2D$.

$$\Rightarrow y = mx + c$$

$\Rightarrow 3D$

$$y = m_1x_1 + m_2x_2 + c$$

$$\therefore y = \beta_0 + \beta_1x_1 + \beta_2x_2$$

need to find $\beta_0, \beta_1, \beta_2$.

$$y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \dots + \beta_nx_n$$

$$y = \beta_0 + \sum_{i=1}^n \beta_i x_i$$

n column = $n+1$ coefficient

coefficient = weight.

less weight means no much contribution

high weight means significant contribution

Mathematical Intuition

cgoq	id	send of package	
x_1	x_2	x_3	y

$$\hat{y} = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3$$

$$Y = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \beta_0 x_{11} & \beta_0 x_{12} & \dots & \beta_0 x_{1m} \\ \beta_1 x_{21} & \beta_1 x_{22} & \dots & \beta_1 x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_n x_{n1} & \beta_n x_{n2} & \dots & \beta_n x_{nm} \end{bmatrix}$$

$$\hat{y}_i = \beta_0 + \beta_1x_{i1} + \beta_2x_{i2} + \dots + \beta_nx_{in}$$

m columns

$$\begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1m} \\ 1 & x_{21} & x_{22} & \dots & x_{2m} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nm} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}$$

$$\boxed{\hat{Y} = X\beta} \quad \overbrace{\quad}^{\text{prediction}} \quad \overbrace{\quad}^{\text{coefficient matrix}}$$

at all rows

(15)

$$X = \begin{array}{c|cc|c} & x_{11} & x_{12} & x_{13} \\ \hline 1 & x_{11} & x_{12} & x_{13} \\ 1 & x_{21} & x_{22} & x_{23} \\ \hline 1 & x_{31} & x_{32} & x_{33} \end{array}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$\epsilon = Y - \hat{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$\epsilon = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$$

In simple linear regression

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

so we can write ϵ as
at transpose e,

$$\begin{bmatrix} y_1 - \hat{y}_1 & y_2 - \hat{y}_2 & \cdots & y_n - \hat{y}_n \end{bmatrix} \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$$

(1 x n) (n x 1)

Scalar

$$= (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \cdots + (y_n - \hat{y}_n)^2$$

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

$$\boxed{\epsilon \doteq \epsilon^T \epsilon}$$

Now:

$$E = e^T e \\ = (e - e^T)^T \cdot (e - e^T)$$

$$(A + B)^T = A^T + B^T \\ (A - B)^T = A^T - B^T$$

$$= (e^T - e^T)^T (e - e^T)$$

$$= [e^T (x_B)^T] (e - x_B)$$

$$= e^T e - \underline{e^T x_B} - \underline{(x_B)^T e} + (x_B)^T x_B$$

To prove



$$= \frac{e^T x_B}{A^T B} = (x_B)^T e$$

$$A^T B = B^T A \quad \text{--- (1)}$$

$$\text{let say } (A^T B)^T = B^T (A^T)^T$$

$$\text{RHS eqn (1)} = B^T A \quad \text{--- (2)} \\ \text{RHS eqn (2)} = \text{eqn (1) RHS}$$

in fact, we say

$$A^T B = (A^T B)^T$$

$$\text{we say, } A^T B = C$$

$$A^T B = C \quad \leftarrow \text{To prove}$$

$$A^T B = \underline{\underline{e^T x_B}}$$

$$(e^T x_B)^T = (e^T x_B) \quad \text{--- (3)}$$

$$\left[\begin{array}{c} (1 \times n) \\ (n \times (m+1)) \end{array} \right] \quad \left[\begin{array}{c} (m+1) \times 1 \\ ((m+1) \times 1) \end{array} \right]$$

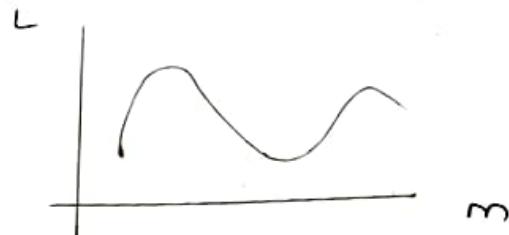
$$\left[\begin{array}{c} 1 \times (m+1) \\ 1 \times (m+1) \end{array} \right] \quad \left[\begin{array}{c} (m+1) \times 1 \\ (1 \times 1) \end{array} \right] = (1 \times 1) \Rightarrow (1 \times 1)$$

(17)

Coming back at eqn v

$$E = \frac{y^T y - 2 y^T \times \beta + \beta^T \times x^T \times \beta}{1055} \text{ function}$$

like in SLR, we need to diff it
to find slope,



similarly, we will find $\frac{\partial E}{\partial \beta}$.

$$\begin{aligned} \frac{\partial E}{\partial \beta} &= \frac{d}{d \beta} [y^T y - 2 y^T \times \beta + \beta^T \times x^T \times \beta] \\ &= 0 - 2 y^T x + \frac{d}{d \beta} [\beta^T \times x^T \times \beta] \end{aligned}$$

matrix differentiation

Going in depth

let us say Ax .

A vector
Under Transformation

or say

$$A = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix}$$

$$Ax = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$f_1(x_1, x_2) = 1(x_1 + 2x_2)$$

$$f_2(x_1, x_2) = 2(3x_1 + 4x_2)$$

How to find

$$\frac{d}{dx}(Ax) = ?$$

$$\frac{\partial F}{\partial x} = \begin{bmatrix} \frac{\partial F_1}{\partial x_1} & \frac{\partial F_1}{\partial x_2} \\ \frac{\partial F_2}{\partial x_1} & \frac{\partial F_2}{\partial x_2} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 2 \\ 3 & 5 \end{bmatrix} = A$$

$$\therefore \frac{\partial}{\partial x} Ax = A$$

so talk about $x^T Ax$

$$= (x_1 \ x_2 \ \dots \ x_n) \begin{pmatrix} a_{11} & a_{12} & a_{1n} \\ a_{21} & a_{22} & \vdots \\ a_{n1} & \vdots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n \end{cases}$$

let now try to just look it as
symmetric. suppose $a_{21} = a_{12}$ say A

so we say

$$(x_1 \ x_2 \ \dots \ x_n) \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & & & \\ a_{n1} & \dots & a_{nn} \end{pmatrix}$$

if we say $n = 2$,

$$a_{11}x_1^2 + a_{12}x_1x_2 + a_{21}x_1x_2 + a_{22}x_2^2$$

(19)

$$= a_{11}x_1^2 + 2a_{22}x_2 + a_{22}x_2^2$$

$$= +(x_1, x_2)$$

$\text{So, } \frac{\partial}{\partial x} (x^T A x) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{pmatrix}$

$$\therefore \begin{pmatrix} 2a_{11}x_1 + 2a_{12}x_2 \\ 2a_{21}x_1 + 2a_{22}x_2 \end{pmatrix}$$

$$= 2 \begin{pmatrix} a_{11}x_1 + a_{12}x_2 \\ a_{21}x_1 + a_{22}x_2 \end{pmatrix}$$

$$= 2 \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$= 2 A X$$

$$\frac{\partial}{\partial x} 2 A X = \frac{\partial}{\partial x} x^T A x$$

come back to step

$$= 0 - 2 Y^T X + \frac{\partial}{\partial \beta} \left[\frac{B^T X^T X B}{X^T X} \right]$$

$$= -2 Y^T X + 2 X^T X \beta^T$$

$$\therefore \beta X^T X \beta^T = 2 Y^T X$$

$$\beta^T = Y^T X \cdot \frac{1}{X^T X}$$

$$\beta^T = Y^T X (X^T X)^{-1}$$

$$(B^T)^T = (\cancel{Y^T X})^{-1}$$

$$[(X^T X)^{-1}]^T (Y^T X)^T$$

$$\beta = [(x^T x)^{-1}]^T x^T y$$

$$\boxed{\beta = (x^T x)^{-1} x^T y}$$

$x = \begin{cases} x_{\text{train}} \\ \vdots \\ x_{\text{train}} \end{cases}$

$y = y_{\text{train}}$

$$\begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{bmatrix} = \frac{[(m+1) \times (m+1)]}{[(m+1) \times n]} [(m+1)(n)]$$

$$[(m+1) \times 1] = [(m+1) \times n]$$

This was OLS.

still why gradient-descent?

Increase time complexity $= n^3$.

a lot of computing. lot of time

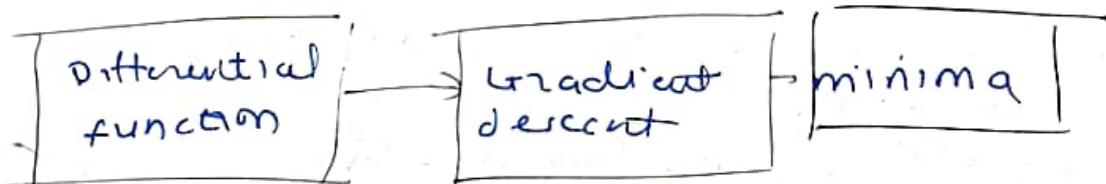
OLS = linear regression

SGD regressor = high dimensions.

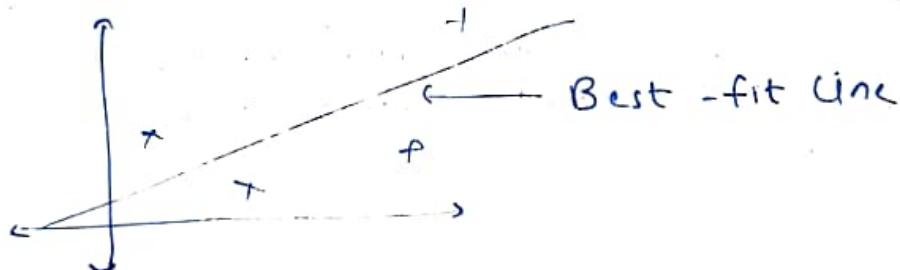
Gradient - descent

A technique to find out minima effectively.

- 1) Linear reg
- 2) Logistic reg
- 3) T-smth reg
- 4) DL backbone



Let I have a dataset



we can do it by OLS

$$\hat{y} = mx + c$$
$$L_i = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$
$$= \sum_{i=1}^n (y_i - mx_i - c)^2$$

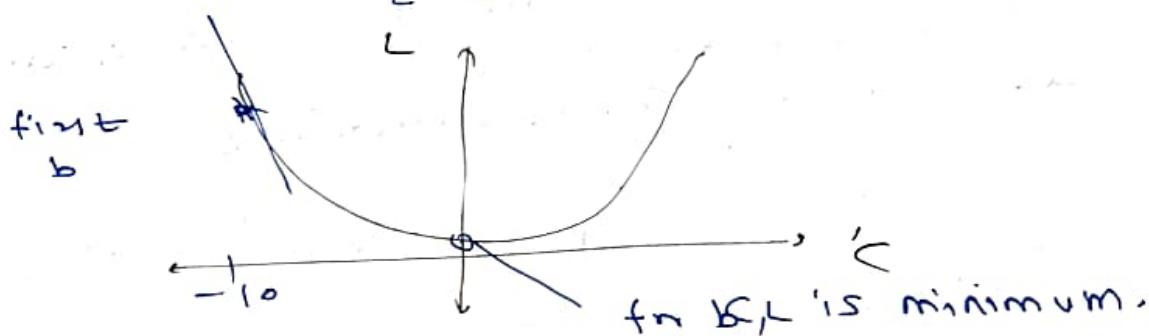
Depends on m & c both.

let I assume that I know m ,
so let it be 10.

$$L = \sum_{i=1}^n (y_i - 10x_i - c)^2$$

Now, I just simplified it,
we need to find such c by
adjusting height that do minimum
error.

$$L \propto c^2$$



we can find it by OLS, but can't in
higher dimension.

so gradient descent

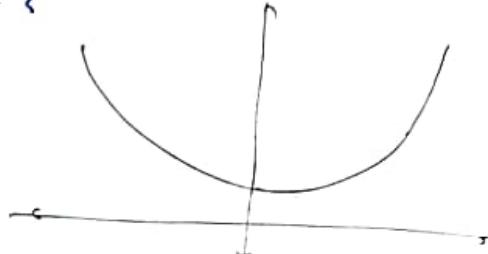
① Select a value for c . (random)

$$or b = -10$$

Algorithm have two direction for going
towards minima either upper or lower.
basically increasing or decreasing c .

we will simply find slope

now to find slope at a point to a
curve?



$$y = mx + c$$

$$\frac{dy}{dx} = m \text{ (slope)}$$

Step 2

so, find slope of that point w/ curve
find $\frac{dy}{dx} = m$, if put x

if

$m < 0 \rightarrow$ increment c

$m > 0 \rightarrow$ decrement c

$$C_{\text{new}} = C_{\text{old}} - m$$

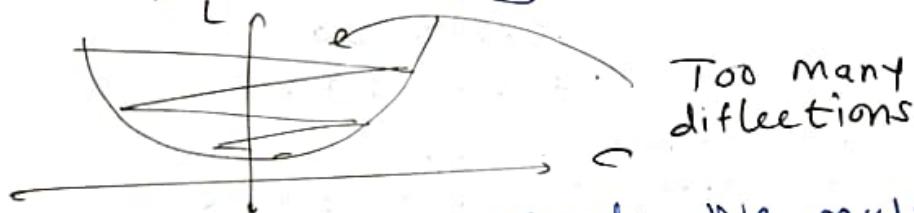
gradient-descent equation.

back to the definition repeated steps
The idea is to take opposite direction
of the gradient. (derivative)

This will happen in loop.

But a point to be noted is,

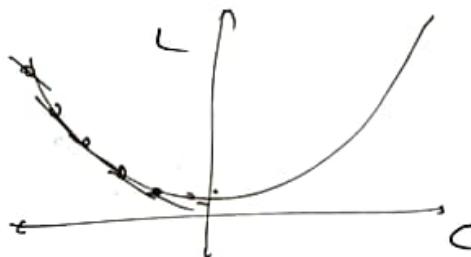
as slope is changing whole value,



To handle this tragedy, we multiply
slope with a thing called learning rate

$$C_{\text{new}} = C_{\text{old}} + \eta m$$

η - learning rate = usually 0.01



if we take learning rate (η) as

so small = slow speed

so large = jumps and can skip eqn.

and this is an Iterative algorithm

when to stop

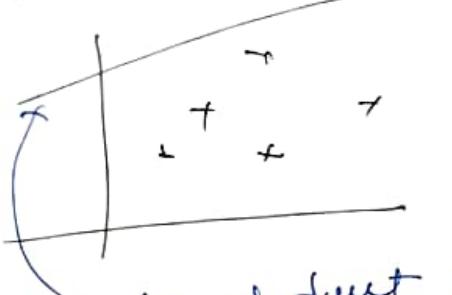
$$C_{\text{new}} - C_{\text{old}} = \frac{0.00001}{\text{A very small number}}$$

Incent / Percent

(23)

so this was approach 1.
approach 2 : limit iteration.
 Iteration = epochs

mathematical intuition



$$m + c.$$

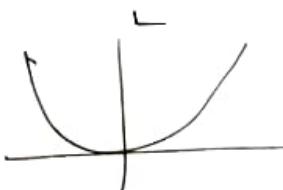
we know m ,
 we don't only need c

by default we

1) start with a random of c .

2) for i in epoch
 $b_{new} = b_{old} - \eta \text{ slope}$
 η - learning rate
 at say 0.01.

In first 100P,
 we have to find b_{new} .
 we have b_{old} , that original value.
 But how to find slope?

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


from this eqn we will find slope eqn
 & let say we have a c value.

$$\frac{dL}{dc} = \frac{d}{dc} \left(\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right)$$

$$= \frac{d}{dc} \left(\sum_{i=1}^n (y_i - mx_i - c)^2 \right)$$

$$= 2(y_i - mx_i - c) \frac{d}{dc}(y_i - mx_i - c)$$

$$= 2(y_i - mx_i - c)(0 - 0 - 1)$$

$$= -2 \sum_{i=1}^n (y_i - mx_i - c)$$

Let us say c is 5.

now, we can write,
as $m = 10$ already assumed

$$= -2 \sum_{i=1}^n (y_i - mx_i - 5)$$

after this, we get our slope at $b=5$.

$$c_{new} = C_{old} - \eta \text{ slope } (c=5)$$

do it until n times.

Now, let us find m .

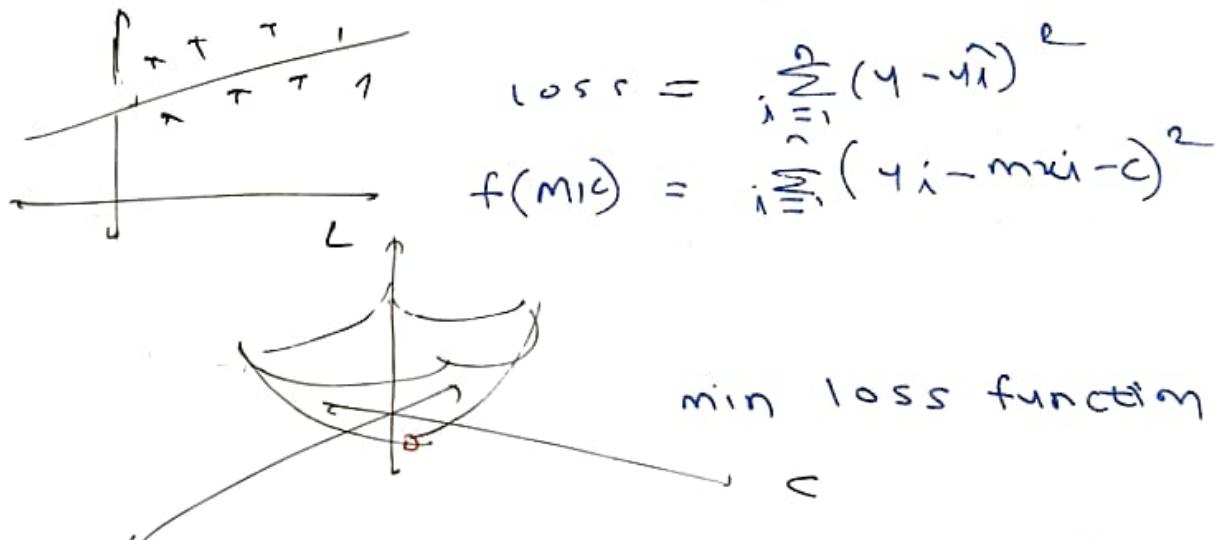
step 1) Initialize for $m \neq c$ as random.

step 2) decide epoch & learning rate

for in in epochs:

$$b \leftarrow b - \eta \cdot \text{slope}$$

$$m \leftarrow m - \eta \cdot \text{slope}$$



we are applying gradient descent in both's terms.

In terms of c and m ,
slope has two components, $m \neq c$.
so we need diff wrt. m & wrt c .

m wrt dr is m direction

c wrt dr is c direction

It is a combination of two diff

variables.

$$\frac{\partial L}{\partial b} \text{ for } b \quad \left| \quad \frac{\partial L}{\partial m} \text{ for } m$$

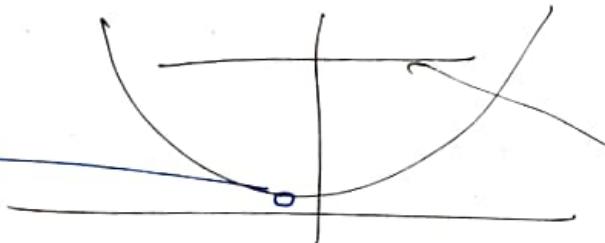
$$\begin{aligned}
 & \sum_{i=1}^n (y_i - w_i - c) \\
 & \text{w.r.t } w_i + c \\
 \frac{\partial L}{\partial c} &= \sum_{i=1}^n (y_i - w_i - c) \\
 &= -2 \sum_{i=1}^n (y_i - w_i - c) \\
 &\text{we can take } \\
 & w=0 \text{ as a start}
 \end{aligned}
 \quad
 \begin{aligned}
 & \frac{\partial L}{\partial w_i} = 2 \sum_{i=1}^n (y_i - w_i - c) \\
 &= -2 \sum_{i=1}^n (y_i - w_i - c) w_i
 \end{aligned}$$

effect of too much high or low learning rate
 less learning rate - time to train
 more learning rate - To jumps, no reach.

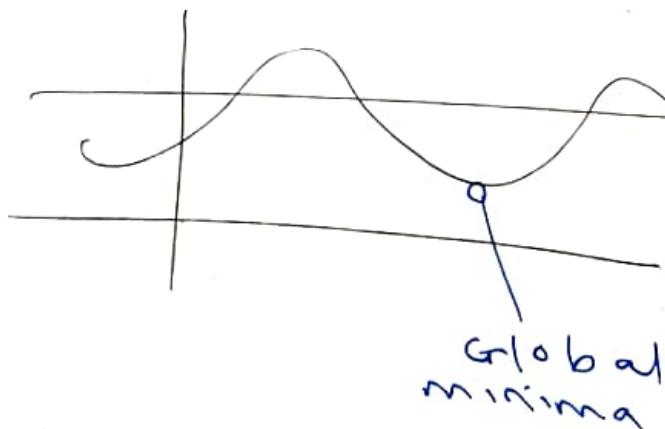
effect of loss function

$$L = \sum (y - \hat{y})^2 \leftarrow \text{convex fun'}$$

minima
only one
of global

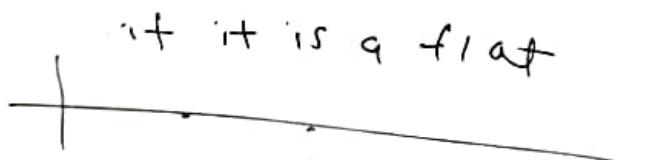


never cross line
so convex



local
minima
non convex
function,
crossed line
Global
minima

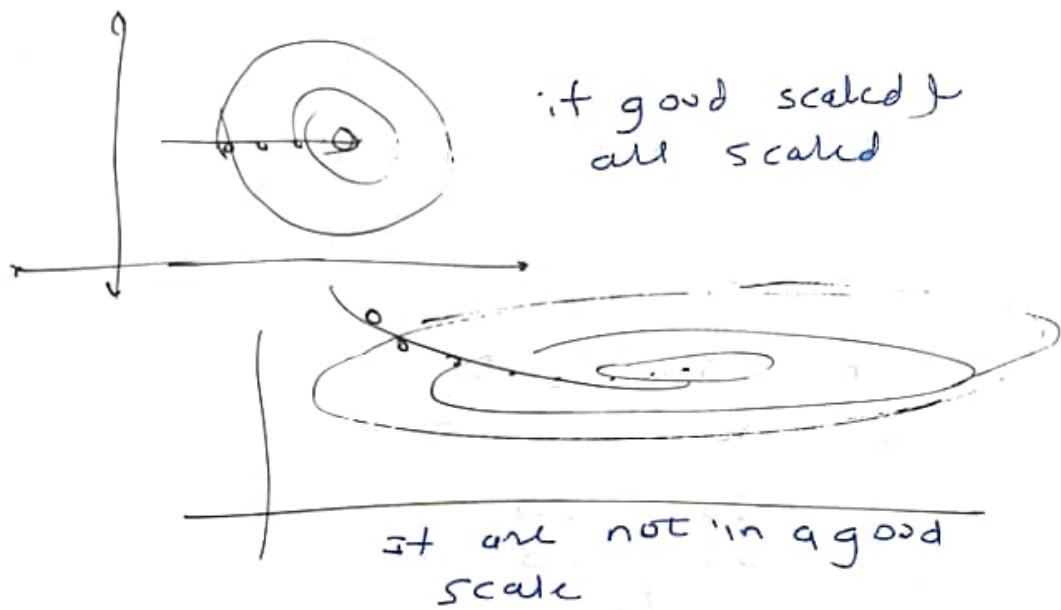
secondly,



less desired
epoches can't let reach you to
place.

effect of data

if columns in are same scale, we will depend on high speed.



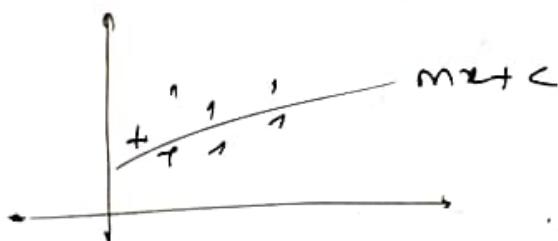
Types of gradient descent

batch GD

stochastic GD

mini batch GD

basis of the types



$$\frac{\partial L}{\partial m}$$

$m = 1, b = 0$
iterate
 $m_n = m_0 - \eta \text{ (slope)}$
 $b_n = b_0 - \eta \text{ (slope)}$

$\frac{\partial L}{\partial b}$

In this we look all rows
literally every row.

This is a batch gradient
descent.

which is slow.

To prevent this computational problem
subset rows we just look a single
row.

fast, suitable for larger datasets.

and if we say, we wish mid of batch
and stochastic then,
it is a mini batch.

we define a batch size then
update m & b according to m batch
size

mini-batch is somewhat rare.
 batch is very rare and sometimes for
 smaller datasets, and convex loss
 function
 stochastic is for large dataset, and is
 less common.

The past was one of batch GD, but it was
 2D.
 Let us study 2D.

Let us start small.

x_{12}	x_{11}	x_{21}	x_{22}	LPA
	cgaq	19		
	2	2		4
	8-1	23		3.2
	7-5	95		3.5

$$y = B_0 + B_1 x_1 + B_2 x_2$$

LPA cgaq Iq

To find: B_0, B_1, B_2

step ① start random

$$B_0 = 0, , B_1 = B_2 = 1$$

step ② decide

$$\text{epoch} = 100$$

$$\text{learning rate} = 0.1$$

$$B_0 = B_0 - \eta \text{slope}$$

$$B_1 = B_1 - \eta \text{slope}$$

$$B_2 = B_2 - \eta \text{slope}$$

$$L(B_0, B_1, B_2) \rightarrow \text{4D graph}$$

3 components there

In direction of B_0, B_1, B_2 .

$$\frac{\partial L}{\partial B_0}, \quad \frac{\partial L}{\partial B_1}, \quad \frac{\partial L}{\partial B_2}$$

n columns $\xrightarrow{n+1}$ derivatives $\xrightarrow{\text{update}}$ in equations

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

↓

$$= \frac{1}{2} \left[(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 \right]$$

$$\hat{y} = B_0 + B_1 x_1 + B_2 x_2$$

$$\hat{y}_1 = B_0 + B_1 x_{11} + B_2 x_{12}$$

$$\hat{y}_2 = B_0 + B_1 x_{21} + B_2 x_{22}$$

$$L = \frac{1}{2} \left[(y_1 - B_0 - B_1 x_{11} - B_2 x_{12})^2 + (y_2 - B_0 - B_1 x_{21} - B_2 x_{22})^2 \right]$$

$$\begin{aligned} \frac{\partial L}{\partial B_0} &= \frac{1}{2} \times \left[2(y_1 - B_0 - B_1 x_{11} - B_2 x_{12}) (0 - 1 - 0 - 0) \right. \\ &\quad \left. + 2(y_2 - B_0 - B_1 x_{21} - B_2 x_{22}) \right] \\ &= \frac{-2}{2} \left[(y_1 - \hat{y}_1) + (y_2 - \hat{y}_2) \right] \end{aligned}$$

= But But But
for n ,

$$= \frac{-2}{n} ((y_1 - \hat{y}_1) + (y_2 - \hat{y}_2) + (y_3 - \hat{y}_3))$$

$$\frac{\partial L}{\partial B_0} = \frac{-2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \quad \underline{\quad \dots (y_n - \hat{y}_n) \quad} \quad (1)$$

$$\begin{aligned} \frac{\partial L}{\partial B_1} &= \frac{1}{2} [2(y_1 - \hat{y}_1)(-x_{11}) + \\ &\quad 2(y_2 - \hat{y}_2)(-x_{21})] \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial B_1} &= \frac{-2}{n} [(y_1 - \hat{y}_1)(x_{11}) + (y_2 - \hat{y}_2)(x_{21})] \quad (29) \\ &\quad + (y_3 - \hat{y}_3)(x_{31}) + \dots + (y_n - \hat{y}_n)(x_{n1}) \end{aligned}$$

$$\frac{\partial L}{\partial \beta_1} = \frac{-2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{i1}$$

what x_{i1} represents?
first column all.

x_{11}	x_{12}	x_{13}
x_{21}	x_{22}	
x_{31}	x_{32}	
x_{41}	x_{42}	

$$\frac{\partial L}{\partial \beta_2} = \frac{-2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{i2}$$

(3)

what " " x_{i2} "
all values of second column.

$$\boxed{\frac{\partial L}{\partial \beta_m} = \frac{-2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{im}}$$

m = no of columns
 n = no of rows

Stochastic gradient descent

problem in batch gradient descent

$$\frac{\partial L}{\partial \beta_n} = \frac{-2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{in}$$

we are calculating n derivatives for each and every coeff + 1.

wt $n = 1000$ (no of rows)
 $col = 5$ (no of columns)
coeff. = $5+1 = 6$

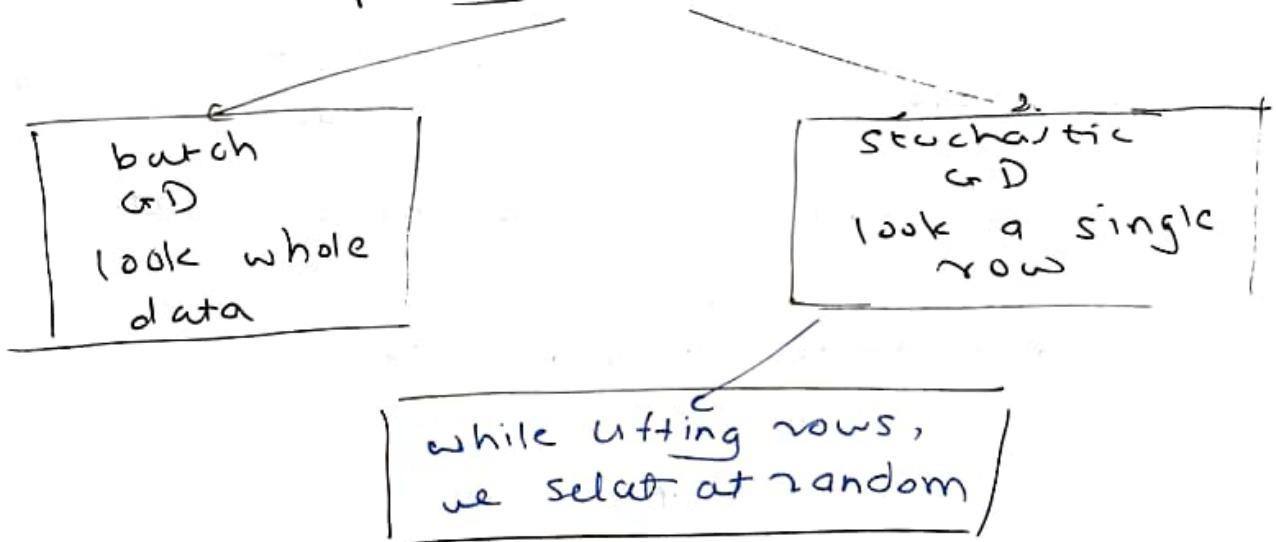
epoch = 50

$$\begin{aligned} & 50 \\ & \left| \begin{array}{l} 1 \rightarrow 1000 \text{ derv} \\ 6 \rightarrow 6000 \text{ derv} \end{array} \right. \\ & = 50 \times 6,000 \\ & = 3,000,000 \text{ derv} \end{aligned}$$

but what if it is a large number?

A lot of computation

Hardware limitation, while calculating \hat{y} , we ~~can't~~ need to load x -train entirely in memory. Can't deal in large dataset. updating coefs and intercept

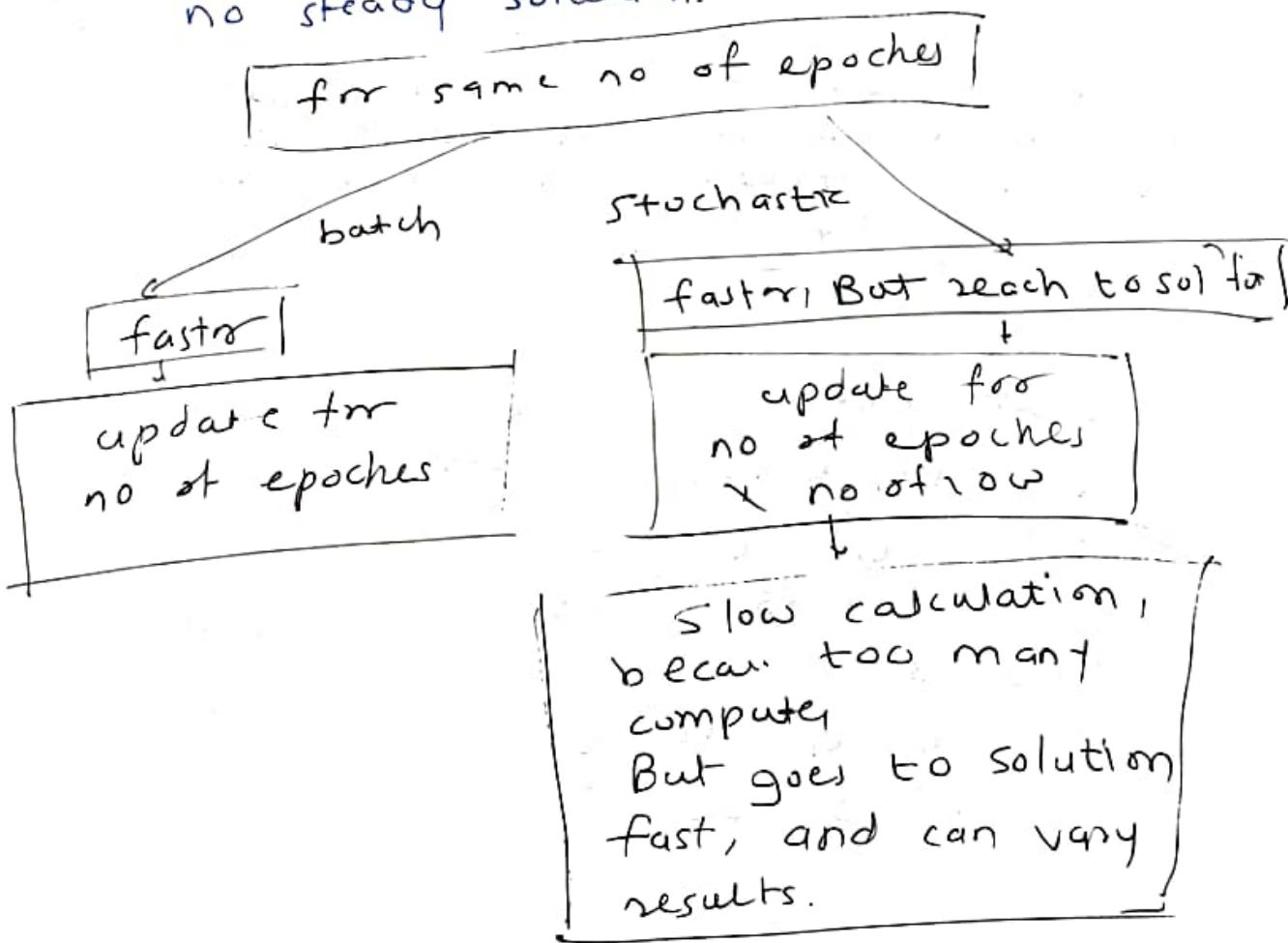


advantage

faster rate of convergence

disadvantage

no steady solution.



when to use stochastic?

→ Big data

② non convex function

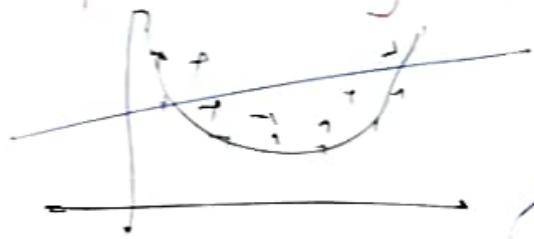
Disadvantages

even after reaching near to actual sol,
we even struggle & fluctuate there
→ solve this,
we can change learning rate later.

mini-batch gradient descent

In between batch and stochastic

Polynomial regression



for a linear regression,
I have eqn $y = mx + c$

$$y = B_0 + B_1 x$$

$$y = B_0 + B_1 x_1 + \dots + B_n x_n$$

multiple linear regression
simple linear regression

If I wish to bring polynomial, what
I will do is I will convert it to
polynomial features in preprocessing.

x	y	(1)
1	1	
2	4	
3	9	
4	16	
5	25	

at degree 2

\downarrow

$1, 35, 1225$

(35, 100)

$1 \quad 35 \quad 1225$

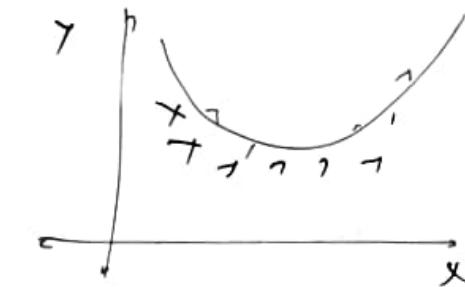
degree
Hyperparameter

for deg 3

x^0	x^1	x^2	x^3
1	35	1225	1500625

$$y_n = B_0 + B_1 x_1 + B_2 x_2 + \dots + B_n x_n$$

Effects of choosing degree
for γ too small,



will not cover
data points properly
and set attribute

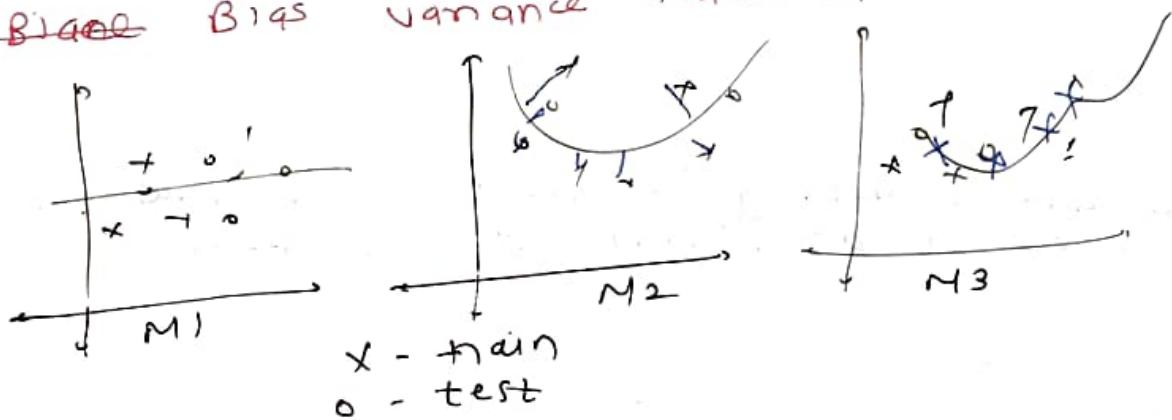


overfitting - need to
find optimal value

Let's say, for 2 cols

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \beta_4 x_2^2$$

Bias vs variance trade-off



Bias: The inability of a ML model to capture relationship between training data

M3 - Low bias

M1 - High bias

Variance
Left her train data error & test data error.

so M1 & M2 - low variance
as they cover train & test both as well. so, LV.

In M3, test error will very high &
train error low.

so High variance

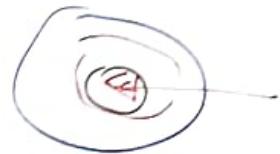
overfitting and underfitting

overfitting
 if ML model works well on train, but
 not on test.
 train error <<< test error
 model M3

underfitting
it model unable to understand data
properly.

Model M1

low bias
+
low variance



To methods to deal with this regularization

bagging

boosting

Ridge Regression

L2 regularization

L_2 regularization
we add something in L_2 model so that our probability of overfitting get decreased.

Techniques available:

→ Ridge regression (regularization)

2) Lasso regression

③ Elastic Net regression

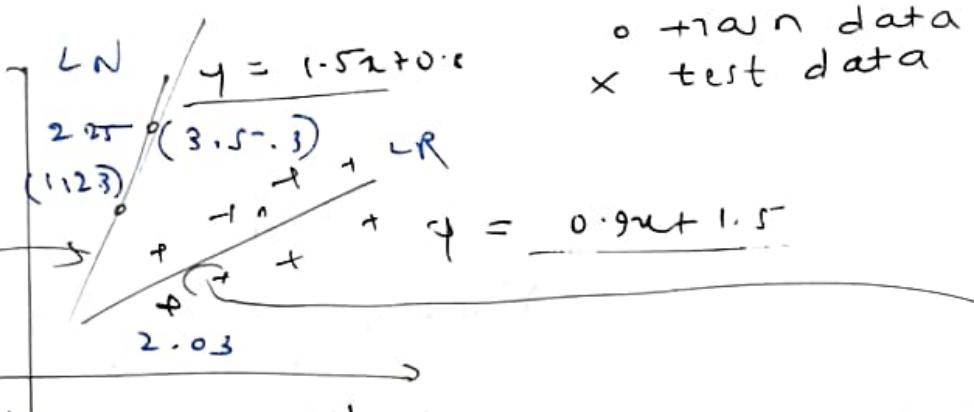
Overfitting in second term,

$$\therefore y = m x + c$$

extreme high value

extreme high value
high dependence of model

so see the next page



I need to convince my model to go from here to here.
overfitting

To minimize error,

I do like thus

$$LF = \sum_{i=1}^n (y_i - \hat{y})^2 + \lambda(m^2)$$

I add this term.

λ = hyper parameter ($0 - \infty$)

m = slope

But how I will convince my ML model that?

We are assuming test data best fit line here

I will calculate errors

Loss LN	Loss LR
$\lambda = 1$	$\lambda = 1$
$0 + (1.5)^2 =$ 2.25	$(2 \cdot 3 - 0.9 - 1.5)^2$ $+(5 \cdot 3 - 2 \cdot 7 - 1.5)^2$ $+(0.9)^2$ $= (0.1)^2 + (1.1)^2 + (0.9)^2$ $= 2.03$
	(L2 regularization)

what + we have made up?
we will add $\lambda(m_1^2 + m_2^2 \dots)$

mathematics of scratch

see LR overview first, and continue

$$L = \sum_{i=1}^n (y - \hat{y}_i)^2 + \lambda m^2$$

$$\text{add } L = \sum_{i=1}^n (y_i - m x_i - b) \quad \frac{\partial L}{\partial m} \quad \frac{\partial L}{\partial b}$$

$$\frac{\partial L}{\partial b} = 2 \sum_{i=1}^n (y_i - m x_i - b) (0 + 0 - 1)$$

$$\sum_{i=1}^n (y_i - m x_i - b) = 0$$

$$\frac{\sum_{i=1}^n (y_i)}{n} - m \frac{\sum_{i=1}^n (x_i)}{n} = \frac{nb}{n}$$

$$b = \bar{y} - m \bar{x} \quad \text{--- (1)}$$

$$L = \sum_{i=1}^n (y_i - m x_i - \bar{y} + m \bar{x})^2 + \lambda m^2$$

$$= 2 \sum_{i=1}^n (y_i - m x_i - \bar{y} + m \bar{x}) (0 - x_i - 0 + \bar{x}) + 2 \lambda m$$

$$= -2 \left(\sum_{i=1}^n (y_i - \bar{y} - m x_i + m \bar{x}) (x_i - \bar{x}) \right) + 2 \lambda m = 0$$

$$= \sum_{i=1}^n (y_i - \bar{y} - m x_i + m \bar{x}) (x_i - \bar{x}) = 2 \lambda m$$

$$\sum_{i=1}^n ((y_i - \bar{y}) - m(x_i - \bar{x})) (x_i - \bar{x}) = 2 \lambda m$$

$$2m - \sum_{i=1}^n (y_i - \bar{y}) (x_i - \bar{x}) - m(x_i - \bar{x})^2 = 0$$

$$\begin{aligned}
 &= \hat{y}_m - \sum_{i=1}^n (y_i - \hat{y}) (x_i - \bar{x}) + m \cdot \sum_{i=1}^n (x_i - \bar{x})^2 \\
 &= \hat{y}_m + m \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (y_i - \hat{y}) (x_i - \bar{x}) \\
 m &= \frac{\sum_{i=1}^n (y_i - \hat{y}) (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2 + \lambda} \quad \text{--- (2)} \\
 &\qquad \qquad \qquad \text{new one added compared to original linear regression}
 \end{aligned}$$

Let us go in n dimensions

$$x_1 \ x_2 \ \dots \ x_n \mid y$$

(n+1) cols, m rows

$$L = \sum_{i=1}^m (y_i - \hat{y})^2$$

$$m \left| \begin{array}{cccc} x_1 & x_2 & \dots & x_n \\ w_1 & w_2 & \dots & w_n \end{array} \right| \text{(coef)}$$

and no intercept

$$= (xw - r)^T (xw - r)$$

$$r = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ y \end{bmatrix} \quad w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \quad x = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}$$

$$L = (xw - r)^T (xw - r) \text{ for normal.}$$

Now, let add λ . term (regularization)

$$L = (xw - r)^T (xw - r) + \lambda \|w\|^2$$

$$\left\{ \lambda w_1^2 + \lambda w_2^2 + \dots + \lambda w_n^2 \right.$$

$$\left. \lambda (w_0^2 + w_1^2 + \dots + w_n^2) \right\}$$

$$[w_0 \ w_1 \ w_2 \dots \ w_n] \left[\begin{array}{c} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{array} \right]$$

(37)

$$L = (xw - y)^T (xw - y) + \lambda w^T w$$

$$L = [(xw)^T - y^T] [xw - y] + \lambda w^T w$$

$$= [w^T x^T - y^T] [xw - y] + \lambda w^T w$$

$$= (w^T x^T x w - w^T x^T y - y^T x^T w + y^T y) + \lambda w^T w$$

$a^T - b^T = (a-b)^T$
 $a^T + b^T = b^T a^T$

$$= \underline{w^T x^T x w} - \underline{2 w^T x^T y} + \underline{y^T y} + \underline{\lambda w^T w}$$

$$\frac{dL}{dw} = \underline{2 x^T x w} - \underline{2 x^T y} + \underline{0} + \underline{2 \lambda w} = 0$$

$$= x^T x + \lambda w + w^T x^T y$$

~~$x^T x + \lambda w$~~ \Rightarrow

$$(x^T x + \lambda I) w = x^T y$$

$$w = \underline{(x^T x + \lambda I)^{-1} x^T y}$$

\Leftarrow in LR form

$$w = (x^T x)^{-1} x^T y$$

by gradient descent

vector form

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

$$L = (xw - y)^T (xw - y) + \lambda \|w\|^2$$

$$L = (xw^T - y)^T (xw - y) + \lambda w^T w$$

we know that we have $w_0, w_1 \dots w_n$

$$w_0 = w_0 + \eta \frac{\partial L}{\partial w_0},$$

$$w_1 = w_1 + \eta \frac{\partial L}{\partial w_1}$$

$$\vdots$$

$$w_n = w_n + \eta \frac{\partial L}{\partial w_n}$$

$$w_{new} = w_{old} - \eta \begin{bmatrix} \Delta L \\ \Delta w \end{bmatrix}$$

Gradient

$$\begin{aligned} L &= (x_w - y)^T (x_w - y) + \lambda w^T w \\ &= \frac{1}{2} (w^T x - y^T) (x_w - y) + \frac{1}{2} \lambda w^T w \\ &= \frac{1}{2} \left[w^T x^T x_w - \cancel{w^T x^T y} - y^T x_w + y^T y \right] \\ &\quad + \frac{1}{2} \lambda w^T w \\ &= \frac{1}{2} \left[w^T x^T x_w - \cancel{\frac{2}{2} w^T x^T y} + y^T y \right] \\ &\quad + \frac{1}{2} \lambda w^T w \end{aligned}$$

$$\frac{\partial L}{\partial w} = \frac{1}{2} \left[2' x^T x_w - \cancel{\frac{2}{2} x^T y} \right] + \frac{1}{2} \cancel{\lambda} \cdot 2(x^T w)$$

$$\frac{\partial L}{\partial w} = x^T x_w - \cancel{x^T y} + w = \boxed{x^T x_w - x^T y - \lambda w}$$

$$w = [0 \ 1 \ 1 \ 1]$$

start point

in epoch

$$w = w - \eta \frac{\partial L}{\partial w}$$

effects of ridge regression.

- 1) effect on coefficient
on higher value of alpha, coefficients goes closer to 0 but not exact 0.
- 2) higher values are impacted more at say,

Linear regression

$$w = w_0 + w_1 x_1 + w_2 x_2$$

| | |
high med low

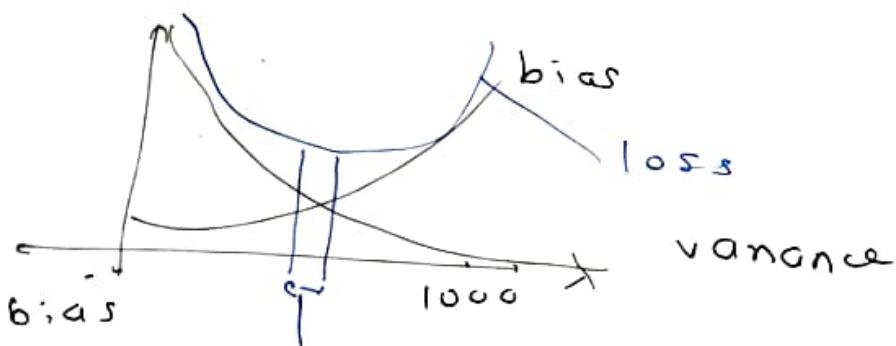
if we take ridge & go on inc value of α , what will impact?

Higher values of coefficients on higher values goes closer to 0.

- 3) Bias - variance tradeoff

-f

		bias	model	variance
x	less ≈ 0	↓	overfit	↑
x	greater	↑	underf	↓



take this value

- 4) Impact of loss function

$$L = \sum_{i=1}^n (y_i - \hat{y})^2 + \lambda \|w\|^2$$

y is min \hat{y} $\rightarrow b$ is constant

$$L = \sum_{i=1}^n (\hat{y} - \mu)^2 + \lambda \|w\|^2$$



Lasso regression

L1 regularization

In ridge regression

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|w\|_2^2 \quad (\text{L2}) \quad (\text{2nd norm square})$$

In Lasso regression

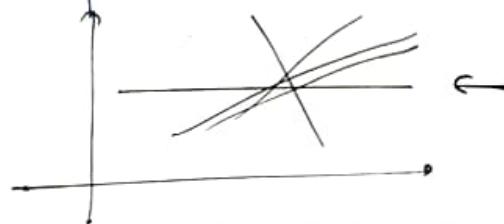
$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \frac{\lambda \|w\|_1}{P} \quad (\text{L1}) \quad (\text{first norm mod})$$

absolute

In case if we increase λ a lot, can underfit.
In case if we decrease λ a lot, can overfit.

$$\lambda > 0$$

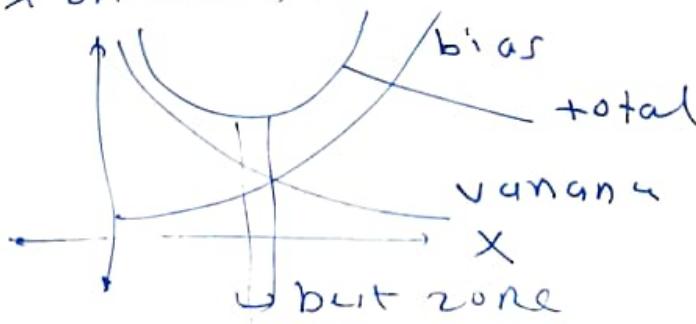
If $\lambda = 0$, then it becomes linear regression
The first point to be noted is,
unlike ridge regression, we can obtain
coefficients equals to 0.



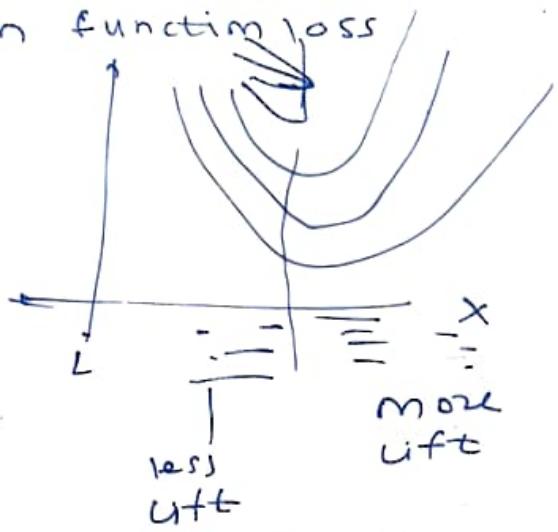
so, In higher dimensions, There is a lot a case of overfitting. If we apply ridge, we would have to consider columns also who are not much important, due to their closer significance to 0. But if we apply lasso, the coefficient of less important features become 0.

key points

- 1) Feature selection
- 2) higher coef's impact more
- 3) λ on bias & varn



a) on function loss



At high values of λ , it becomes 0 but will not go behind 0.

Why sparsity in lasso?

I mean why some values become 0?
why lasso can do this & why don't ridge?

Simple example,

$$x, y \rightarrow y = mx + b, b = \bar{y} - m\bar{x} \quad (\text{sg one})$$

$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Simple linear =

$$\text{ridge} \quad m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2 + \lambda}$$

Let us derive for lasso.

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda|m|$$

$$L = \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2 + \lambda|m|$$

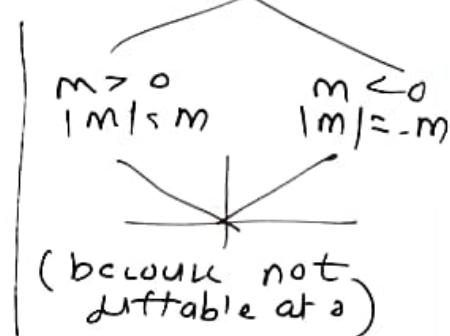
for, $|m| = m$ (positive)

$$L = \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2 + 2\lambda|m|$$

$$\frac{\partial L}{\partial m} = 2 \cdot \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x}) (-x_i + \bar{x}) + 2\lambda$$

$$= -2 \left[(y_i - \bar{y}) - \frac{m}{n} (x_i - \bar{x}) \right] (x_i - \bar{x}) + 2\lambda = 0$$

$$= - \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + m \sum_{i=1}^n (x_i - \bar{x})^2 + \lambda = 0$$



$$3. \frac{\sum_{i=1}^n (x_i - \bar{x})^2 = \sum (x_i - \bar{x})(x_i - \bar{x}) - \lambda}{\begin{cases} \lambda = \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x}) - \lambda}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{cases}} \quad \text{--- } \textcircled{1}$$

for $\lambda = 0$,

$$\begin{cases} \lambda = \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{cases} \quad (\text{linear}) \quad \text{--- } \textcircled{2}$$

for $\lambda < 0$,

$$L = \sum_{i=1}^n (x_i - \bar{x})^2 - \lambda m$$

$$L = \sum_{i=1}^n (x_i - mx_i - \bar{x} + m\bar{x})^2 - 2\lambda m$$

$$\frac{\partial L}{\partial m} = 2 \sum_{i=1}^n (x_i - mx_i - \bar{x} + m\bar{x})(x_i - \bar{x}) - 2\lambda$$

$$\frac{\partial L}{\partial \lambda} = 2 \sum_{i=1}^n (x_i - mx_i - \bar{x} + m\bar{x})(x_i - \bar{x}) - 2\lambda = 0$$

$$0 = -2 \sum_{i=1}^n [(x_i - \bar{x}) - m(x_i - \bar{x})] (x_i - \bar{x}) - 2\lambda = 0$$

$$\begin{cases} = -2 \sum_{i=1}^n \left[(x_i - \bar{x})(x_i - \bar{x}) - m(x_i - \bar{x})^2 \right] - 2\lambda = 0 \\ = m \sum_{i=1}^n (x_i - \bar{x})^2 - \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x}) + \lambda \end{cases} \quad \text{--- } \textcircled{3}$$

According to formula,

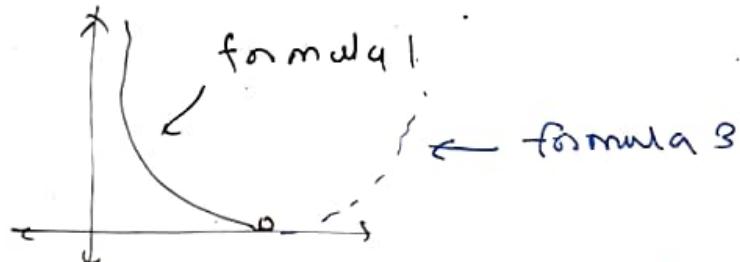
if $\frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x}) - \lambda}{\lambda}$ travels until, $\textcircled{1}$ $\text{--- } \text{This and }$ this become equal at negative we add λ . $\text{--- } \textcircled{3}$

so, If we decrease on and on, (Increase) λ , then algorithm stop there, due to bad case situation.

means. In case of positive slope

$$\text{if } \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + \lambda}{\sum_{i=1}^n (x_i - \bar{x})^2} = m$$

if we increase λ , numerator start decreasing
see the black line.



as soon as for the terms of,

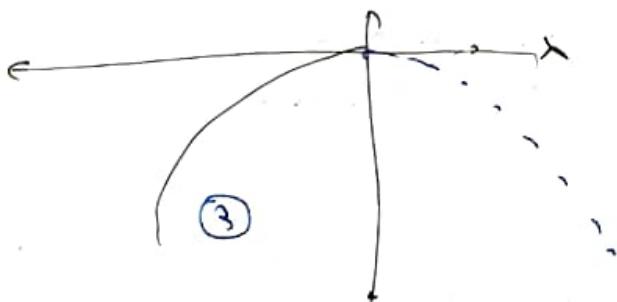
$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \leq \lambda,$$

The formula changes & becomes ③.
That's why, It don't goes down at 0,
and start increasing again

for negative slope

$$\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + \lambda}{\sum_{i=1}^n (x_i - \bar{x})^2} = m$$

if we decrease increase λ and then : if
I decrease λ , It goes towards 0.



as soon as for m + term of

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \geq \lambda,$$

the slope comes +ve & becomes ③

that's why, It goes up at 0. so algo.
doubt
stops it.

so why lasso create sparsity? why equal to 0?

λ comes in numerator in lasso
 λ comes in denominator in ridge
 so it is able to ride to 0.
 so it is not able to ride to 0.
 so it closes to 0 but don't get 0.

Elastic Net regression

If we don't know, in a large dataset, for a column in it, its importance to predict. In short, we confuse to choose between ridge & lasso.

we use elasticnet regression.

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \alpha \|w\|^2 + \beta \|w\|$$

In sklearn,
 we get hyperparameters,

$$\lambda = \alpha + \beta$$

$$\text{L1 ratio} = \frac{\beta}{\alpha + \beta}$$

we decide α and β from here,
 so usually, $\lambda \leq 1$. $\text{L1 ratio } 0-5$ by def.
 $\therefore \alpha = 0.5, \beta = 0.5$.

means 50% ridge & 50% lasso.

if L1 ratio is ~~go~~ 0.9
 means 90% ridge & 10% lasso

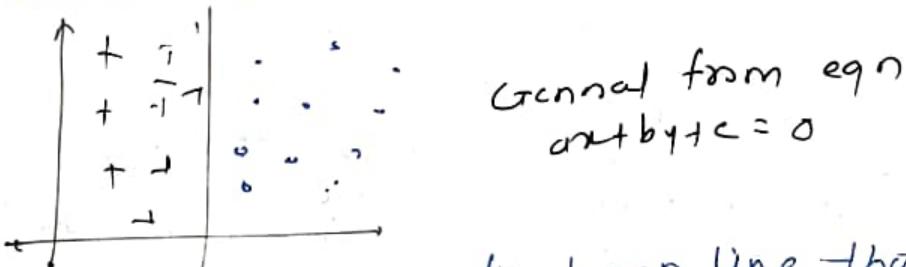
If our dataset are multi-collinear columns
 at at,

2 input cols depends on each other.
 we should use elasticNet regression.

Logistic regression

two approaches
probabilistic
geometric

scenario of requirements
we should have only separable data
or almost linearly separable data

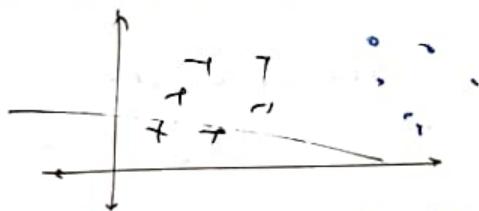


Now we just need to find eqn line that can classify both data in parts.

Perception Trick

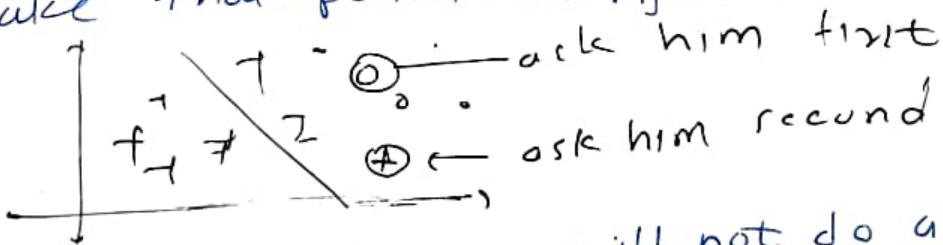
This algorithm does like this.

It starts with first random a, b & c .



Then it in a loop select a random point, and asks it if it is correctly directed on the line, towards the line. If points say no, then it will do nothing.

But if it is not, it has to move such that I can take that point on right side.



In first ask, you will not do anything. In second ask, you will come just aside of that point.

How much time to travel for loops.

→ either a fixed no iteration

→ until a converge.

How this transformation a,b,c does
to move line?

I think at all, How I will classify & is how
I know side of line. We say +ve & -ve



so, How I will know that point direction,
ie it "+ve" or "-ve"?
well, I know line eqn, what it

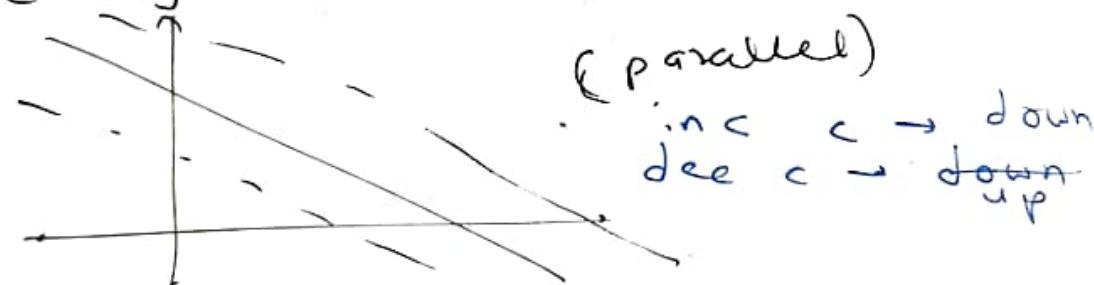
out but $c = 0$
and we say, I have a point (10, 20)
so I will put

$a(10) + b(20) + c = 0$
if it is > 0 , +ve region.

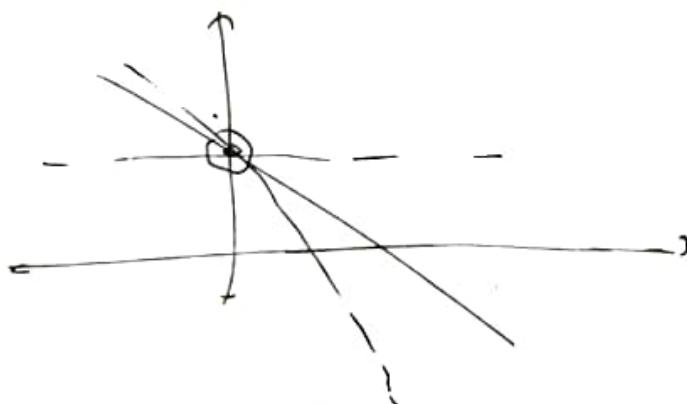
if it is < 0 , -ve region

if it is 0, on the line

The transformation changes comes from
changing a,b & c.
If we change c. It move like this

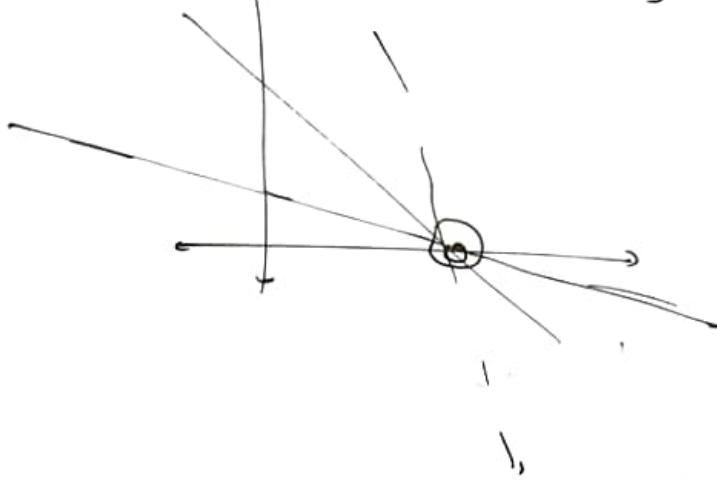


if I change a, I rotate on axis



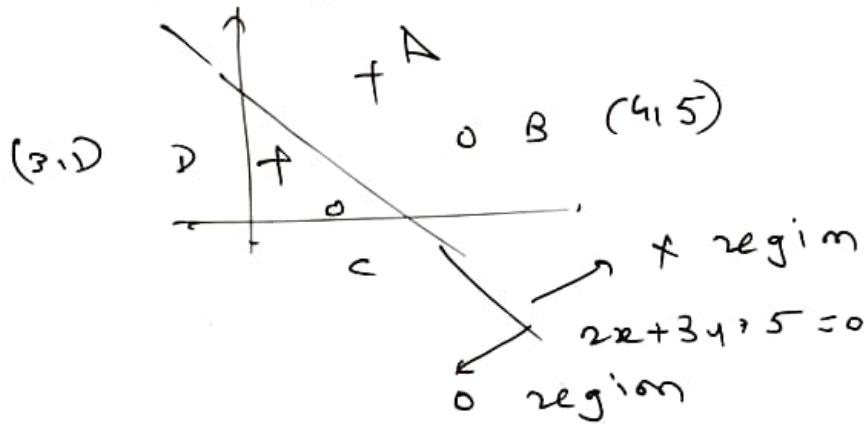
inc $c \rightarrow$ down
dec $c \rightarrow$ down up

similarly for changing b ,



Now we know how transformation are getting applied.

But how will I know how to & what to apply?



point A \rightarrow correct else incorrect

point C \rightarrow correct
I will first ask point B.

So I is not. So I wish to move

these what pupkin trick is,

so, I have line

$$(2, 3, 5)$$

I will consider 1 at end of point.

$$(4, 5, 1)$$

if a negative pt subtract to move towards the dire
it is in pos regm $(-2, -2, 4)$

this, $-2x - 2y + 4 = 0$, will move

point B in negative direction.

or say, point D,

$$(2, 3, 5)$$

at our pos pt it is in neg regm odd $(3, 3, 1)$
 $3x + 3y + 1 = 0$ move in negative.

(49)

But in ML,
we can't do this kind of large transforms.
So, To control, we use learning rate (η)
we multiply each coef in point & then
add pow- to coeffs of line / subtract.

i.e. $\begin{matrix} -2 & 3 & 5 \\ 0.02 & 0.02 & 0.01 \end{matrix}$

like this.

$$\text{new coef} = \text{old coef} - \eta \times (\text{err})$$

so change happens gradually.

Algorithm:

x_0	$c g p q$	I_q	placed
1	7.5	61	1
1	5.9	109	1
1	7.0	61	0

$$ax + bu + c = 0$$



$$w_0 + w_1 x_1 + w_2 x_2$$

$$\left. \begin{array}{l} w_0 = c, \quad w_1 = A, \quad w_2 = B \\ w_0 x_0 + w_1 x_1 + w_2 x_2 = 0 \\ (\because x_0 = 1) \end{array} \right\}$$

$$n = \sum_{i=0}^2 w_i x_i = 0$$

say this is my model.

But how I will predict?

so we put

$$w_0(1) + w_1(7.5) + w_2(5.9) = 0$$

so, if this number n ,

if $n \geq 0$ predict 1.

& $n < 0$ predict 0

$$[w_0 \quad w_1 \quad w_2] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = n$$

$$\begin{array}{ll} n \geq 0 & \rightarrow 1 \\ n < 0 & \rightarrow 0 \end{array}$$

algorithm

epoch define

define learning rate

for i in range(epoch)

random select a point.

check if $x_i \cdot (\text{row}) \in \text{negative region}$

$\sum_{i=1}^n w_i x_i \geq 0$ (according to mod)

(so, it is a negative point, stuck in +ve region)

$$w_{\text{new}} = w_{\text{old}} - \eta [x_0 + x_i]$$

check if $x_i \cdot (\text{row}) \in \text{positive region}$

$\sum_{i=1}^n w_i x_i < 0$ (according to mod)

(so, it is a positive point, stuck in -ve region)

$$w_{\text{new}} = w_{\text{old}} + \eta (x_i)$$

means, simply

$$\begin{cases} w_i \in N & w_i < 0 \\ x_i \in P & w_i \geq 0 \end{cases} \quad \left. \begin{array}{l} \checkmark \\ \checkmark \end{array} \right\} \text{change}$$

simplifying algorithm!

only one cond'n instead of two

other than this,

$w_i \in N$ and $\sum_{i=1}^n w_i x_i > 0$

$$w_n = w_{\text{old}} - \eta x_i$$

$x_i \in P$ and $\sum_{i=1}^n w_i x_i < 0$

$$w_n = w_{\text{old}} + \eta x_i$$

for i in epoch

select random row

$$w_n = w_0 + \eta (y_i - \hat{y}_i) x_i$$

(update rule directly)

if model & I am same

this will 0.

$$w_n = w_0 + 0. \quad (\text{same})$$

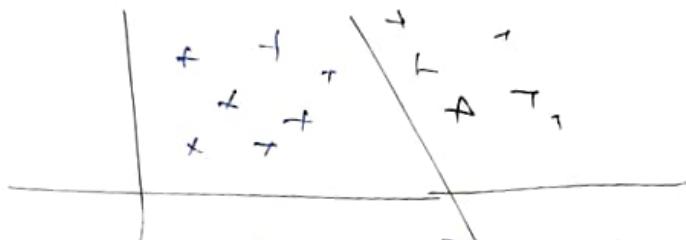
if model & I missmatch
this will $\neq 0$

$$w_n = w_0 \pm (n)(x_i - \hat{y}_i)$$

(change occurs) (51)

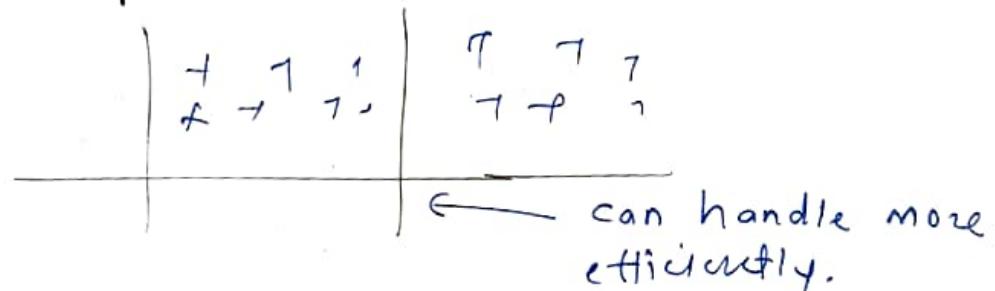
act-pred	$y_i - \hat{y}_i$	mark
0	0	✓
1	0	✓
1	-1	✗
1	1	✗

Issue with LWS



can't handle test data more efficiently → line by algorithm

It stops as soon as it classify all points. But it is not more practical.
more practical should be



upgraded algorithm

until now, missclassified points used to attract (pull) towards it.

Now we will change the approach

The correctly classified points will send along the line. (push) from it. and also missclassified point - each and every point miss - pull
correct - push

The magnitude of pushing / pulling line will depend on the distance.

point	action	magnitude	distance
missclass	pull	less	near
missclass	pull	high	long
correctclass	push	less	near
correctclass	push	high	long

so this is the case that will help us.

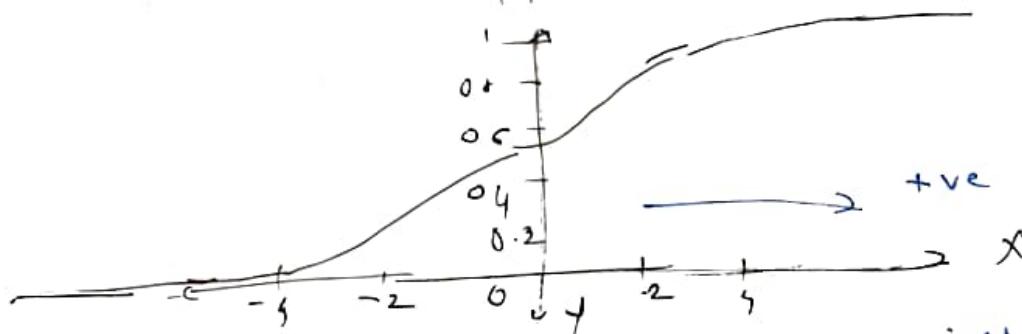
$$w_n = w_0 + \eta (y_i - \hat{y}_i) x_i$$

y_i	\hat{y}_i	$y_i - \hat{y}_i$
1	1	0
0	0	0
1	0	1
0	1	-1

if we wish to change in both cases avoid to create $(y_i - \hat{y}_i) = 0$
So mean $(y_i - \hat{y}_i) \neq 0$.
we can change only y_i , can't \hat{y}_i .
so one eqn at 2-D can be
 $w_0 + w_1 x_1 + w_2 x_2 = 0$
I have to change step function
so we have to understand sigmoid

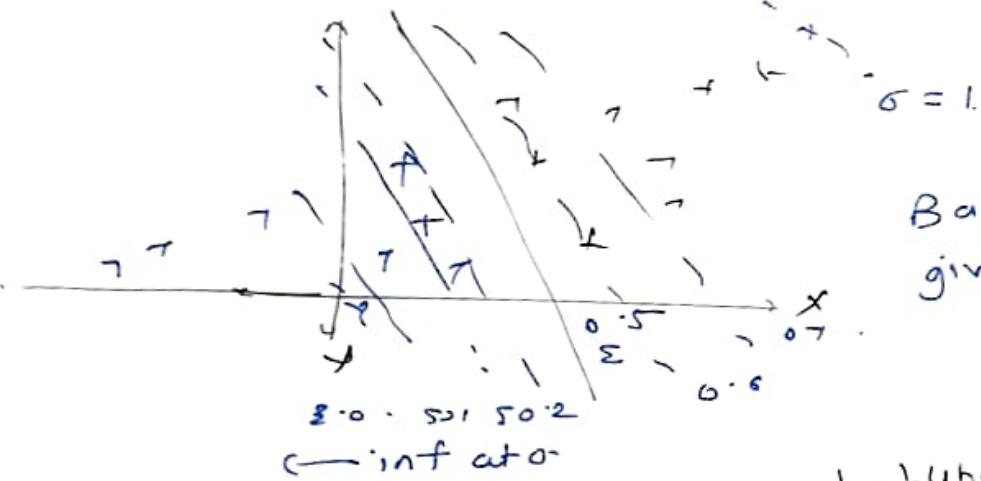
Sigmoid function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



so we were calculating $\frac{\sum w_i x_i}{-1}$
this if positive no $\geq 0.5 \rightarrow 1$
if negative no $< 0.5 \rightarrow 0$

But what is disadvantage of that?



Basically it give a gradient

also sigmoid states probability also if σ is 0.5, It means 50% chance probability.

so, sigmoid gives a probabilistic interpretation.

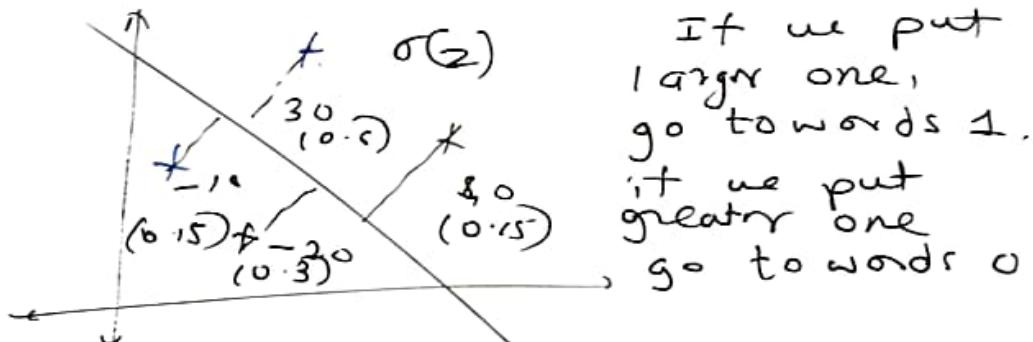
Derivation

effect of sigmoid function

$$w_n = w_0 + \gamma (y_i - \hat{y}_i) x_i$$

$$\hat{y}_i = \sigma(z_i) \quad (\text{sigmoid})$$

$$\gamma = \sum w_i x_i$$



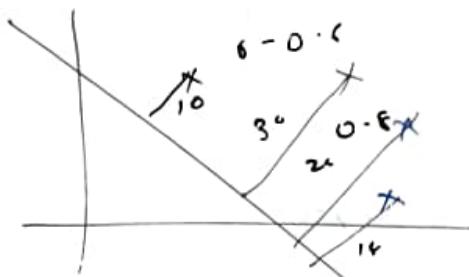
If we put larger one, go towards 1. If we put greater one, go towards 0

y_i	\hat{y}_i	$w_i x_i$
1	0.8	0.2
0	0.15	-0.65
1	0.3	0.7
0	0.15	-0.15

$$x_{n+1} = x_n + \gamma w_i x_i \quad \text{push down}$$

$$x_n = x_0 - \gamma w_i x_i \quad \text{lift up}$$

magnitude of pushing and pulling



t_i	\hat{w}	w_{i-j}
1	0.4	0.4
1	0.2	0.2

$$w_n = w_0 + n \cdot 0.4xi \rightarrow w_1$$

$$w_n = w_0 + n \cdot 0.2xi \rightarrow w_2$$

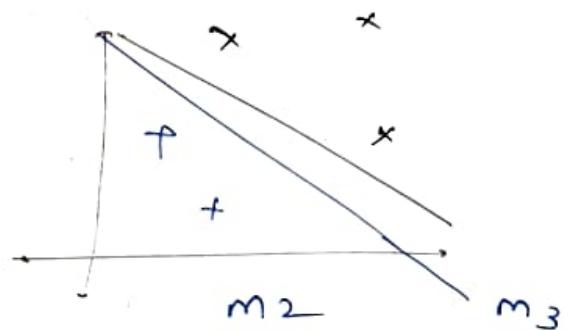
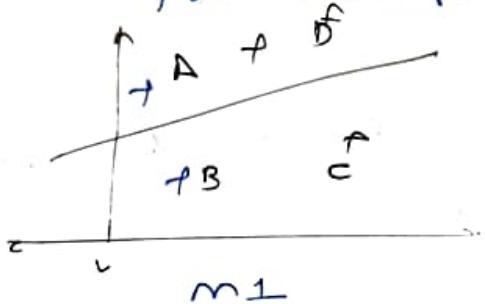
$$w_1 > w_2$$

In this approach As we are towards solution, and But not to a optimal soln and secondly we are generating points randomly. So we cannot tell if we are right there on optimal solution.

In machine learning, things doesn't work this way. we always need a loss function. we take minima, where loss function is minimum & that coef's are our actual solution.

find a loss function

for example



for now, M_2 is better.

But now, I can't tell if M_2 and M_3 , which one better.

To calculate that,

we need a loss function.

So, we will use maximum likelihood function these related to probability.

so, for two modus,
let us calculate \hat{z} first.

$$\hat{z} = \sigma(z)$$

so, we will calculate green prob and
black prob

points	$P(\text{blue})$	$P(\text{black})$
A	0.4	0.6
B	0.2	0.8
C	0.9	0.1
D	0.7	0.3

points	$P(\text{blue})$	$P(\text{red})$
A	0.7	0.3
B	0.4	0.6
C	0.9	0.1
D	0.3	0.7

maximum likelihood says to multiply all the probabilities. One of the better prob will be better modus.

$$m_1 = P(\text{blue}) * P(\text{blue}) * P(\text{black}) * P(\text{blue})$$

$$= 0.7 * 0.4 * 0.9 * 0.7 = 0.089$$

$$m_2 = \underbrace{0.7 * 0.6 * 0.1 * 0.3}_{\text{not}} = 0.176$$

$$m_2 > m_1$$

But this way good; problematic for large datasets.

so we will take log to handle it.

$$m_1 = \log(ab) = \log(a) + \log(b)$$

But still a problem.

log is always between 0 & 1 & if I calculate log between (0-1), then it is always negative

To tackle it, I will take -

$$= -\log(a) - \log(b)$$

It is called cross entropy.

When I do summation of all negative terms of logs of maximum likelihood,
It is called cross entropy.

Like we have to maximise ~~cross entropy~~ likelihood,
we need to minimize cross entropy

We need to consider \hat{y}_i s as that blue/black points probability.

For a moment, can y for 4 points, can we do thus?

$$-\log(\hat{y}_1) - \log(\hat{y}_2) - \log(\hat{y}_3) - \log(\hat{y}_4)$$

No.

Because, we need to see probability of classification, of that original point.

Sol We convert it.

$$\boxed{L = -\sum_i y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)}$$

= for point 1, target is 1. (y_1)

$$-y_1 \log(0.7) - (0)$$

In 2nd point, target 0

$$-0 - \log(1-0.7)$$

for 3rd point target 1

$$-0 - \log(0.9)$$

for 4th point, target 0

$$-0 - \log(0.2)$$

$$-0 - \log\left(\frac{\log(0.7)}{0.2}\right)$$

Loss function =

$$\boxed{L = \sum_{i=1}^4 -y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)}$$

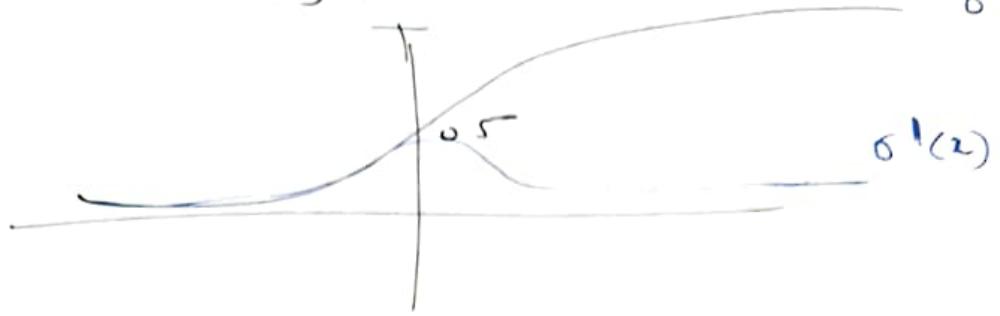
no
closed
form
solution

\log - loss error

binary - cross entropy fun?

(7)

Derivative of sigmoid function



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

$$\begin{aligned}\frac{d}{dx} \left(\frac{1}{1 + e^{-x}} \right) &= \frac{d}{dx} (1 + e^{-x})^{-1} \\&= \frac{-1}{(1 + e^{-x})^2} \cdot \frac{d}{dx} (1 + e^{-x}) \\&= \frac{-1}{(1 + e^{-x})^2} \cdot \frac{d}{dx} (e^{-x}) \\&= -\frac{e^{-x}}{(1 + e^{-x})^2} \cdot \frac{d}{dx} (-x) \\&= -\frac{e^{-x}}{(1 + e^{-x})^2} \\&= \frac{1}{(1 + e^{-x})} \cdot \frac{e^{-x}}{(1 + e^{-x})} \\&= \frac{\sigma(x)}{\sigma(x)} \cdot \frac{e^{-x}}{(1 + e^{-x})}\end{aligned}$$

$$\therefore \sigma(x) \left[\frac{1 + e^{-x}}{1 + e^{-x}} - 1 \right]$$

$$= \sigma(x) \left[\frac{1 + e^{-x}}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right]$$

$$\boxed{\sigma'(x) = \sigma(x)(1 - \sigma(x))} \quad \text{--- ①}$$

gradient - Ducci approach

dataset

$$\begin{matrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{matrix} \quad \begin{matrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{matrix}$$

$$[w_0 \ w_1 \ w_2 \ w_3 \ \dots \ w_n]$$

coefficients

for first one,

$$\sigma(w_0 + w_1 x_{11} + w_2 x_{12} + w_3 x_{13} + \dots + w_n x_{1n}) = \hat{y}_1$$

$$\sigma(w_0 + w_1 x_{21} + w_2 x_{22} + w_3 x_{23} + \dots + w_n x_{2n}) = \hat{y}_2$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \sigma(w_0 + w_1 x_{11} + w_2 x_{12} + \dots + w_n x_{1n}) \\ \sigma(w_0 + w_1 x_{21} + w_2 x_{22} + \dots + w_n x_{2n}) \\ \vdots \\ \sigma(w_0 + w_1 x_{m1} + w_2 x_{m2} + \dots + w_n x_{mn}) \end{bmatrix}$$

$$= \sigma \begin{bmatrix} w_0 + w_1 x_{11} + w_2 x_{12} + \dots + w_n x_{1n} \\ w_0 + w_1 x_{21} + w_2 x_{22} + \dots + w_n x_{2n} \\ \vdots \\ w_0 + w_1 x_{m1} + w_2 x_{m2} + \dots + w_n x_{mn} \end{bmatrix}$$

$$= \sigma \left(\underbrace{\begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}}_X \underbrace{\begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}}_w \right)$$

$$\hat{y} = \sigma(xw)$$

from the equation,

$$L = -\frac{1}{m} \cdot \sum_{i=1}^m \left[y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i) \right]$$
$$= \frac{-1}{m} \underbrace{\left[\sum_{i=1}^m y_i \log(\hat{y}_i) \right]}_{\textcircled{1}} + \frac{-1}{m} \underbrace{\left[\sum_{i=1}^m (1-y_i) \log(1-\hat{y}_i) \right]}_{\textcircled{2}}$$

subject to \textcircled{1}

$$S_1 = \sum_{i=1}^m y_i \log(\hat{y}_i)$$

$$= y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + y_3 \log \hat{y}_3 + \dots + y_m \log \hat{y}_m$$

$$= [y_1 \ y_2 \ \dots \ y_m] \begin{bmatrix} \log \hat{y}_1 \\ \log \hat{y}_2 \\ \vdots \\ \log \hat{y}_m \end{bmatrix}$$

$$= [y_1 \ y_2 \ \dots \ y_m] \log \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{pmatrix}$$

$$= \gamma \log \hat{y}$$

$$= \gamma \log(\sigma(x_w)) \quad \textcircled{1}$$

$$S_2 = \sum_{i=1}^m (1-y_i) \log(1-\hat{y}_i)$$

$$= (1-y_1) \log(1-\hat{y}_1) + (1-y_2) \log(1-\hat{y}_2) + \dots + (1-y_m) \log(1-\hat{y}_m)$$

$$= [(1-y_1)(1-y_2) \dots (1-y_m)] \begin{bmatrix} \log(1-\hat{y}_1) \\ \log(1-\hat{y}_2) \\ \vdots \\ \log(1-\hat{y}_m) \end{bmatrix}$$

$$= \left[(1-y_1) \quad (1-y_2) \quad \dots \quad (1-y_m) \right] \left(\log \begin{bmatrix} 1-\hat{y}_1 \\ 1-\hat{y}_2 \\ \vdots \\ 1-\hat{y}_m \end{bmatrix} \right)$$

$$= (1-y) \log(\sigma(1-x_w)) \quad \text{--- (2)}$$

from s_1 and s_2

$$L = -\frac{1}{m} \left[Y \log \hat{Y} + (1-Y) \log(1-\hat{Y}) \right]$$

where $\hat{Y} = \sigma(x_w)$
loss in matrix form

gradient-descent approach

define w and epoch

for i in epoch:

$$w_n = w_0 + \gamma \frac{\Delta L}{\Delta w} \quad (\text{gradient})$$

$$\frac{\Delta L}{\Delta w} = \left[\frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_n} \right]$$

$$L = \frac{1}{m} \left[Y \log \hat{Y} + (1-Y) \log(1-\hat{Y}) \right]$$

$$\frac{dL}{dw} = \gamma Y \cdot \frac{d}{dw} \log \hat{Y} = \frac{\gamma}{\hat{Y}} \frac{d(\hat{Y})}{dL}$$

$$= \frac{\gamma}{\hat{Y}} \frac{d}{d\hat{Y}} (\sigma(x_w))$$

$$= \frac{\gamma}{\hat{Y}} \sigma(w_2) [1 - \sigma(w_2)]$$

$$\frac{d}{dw} (w_2)$$

$$= \frac{y}{\hat{y}} \hat{x}(1-\hat{y})x$$

$$= y(1-\hat{y})x \quad \text{---} \textcircled{1}$$

for the second term

$$= (1-y) \log(1-\hat{y})$$

$$\frac{\partial L}{\partial w} = \frac{\partial}{\partial w} [(1-y) \log(1-\hat{y})]$$

$$= (1-y) \frac{\partial}{\partial w} [\log(1-\hat{y})]$$

$$= \frac{(1-y)}{(1-\hat{y})} \frac{\partial}{\partial w} (1-\hat{y})$$

$$= \frac{(1-y)}{(1-\hat{y})} \left[-\frac{\partial \hat{y}}{\partial w} \right]$$

$$= \frac{(1-y)}{(1-\hat{y})} \left[-\frac{\partial}{\partial w} (\sigma(xw)) \right]$$

$$= -\left(\frac{1-y}{1-\hat{y}}\right) \left[\sigma(xw) \cdot (1-\sigma(xw)) \right] \cdot \frac{\partial}{\partial w} [wx]$$

$$= -\frac{(1-\hat{y})\hat{y}}{(1-\hat{y})} (1-\hat{y})x$$

$$= -\hat{y}(1-\hat{y})x \quad \text{---} \textcircled{2}$$

$$\frac{\partial L}{\partial w} = -\frac{1}{m} \left[y(1-\hat{y})x + -\hat{y}(1-y)x \right]$$

$$= -\frac{1}{m} \left[y(1-\hat{y}) - \hat{y}(1-y) \right] x$$

$$= -\frac{1}{m} \left[y - \cancel{y}\hat{y} - \hat{y} + \cancel{y}\hat{y} \right] x$$

$$\boxed{\frac{\partial L}{\partial w} = -\frac{1}{m} [y - \hat{y}] x}$$

so update frequency can be

$$w = w_0 + \eta \frac{1}{m} (y - \hat{y}) x$$

$$\underline{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}, \quad \underline{x} = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,n} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,n} \\ \vdots & x_{m,1} & x_{m,2} & \dots & x_{m,n} \end{bmatrix}_{(m+1) \times (n+1)}$$

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}, \quad \hat{\underline{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix}_{(m) \times 1}$$

$$\underline{w} = \underline{w}_0 + \eta \frac{1}{m} (\underline{y} - \hat{\underline{y}}) \underline{x}$$

$$\underline{\underline{w}}_{+1} = \underline{\underline{w}}_{+1} + \frac{(m,1)}{\downarrow} \underline{\underline{w}}_{m,m+1}$$

$$\underline{\underline{w}}_{m+1,1} = \underline{\underline{w}}_{m+1,1} + \underline{\underline{w}}_{1,m} \underline{\underline{w}}_{m,m+1}$$

= $\underline{\underline{w}}_{1,m+1}$ proved and verified

Classification Matrix

Accuracy score

$$\frac{\text{Total no of correct predictions}}{\text{Total no of predictions}}$$

How much should accuracy score?
Depends on data & problem.

example 1:

A model predicting by 99% accuracy
that a patient has cancer. We
cannot deploy it due to one life.

example 2

If we want to predict that a
customer this weekend will order
food or not. In case, In 80%. (3)

But this is not much accurate. matrix
 It is not like best,
 Because a number, But not tells type
 and detail exactly.
 To solve this, We can use confusion
 matrix.

confusion matrix

Predicted	
True positive	false negative (Type II)
False positive (Type I)	True negative
J J J A	conect

True / false	positive / negative
from actual data	from model prediction

confusion matrix → accuracy score

$$\text{accuracy score} = \frac{\text{True pos} + \text{True neg}}{\text{True pos} + \text{True neg} + \text{False pos} + \text{False neg}}$$

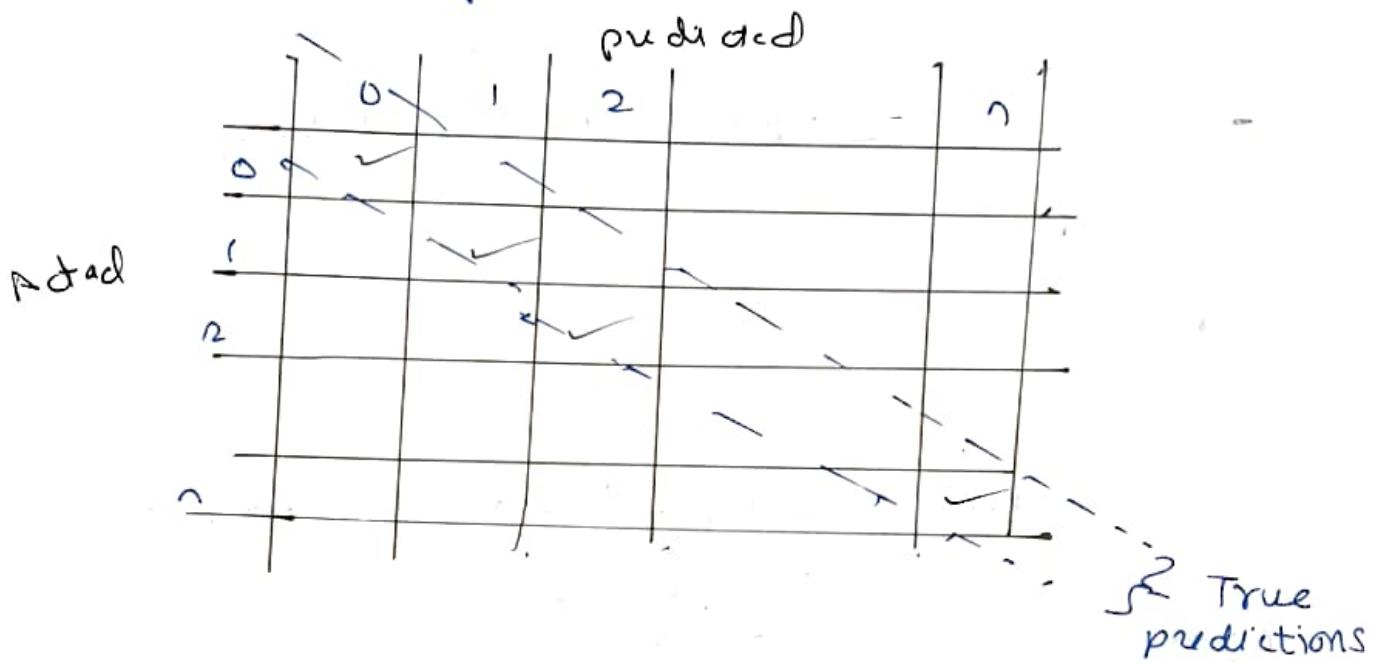
Type I - error

No of false positive

Type - II error

No of false negative

for multiple

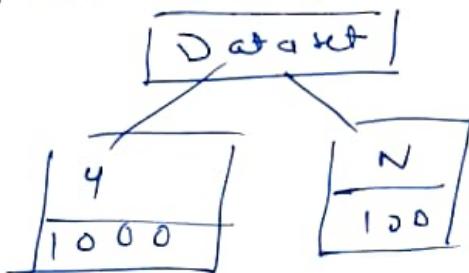


$$\text{Accuracy} = \frac{\text{Total correct predictions} \quad (\text{Diagonal elements})}{\text{Total number of predictions}}$$

Accuracy can mislead too!!

Due to imbalance dataset,

for example



Imbalanced

Accuracy give wrong results

In case

As a manager you have two models.

		Scnt to spam	not scnt to spam
		spam	not spam
model 1	spam	100	170
	not spam	30	700
model 2		Scnt to spam	Not scnt to spam
model 2	spam	100	190
	not spam	10	700

most safe model is, model 2,
as due to false negative

Precision

The ratio of true positives to total positives

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

Again In case

To decide as a manager to decide

		Dected cancer	Not detected cancer
		has cancer	no cancer
model 1	has cancer	1000	200
	no cancer	800	8000
model 2		Dected cancer	Not detected cancer
model 2	has cancer	1000	500
	no cancer	500	8000

most safe model is, model 1

Recall

A mong of all positive, how many positive detected.

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

In case In a problem, we can't decide a that if type I error is more problematic or type II.

example → classify cat/dog

we don't know which is more dangerous.

Type I - said a cat a dog.

Type II - said a dog a cat.

$$\frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

so F1 score

(harmonic mean)

F1 score

$$\frac{2(\text{precision})(\text{recall})}{(\text{precision}) + (\text{recall})}$$

f1 score tends always to lower side

Because to stabilize the model

multi-class

predicted

	Dog	Cat	Rabbit	Total
Dog	25	5	10	40
Cat	20	30	4	34
Rabbit	4	10	20	34
total	29	45	34	

actual

predicted by model Original numbers in data

So let now we will calculate precision

$$\text{pre dog} = \frac{\text{Actual}}{\text{predicted}} = \frac{25}{29} = 0.85$$

$$\text{pre cat} = \frac{30}{45} = 0.66$$

12. n. pre Rab = $\frac{20}{39} = 0.51$

$$\text{macro-precision} = \text{Average}$$

$$= \frac{0.85 + 0.66 + 0.51}{3}$$

$$= 0.70$$

$$\text{weighted precision} = \sum_{\text{classes}} (\text{weight of class}) (\text{pre})$$

$$= \sum_{\text{classes}} \left(\frac{\text{class actual}}{\text{total actual}} \right) (\text{pre})$$

$$= \left(\frac{40}{108} \right) (0.85) + \left(\frac{34}{108} \right) (0.66) + \left(\frac{34}{108} \right) (0.51)$$

$$= 0.71$$

Balance classes - macro

Inbalance class - weighted

$$\text{recall} = \frac{\text{correct dogs}}{\text{actual dogs}} = \frac{25}{40}$$

$$\text{recall dog} = \frac{\text{actual dog}}{\text{correct dogs}} = \frac{25}{40} = 0.625$$

$$\text{recall cat} = \frac{30}{39} = 0.769$$

$$\text{recall rab} = \frac{20}{34} = 0.588$$

macro recall

$$0 \cdot \frac{625 + 852}{3} + 0 \cdot 568 = 0 \cdot 69$$

weighted recall

$$\begin{aligned} &= \left(\frac{40}{108} \right) (0.425) + \left(\frac{34}{108} \right) (0.552) + \left(\frac{34}{108} \right) (0.568) \\ &= (0.37)(0.425) + (0.31)(0.552) + (0.31) \\ &\quad (0.568) \\ &= 0.6869 \end{aligned}$$

$$F_1, \text{dog} = \frac{2 \times \text{Predog Recdog}}{\text{Predog} + \text{recall dog}} = \frac{2 \times 0.86 \times 0.62}{0.86 + 0.62} = \frac{1.064}{1.48} = 0.7202$$

$$F_1, \text{cat} = \frac{2 \times 0.88 \times 0.64}{0.85 + 0.64} = \frac{1.161}{1.49} = 0.7538$$

$$F_1, \text{rab} = \frac{2 \times 0.55 \times 0.58}{0.55 + 0.58} = \frac{0.672}{1.13} = 0.42$$

$$\text{macro } F_1 = \frac{0.7202 + 0.7538 + 0.42}{3} = 0.611$$

weighted F_1

$$\begin{aligned} &= \left(\frac{40}{108} \right) (0.7538) + \left(\frac{40}{108} \right) (0.7202) \\ &\quad + \left(\frac{34}{108} \right) (0.42) \\ &= (0.37)(0.7538) + (0.31)(0.7202) \\ &\quad + (0.31)(0.42) \\ &= 0.2183 \end{aligned}$$

softmax regression (multinomial logistic regression)
going ahead from binary classification.

Sample data

cgpa	Iq	placnet
-	-	yes
-	-	No outout

softmax function

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

k = no of classes

or individual value < 1
add probs equals to 1.

Training part

data		onehot encoding		
cgpa	Iq	0 (y)	1 (n)	2 (No p)
-	-	1	0	0
-	-	0	1	0
-	-	0	0	1

multi to binary (0,1)	
cgpa	Iq
-	-
-	-

split into 3 datasets

master (original) dataset
SPUT

Dataset 1		0
GPA	IQ	
...	...	

logistic regression model 1

$$w_0^{(0)} \quad w_1^{(0)} \quad w_2^{(0)}$$

coeff

Dataset 2

GPA	IQ	1
...	...	

logistic regression model 2

$$w_0^{(1)} \quad w_1^{(1)} \quad w_2^{(1)}$$

coeff

Dataset 3

GPA	IQ	2
...	...	

logistic regression model 3

$$w_0^{(2)} \quad w_1^{(2)} \quad w_2^{(2)}$$

coeff

For just Intuition, In backend process is different

prediction

$$w_0^{(0)} \quad w_1^{(0)} \quad w_2^{(0)}$$

D - yes
 $1, 2 \rightarrow$ other classes

$$w_0^{(1)} \quad w_1^{(1)} \quad w_2^{(1)}$$

1 - No
0, 2, other classes

$$w_0^{(2)} \quad w_1^{(2)} \quad w_2^{(2)}$$

2 - No
0, 1 - other classes

prediction tut, x_1, x_2 .

calculating z_1

$$z_1 = w_0^{(1)} + \sum_{j=1}^2 w_j^{(1)} x_j$$

prob for
yes

$$\sigma(y) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

calculating z_2

$$z_2 = w_0^{(2)} + \sum_{j=1}^2 w_j^{(2)} x_j$$

prob for
no

Apply sigmoid function

$$\sigma(u) = \frac{e^u}{e^u + e^{2u} + e^{3u}}$$

$$z_3 = w_0^{(3)} + \sum_{j=1}^2 w_j^{(3)} x_j$$

prob for
opt out

The largest (probability) will consider (as prediction)

But THIS is seriously time consuming!

In situation

loss function

$$L = \frac{1}{m} \sum_{i=1}^m y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)$$

In logistic regression,

$$L = \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_{ik} \log(\hat{y}_{ik})$$

In softmax.

Sample data

x_1	x_2	y	$y=0$	$y=1$	$y=2$
x_{11}	x_{12}	0	1	0	0
x_{21}	x_{22}	1	0	1	0
x_{31}	x_{32}	2	0	0	1

$$L = \frac{1}{m} \sum_{i=1}^m \left[\sum_{k=1}^K \left(y_k^{(i)} \right) \log \left(\hat{y}_k^{(i)} \right) \right]$$

first row. $k = \{0, 1, 2\}$.

$$\begin{aligned}
 L &= \frac{1}{m} \sum_{i=1}^m \left\{ \left[\left(y_0^{(i)} \right) \log \left(\hat{y}_0^{(i)} \right) \right] + \right. \\
 &\quad \left. \left[\left(y_1^{(i)} \right) \log \left(\hat{y}_1^{(i)} \right) \right] + \right. \\
 &\quad \left. \left[\left(y_2^{(i)} \right) \log \left(\hat{y}_2^{(i)} \right) \right] \right\} \\
 &= \frac{1}{m} \left\{ \left[\left(y_0^{(1)} \right) \log \left(\hat{y}_0^{(1)} \right) \right] + \right. \\
 &\quad \left. \left[\left(y_0^{(2)} \right) \log \left(\hat{y}_0^{(2)} \right) \right] \right\} + \\
 &\quad \left. \left[\left(y_0^{(3)} \right) \log \left(\hat{y}_0^{(3)} \right) \right] \right\} + \\
 &\quad \left. \left[\left(y_1^{(1)} \right) \log \left(\hat{y}_1^{(1)} \right) \right] + \right. \\
 &\quad \left. \left[\left(y_1^{(2)} \right) \log \left(\hat{y}_1^{(2)} \right) \right] \right\} + \\
 &\quad \left. \left[\left(y_1^{(3)} \right) \log \left(\hat{y}_1^{(3)} \right) \right] \right\} + \\
 &\quad \left. \left[\left(y_2^{(1)} \right) \log \left(\hat{y}_2^{(1)} \right) \right] + \right. \\
 &\quad \left. \left[\left(y_2^{(2)} \right) \log \left(\hat{y}_2^{(2)} \right) \right] \right\} + \\
 &\quad \left. \left[\left(y_2^{(3)} \right) \log \left(\hat{y}_2^{(3)} \right) \right] \right\} + \\
 &\quad \left. \left[\left(y_3^{(1)} \right) \log \left(\hat{y}_3^{(1)} \right) \right] + \right. \\
 &\quad \left. \left[\left(y_3^{(2)} \right) \log \left(\hat{y}_3^{(2)} \right) \right] \right\} + \\
 &\quad \left. \left[\left(y_3^{(3)} \right) \log \left(\hat{y}_3^{(3)} \right) \right] \right\}
 \end{aligned}$$

$$L = \underbrace{(y_1^{(1)})}_{\text{so, thus for 3 rows.}} \log(\hat{y}_1^{(1)}) + \underbrace{(y_2^{(2)})}_{\text{for 3 rows.}} \log(\hat{y}_2^{(2)}) + \underbrace{(y_3^{(3)})}_{\text{for 3 rows.}} \log(\hat{y}_3^{(3)})$$

$$\hat{y}_1^{(1)} = \sigma(w_1^{(1)}x_{11} + w_2^{(1)}x_{12} + w_0^{(1)})$$

$$\hat{y}_2^{(2)} = \sigma(w_1^{(2)}x_{21} + w_2^{(2)}x_{22} + w_0^{(2)})$$

$$\hat{y}_3^{(3)} = \sigma(w_1^{(3)}x_{31} + w_2^{(3)}x_{32} + w_0^{(3)})$$

So, how are coeffs of us model.

$$\begin{bmatrix} w_0^{(1)} & w_1^{(1)} & w_2^{(1)} \\ w_1^{(1)} & w_2^{(1)} & w_0^{(1)} \\ w_1^{(2)} & w_2^{(2)} & w_0^{(2)} \\ w_1^{(3)} & w_2^{(3)} & w_0^{(3)} \end{bmatrix}$$

while applying gradient descent, I have to take derivative w.r.t everyone

$$\frac{\partial L}{\partial w_0^{(1)}}, \frac{\partial L}{\partial w_1^{(1)}}, \frac{\partial L}{\partial w_2^{(1)}} \dots$$

9 values.

Initialise value of matrix
epochs for loop

update everyone coeff.

$$w_1^{(1)} = w_1^{(1)} - \eta \frac{\partial L}{\partial w_1^{(1)}}$$

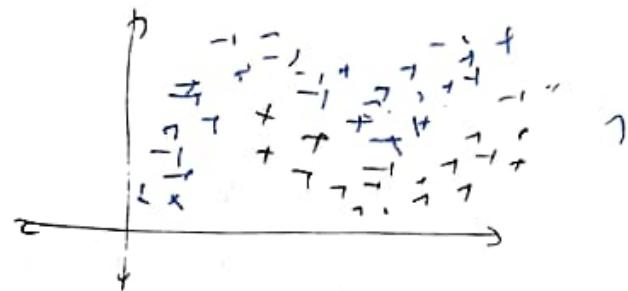
$$w_2^{(1)} = w_2^{(1)} - \eta \frac{\partial L}{\partial w_2^{(1)}}$$

⋮

9 classes.

Polynomial Logistic Regression

We can convert to polynomial features



For 2 features of 3 degree,

x_0, x_1, x_2 ,
we have

$$x_0^0, x_1^0, x_2^0, x_1^1, x_2^1, x_1^2, x_2^2$$

Output \rightarrow Input columns
= degree + 1

Decision Tree

example

Gender	Occupation	Suggestion
F	stu	PUBG
F	prog	Github
M	prog	whatsapp
F	prog	Github
M	stu	PUBG
M	stu	PUBG

if occup == student:

 print(PUBG)

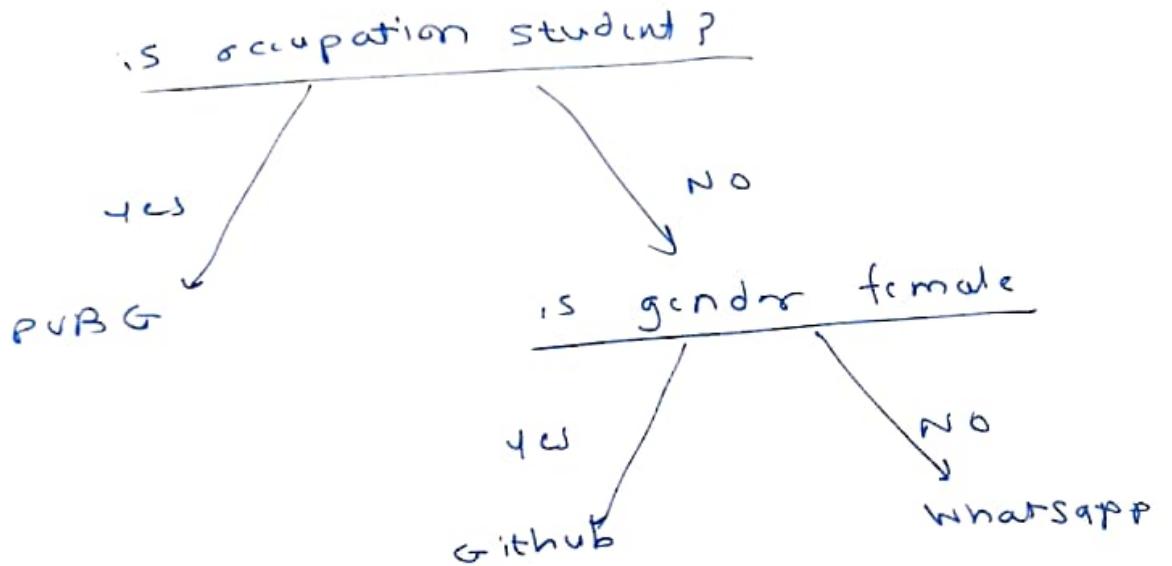
else if gender == female

 print(Github)

else

 print(whatsapp)

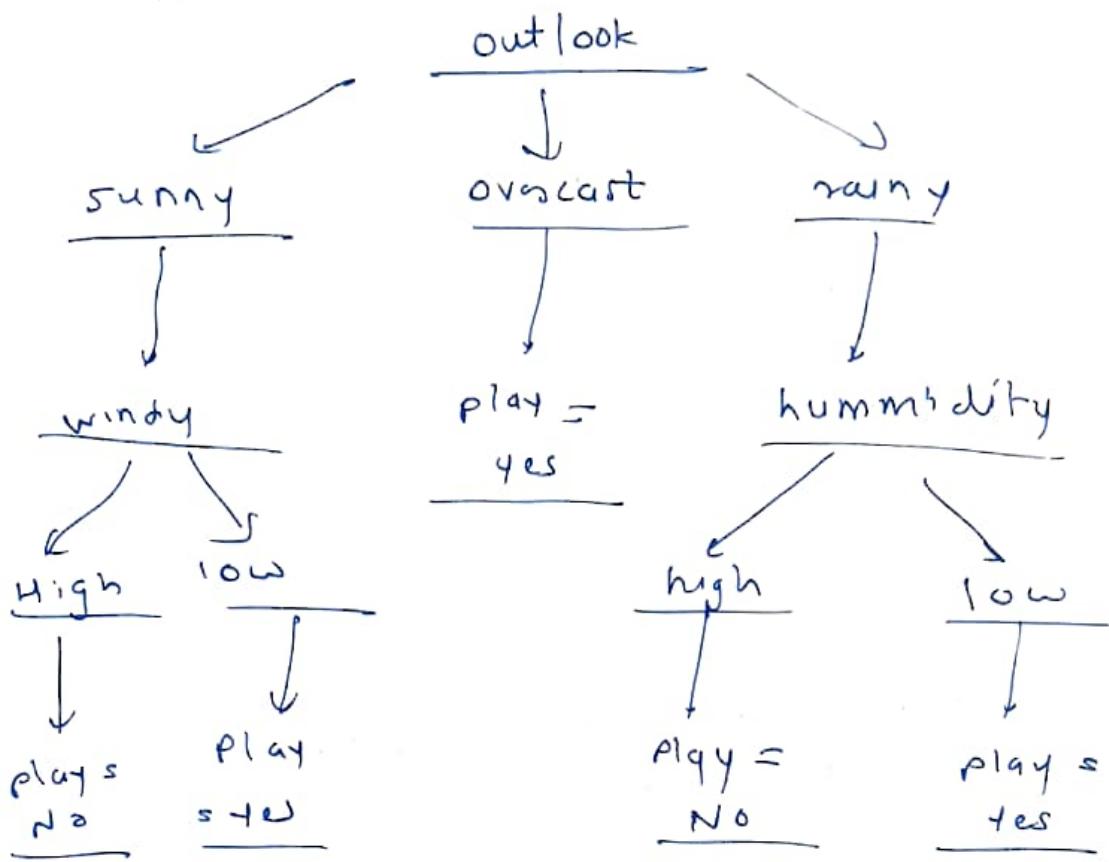
Tree



example 2

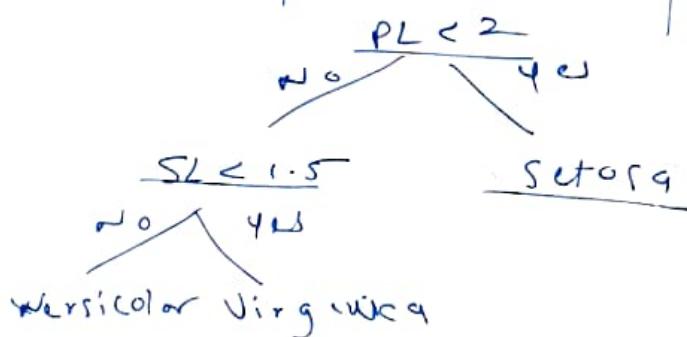
Day	Outlook	Temp	Humid	Wind	play
1	sunny	Hot	High	weak	No
2	Sunny	Hot	High	strong	No
3	Overcast	Hot	High	weak	Yes
4	Rain	mild	High	weak	Yes
5	Rain	cool	Nor	weak	Yes
6	Rain	cool	Nor	strong	No
7	overcast	cool	Nor	strong	Yes
8	Sunny	mild	High	weak	No
9	sunny	cool	Nor	weak	Yes
10	rain	cloudy	Nor	weak	Yes
11	Sunny	mild	Nor	strong	Yes
12	overcast	mild	High	strong	Yes
13	overcast	Hot	Nor	weak	Yes
14	rain	mild	High	strong	No

Tree

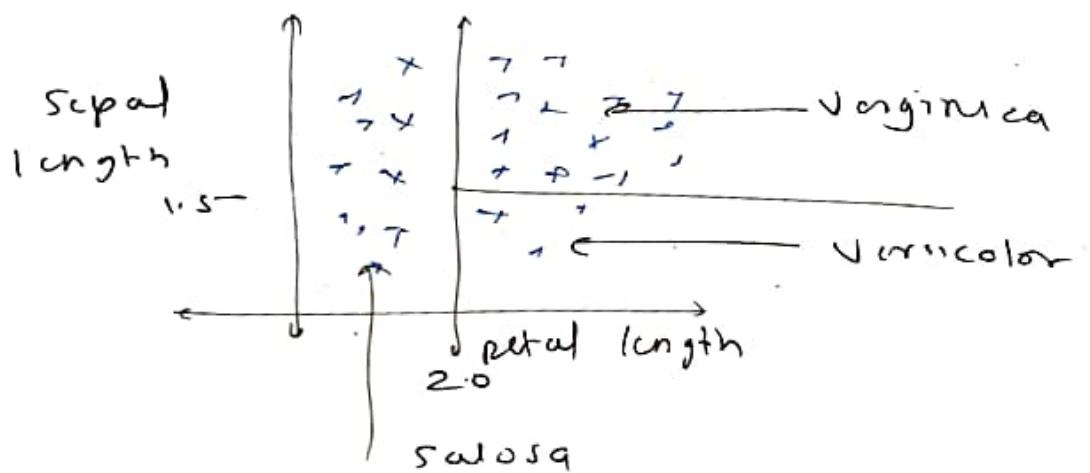


example 3

petal length	sepal length	Type
1.34	0.34	Setosa
3.45	1.35	Versicolor
1.09	0.90	Setosa
2.56	1.79	Versicolor
3.00	1.13	Setosa
1.3	0.86	



Geometric Intuition



Algorithm

A training data with some feature variables and classification or regression output.

Determine best feature to split data, & later on also.

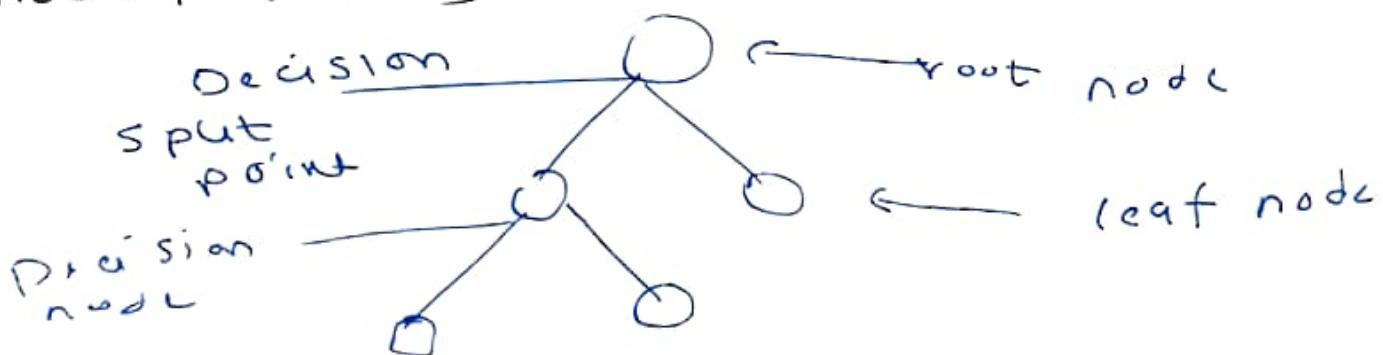
Split data in features for bwt subset contain correct values for feature. It is a node & tree structure.

Recursively generate nodes.

Mathematically, DT's are using hyperplanes which run parallel to any axes to cut your coordinate system in hyper cuboids.

Programmatically, DT's are if-else statements

How to decide that which node to be a root node? (Next)



Advantages

Intuitive & easy to understand
minimal data preparation
log time complexity

Disadvantages

overfitting

prone to error for imbalance datasets

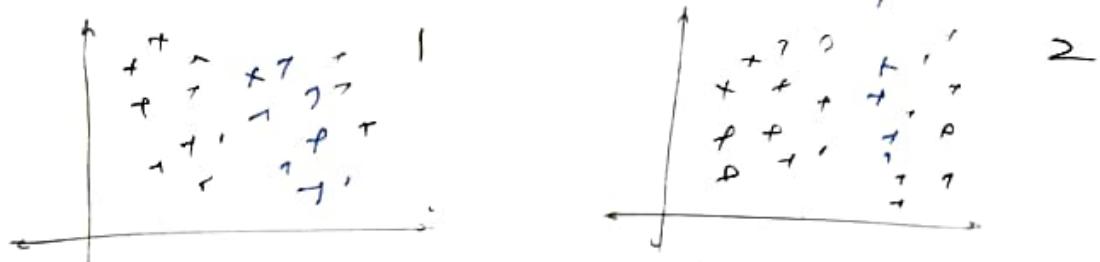
CART - Classification And Regression Tree

Entropy

measure of disorder / purity - impurity

ice	water	vapour
less entropy		more entropy
more order		less order

more knowledge = less entropy



model 1 → more diff for confidence
dataset → high entropy

The mathematical formula for entropy is

$$E(S) = - \sum_{i=1}^c p_i \log_2 p_i$$

p_i - frequentist probability of an element/
class in our data

e.g. if our data has only two classes
yes and no.

$$E(D) = - P_{\text{yes}} \log_2 (P_{\text{yes}}) - P_{\text{no}} \log_2 (P_{\text{no}})$$

sample is dataset

salary	age	phrasal
20000	21	Yes
10000	45	No
60000	27	Yes
15000	31	No
12000	18	No

$$\begin{aligned}
 H(\hat{\alpha}) &= -P_1 \log_2(P_1) - P_0 \log_2(P_0) \\
 &= (2/5) \log_2(2/5) - (3/5) \log_2(3/5) \\
 &= 0.97
 \end{aligned}$$

salary	age	phrasal
34000	31	No
15000	25	No
69000	57	Yes
25000	21	No
32000	28	No

$$\begin{aligned}
 P(\hat{\alpha}) &= (-1/5) \log_2(-1/5) + (-4/5) \log_2(-4/5) \\
 &= 0.72
 \end{aligned}$$

salary	age	phrasal
20000	45	No
15000	45	No
60000	34	No

$$f(\theta) = (0.5) \log_2(0.5) - (0.4) \log_2(0.4)$$

$$= 0$$

Calculate entropy for 3 classes

Salary	Age	Pharmacy
20000	21	44
10000	45	No
60000	27	44
15000	31	No
30000	30	maybe
12000	18	No
40000	40	maybe
20000	20	maybe

$$\begin{aligned}
 H(\theta) &= p_Y \log_2 p_Y + p_N \log_2 p_N + p_{MB} \log_2 p_{MB} \\
 &= -\left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right) - \left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) + \left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) \\
 &= 1.56
 \end{aligned}$$

Observation

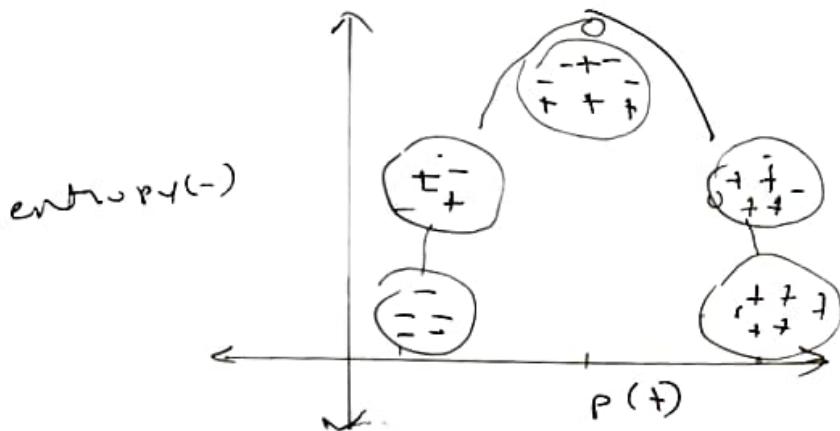
more uncertainty - more entropy

for 2 class problem, min entropy = 0
max entropy = 1.

for > 2 classes, min entropy = 0
max entropy > 1 .

Can we use $\log_2 a$ or $\log_e a$

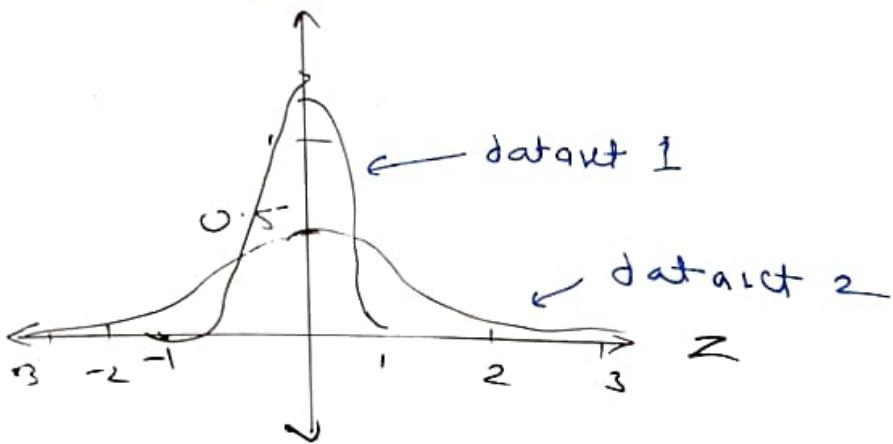
Entropy vs probability



entropy for continuous variables

Area	Built-in	Price	Area	Built-in	Price
1200	1999	3.5	2200	1989	4.6
1800	2011	5.6	600	2018	6.5
1400	2000	7.3	1100	2005	12.8

find most entroped dataset.



Due to less range (-1, 1), we know more about dataset 1. Hence less entropy.

entropy
dataset 1 < entropy
dataset 2

Information gain

Measures the quality of a split.

It's based on entropy + the decrease in entropy after a data set is split or a column. Constructing a decision tree is all about finding attribute that returns the highest information gain.

$$\text{Information gain} = E(\text{Parent}) - \frac{\text{weighted average}}{E(\text{children})}$$

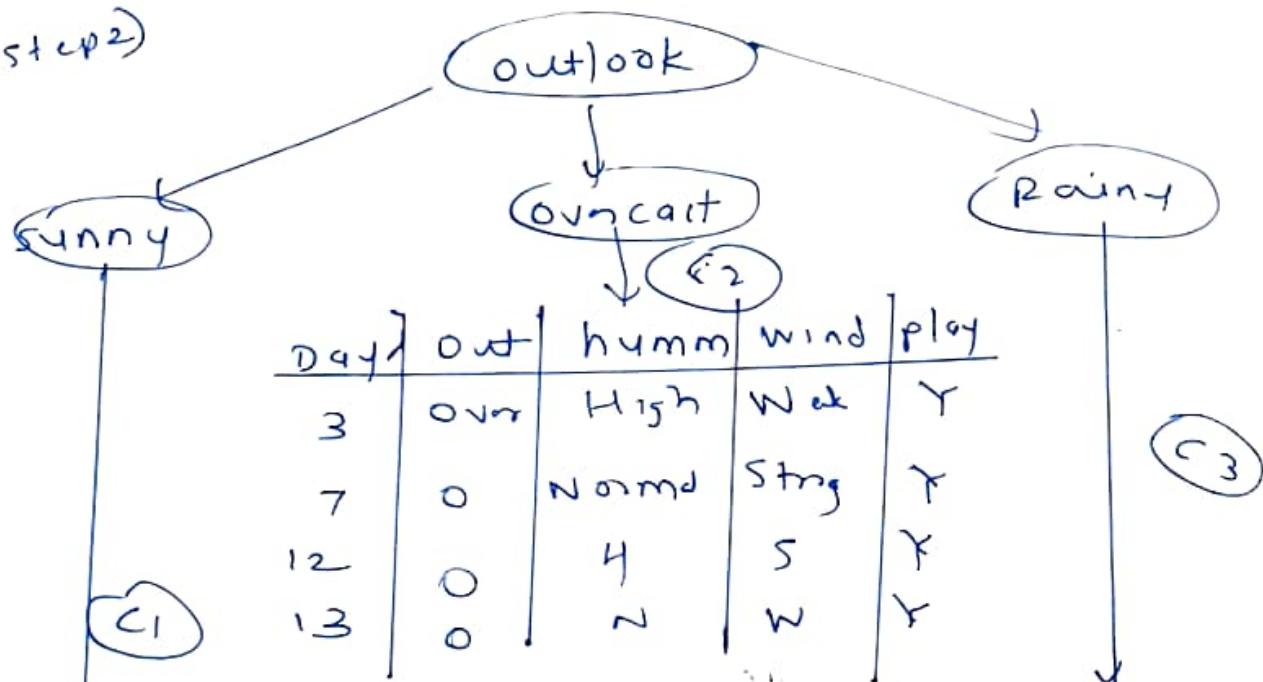
example: see tennis dataset again.

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	H	H	True	No
Overcast	H	H	F	Yes
Rainy	Mid	H	F	F
Rainy	Cool	Normal	F	Yes
Rainy	C	N	T	No
Overcast	C	N	T	F
Sunny	N	H	F	F
Sunny	C	N	F	T
Rainy	M	N	F	Yes
Sunny	M	N	T	Yes
Overcast	M	H	T	Yes
Overcast	H	N	F	Yes
Rainy	M	H	T	No

step 1

$$\begin{aligned} \text{root node entropy } E(P) &= -P_Y \log_2(P_Y) - P_R \log_2(P_R) \\ &= \left(\frac{9}{14}\right) \log_2\left(\frac{9}{14}\right) - \left(\frac{5}{14}\right) \log_2\left(\frac{5}{14}\right) \\ &= 0.94 \end{aligned}$$

step 2



Day	outlook	Hummm	wi	play	Day	outlook	hummm	Wind
1	Sunny	High	Weak	No	4	Rainy	High	Weak
2	S	H	Strong	No	5	R	Normal	W
8	S	H	W	N	6	R	N	Strong
9	S	Normal	W	Yes	10	R	N	W
11	S	N	S	Y	14	R	H	S

$$\text{for } c_1 \rightarrow (-2/5) \log(2/5) + (3/5) \log(3/5) = 0.97$$

$$\text{for } c_2 = (-5/15) \log(5/15) + (10/15) \log(10/15) = 0$$

$$\text{for } c_3 = (-3/15) \log(3/15) - (2/15) \log(2/15) = 0.97$$

Step 3) - calculate weighted entropy

$$we\ ent = \left(\frac{5}{14} \right) (0.97) + \left(\frac{4}{14} \right) (0) + \left(\frac{5}{14} \right) (0.97)$$

$$we\ (child) = 0.49$$

Step 4)

$$\text{Info gain} = E(\text{parent}) - \frac{\text{weight}_{avg}}{E} (\text{children})$$

$$= 0.97 - 0.49 = 0.48$$

so information gain (Increase / Decrease) in entropy / Impurity on the basis of outlook column is 0.48.

Step 5) calculate Information Gain for all the columns.

which ever column has the highest info gain (maximum decrease in entropy)

the algo will select that column to split data

Step 6)

DT then apply a greedy search algorithm in top-bottom fashion to find the info gain at every level of the tree.

Once a Leaf node get reached. ($\text{Entropy} = 0$) no more splitting is done.

Gini - impurity - For scikit-learn formula:

$$GI = 1 - (p_1^2 + p_n^2)$$

Go on next page for example for gini - impurity.

salary	σ_1^2 age	phrases	salary	σ_2^2 age	phrases
20k	21	+	34k	31	-
10k	45	-	15k	25	-
60k	27	+	69k	57	+
15k	31	-	25k	21	-
12k	18	-	32k	28	-

$$\begin{aligned}
 G_{i1} &= 1 - (\rho_{y_1}^2 + \rho_{n_1}^2) \\
 &= 1 - ((2/5)^2 + (3/5)^2) \\
 &= 1 - (4/25 + 9/25) \\
 &= 1 - \frac{13}{25} = \frac{12}{25} = 0.48
 \end{aligned}$$

$$\begin{aligned}
 G_{i2} &= 1 - (\rho_{y_2}^2 + \rho_{n_2}^2) \\
 &= 1 - ((1/5)^2 + (4/5)^2) \\
 &= 1 - (1/25 + 16/25) \\
 &= 1 - \frac{17}{25} \\
 &= \frac{8}{25} \\
 &\approx 0.32
 \end{aligned}$$

$$G_{i2} < G_{i1}$$

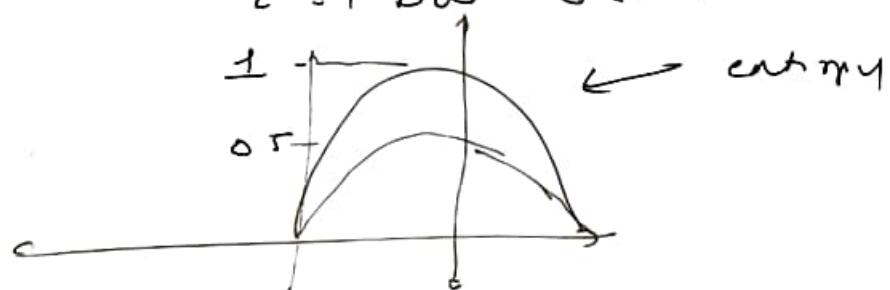
Difference between Gini and entropy

case 1: $P_1 = 0$ & $P_N = 1$.

$E = 1$ and $G_i = 1$.

case 2: $P_1 = 0.5$ & $P_N = 0.5$

$E = 1$ but $G_i = 0.5$



Why gini?

fast on board computing

But in some cases, Entropy can give better splits.

Handle numerical

example

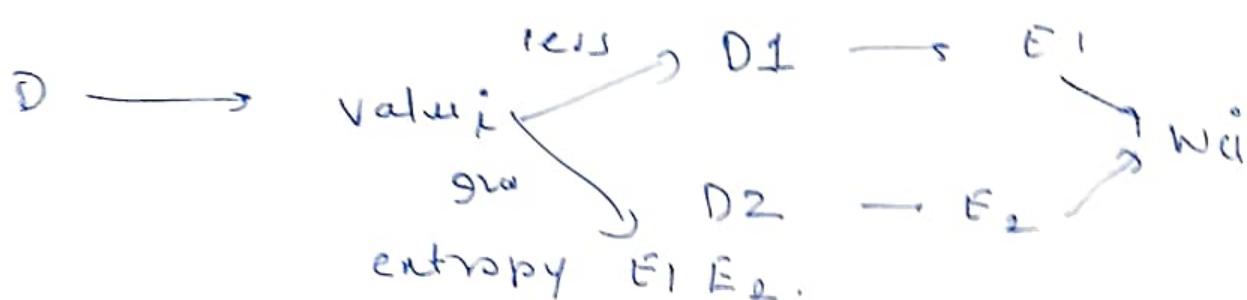
No	User rating	Download
1	3.5	Y
2	4.6	Y
3	2.2	N
4	1.6	N
5	4.1	N
6	3.9	N
7	3.2	Y
8	2.9	Y
9	4.5	N
10	3.3	Y
11	2.5	Y
12	1.0	
13		

Step 1) Sort data on basis of numerical column

No	rating	Download
1	1.6	Y
2	1.9	N
3	2.2	N
4	2.5	Y
5	2.9	Y
6	3.2	N
7	3.3	N
8	3.5	Y
9	3.9	N
10	4.1	N
11	4.6	Y
12	4.8	Y
13		

Step 3) Split in parts, for every row value
data < value < data
1 2

Step 4) for every row element in column



for all values in column.

$w_i = \text{weighted entropy}$

D → Dataset

Now we will calculate

$$E(\text{parent}) - E(\text{child}) = \text{Information Gain}$$
$$IG_i$$

find $\max(IG_1, IG_2, \dots, IG_n)$

let us I assume it is IG_3 .
Now, use it as splitting

criterion



step 9) Do it recursively until to get all leaf nodes.

This thing looks computationally too expensive. Is it right way of splitting criterion?

Yes.

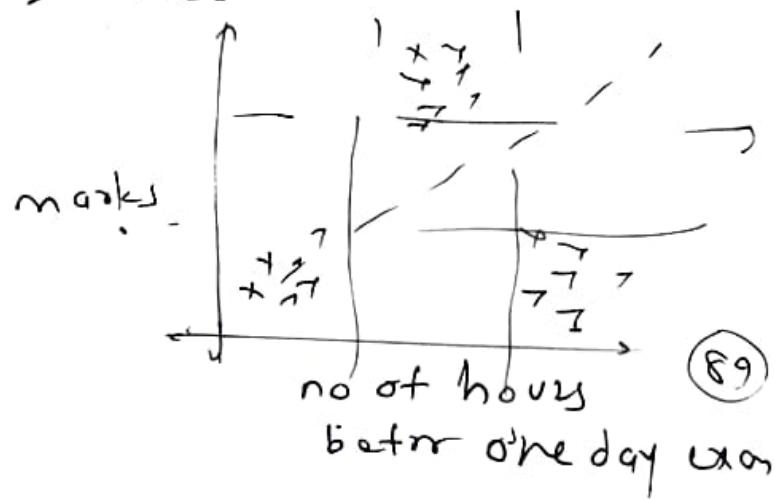
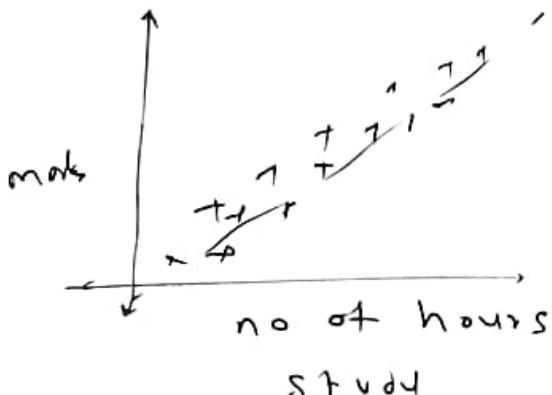
It's worth it.:)

Because compute is in only training part, & not in testing part much.

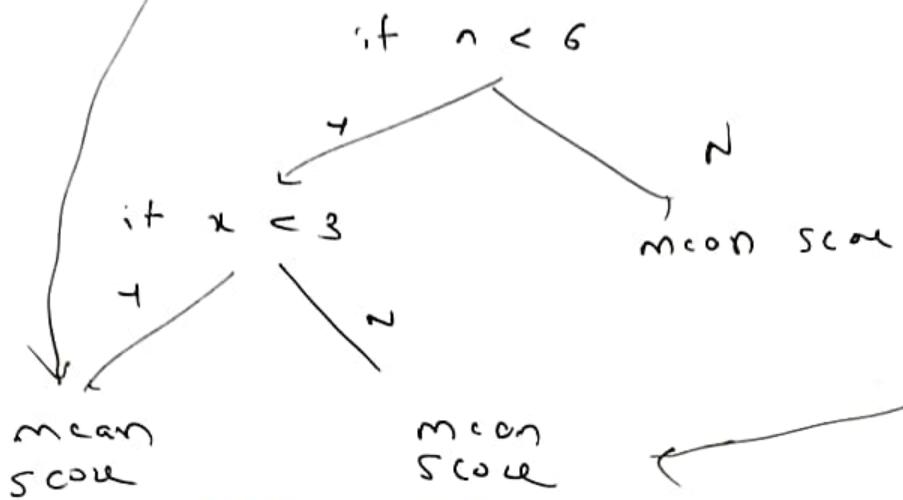
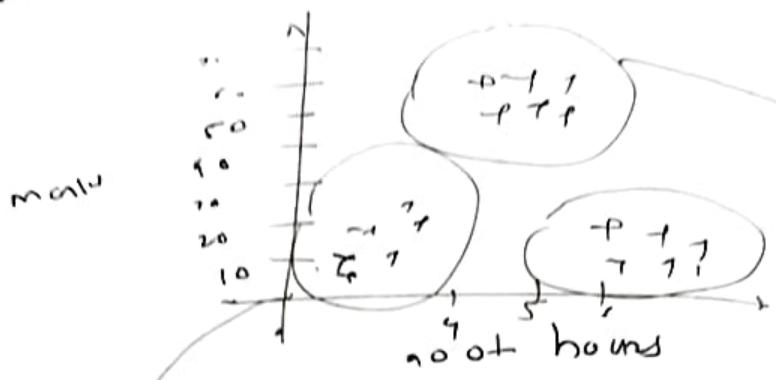
for too much numerical column, Avoid decision trees.

Regression trees

In non linear situation, regression tree gives us better results.



(89)

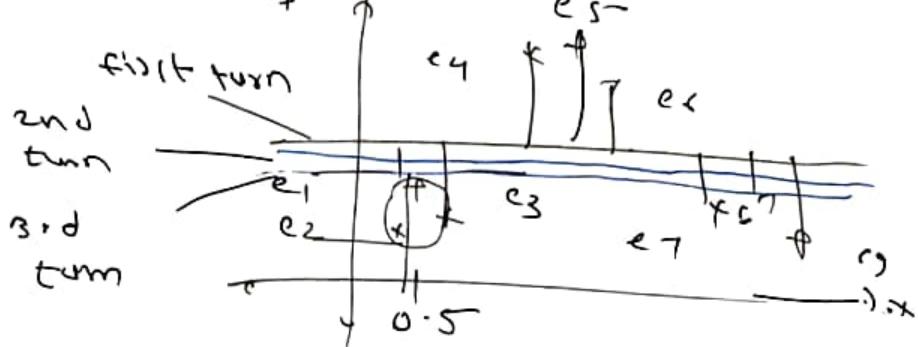


But what if the curve is like this

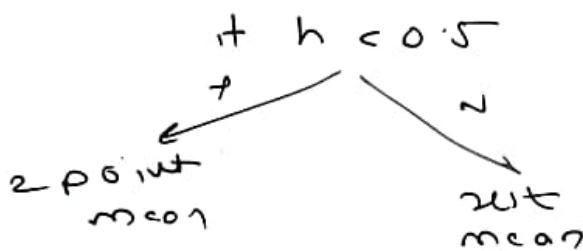


splitting criteria

let us go with an example

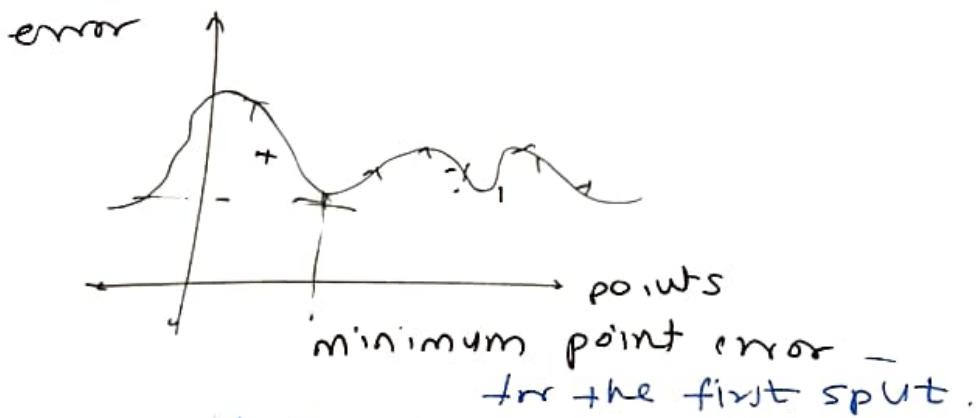


consider first two points, find m_1
 consider next set of points find m_2
 calculate distances.

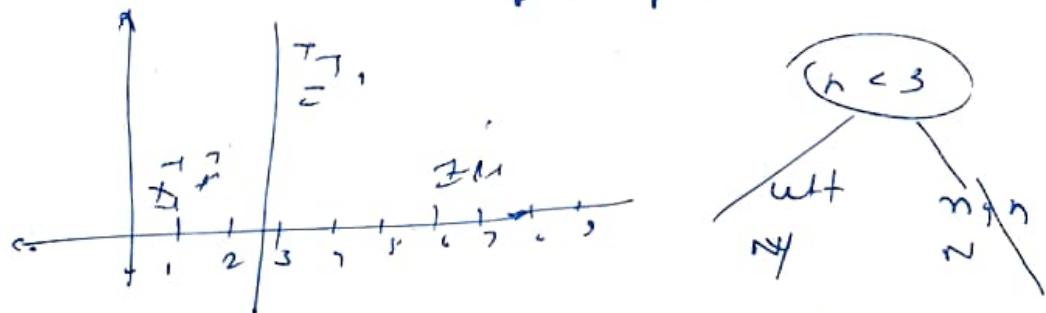


$$SSE = e_1^2 + e_2^2 + \dots + e_n^2$$

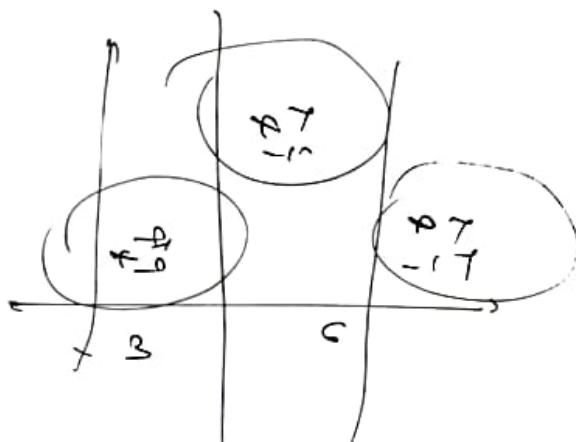
Now consider next 2 points. Calculate both "left" & "right" means & so on. and plot them.



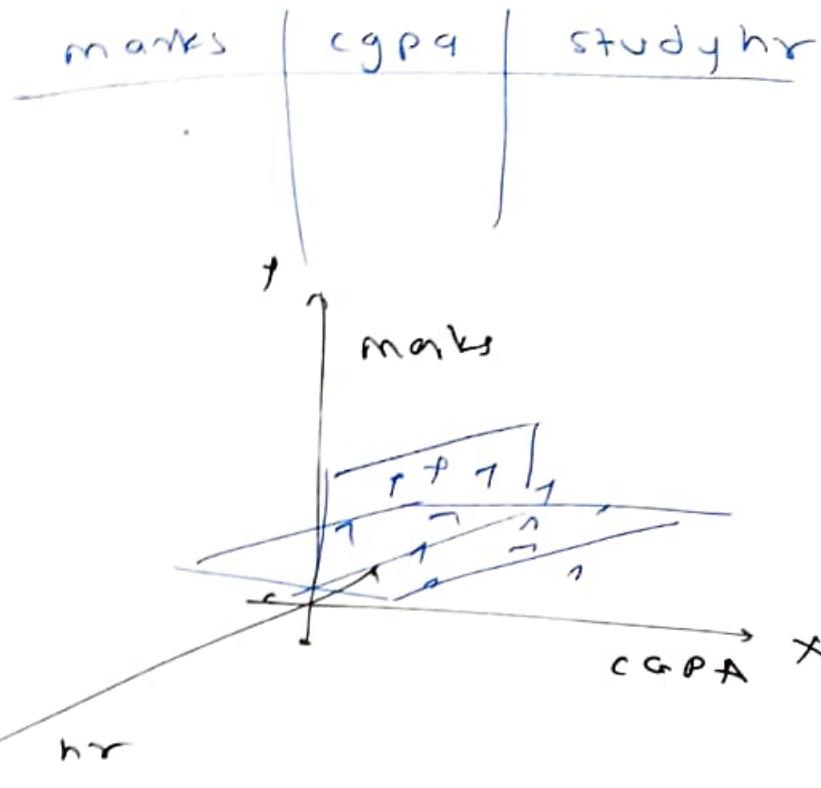
Now, first one to split and repeat process.



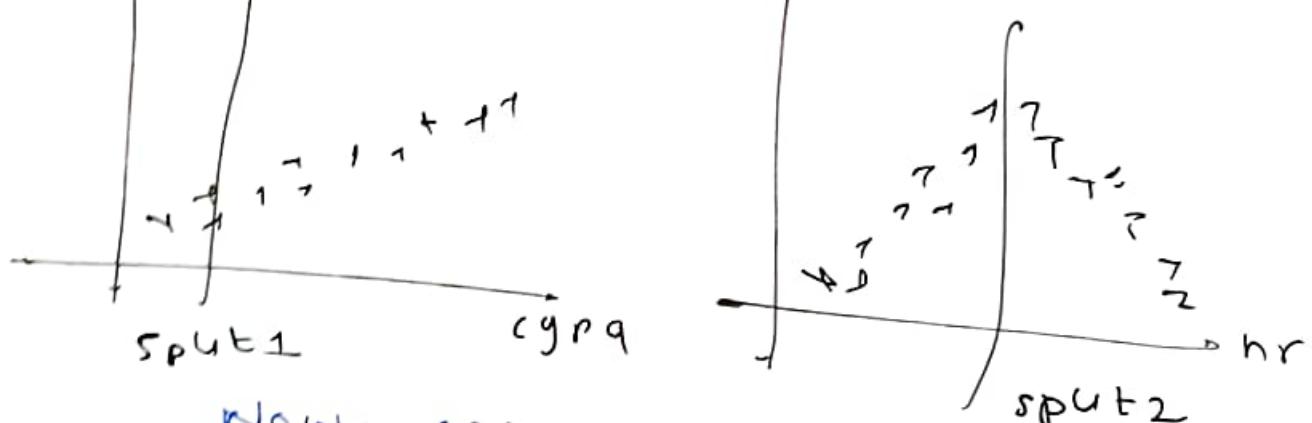
But if we split on and on, then, we seriously split on every point (overfit). So, we will set a threshold like if only points $<$ threshold, we will not split. or continue to next points,



still this is in 2D what if higher dimension
sample data



marks so, we will do it independently.



Now, compare the errors

it $SE_{hr} < SE_{CGPA}$
what ever small, cut parallel to
threshold, and follow on

Ensemble learning

A collection of multiple models

Wisdom of the crowd

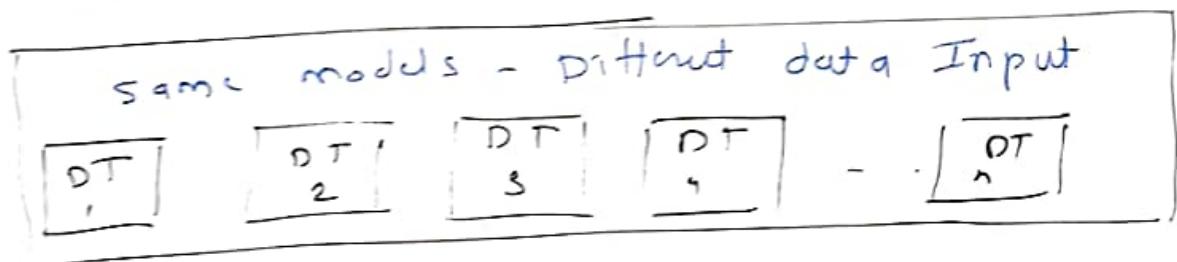
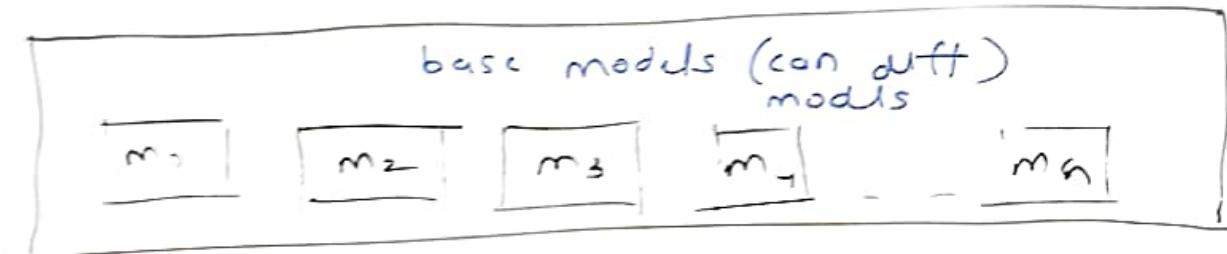
Audience poll in KBC.

No of reviews roles on Amazon.

→ trust of people on a democratic party.

Core idea

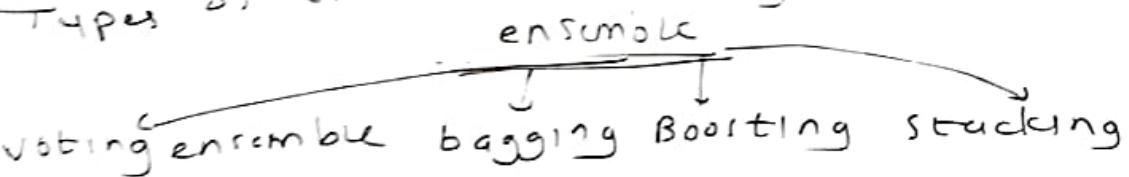
Prediction



for classification
Data will be given to models. Each model will predict, and majority will be considered as a class.

for regression
Data will be given to models. Each model will predict, and we will calculate mean.
Simple but effective

Types of ensemble learning



random forest

- ada boost

- gradient - boosting

- xgBoost (extreme gradient boosting)

(93)

why so many types?

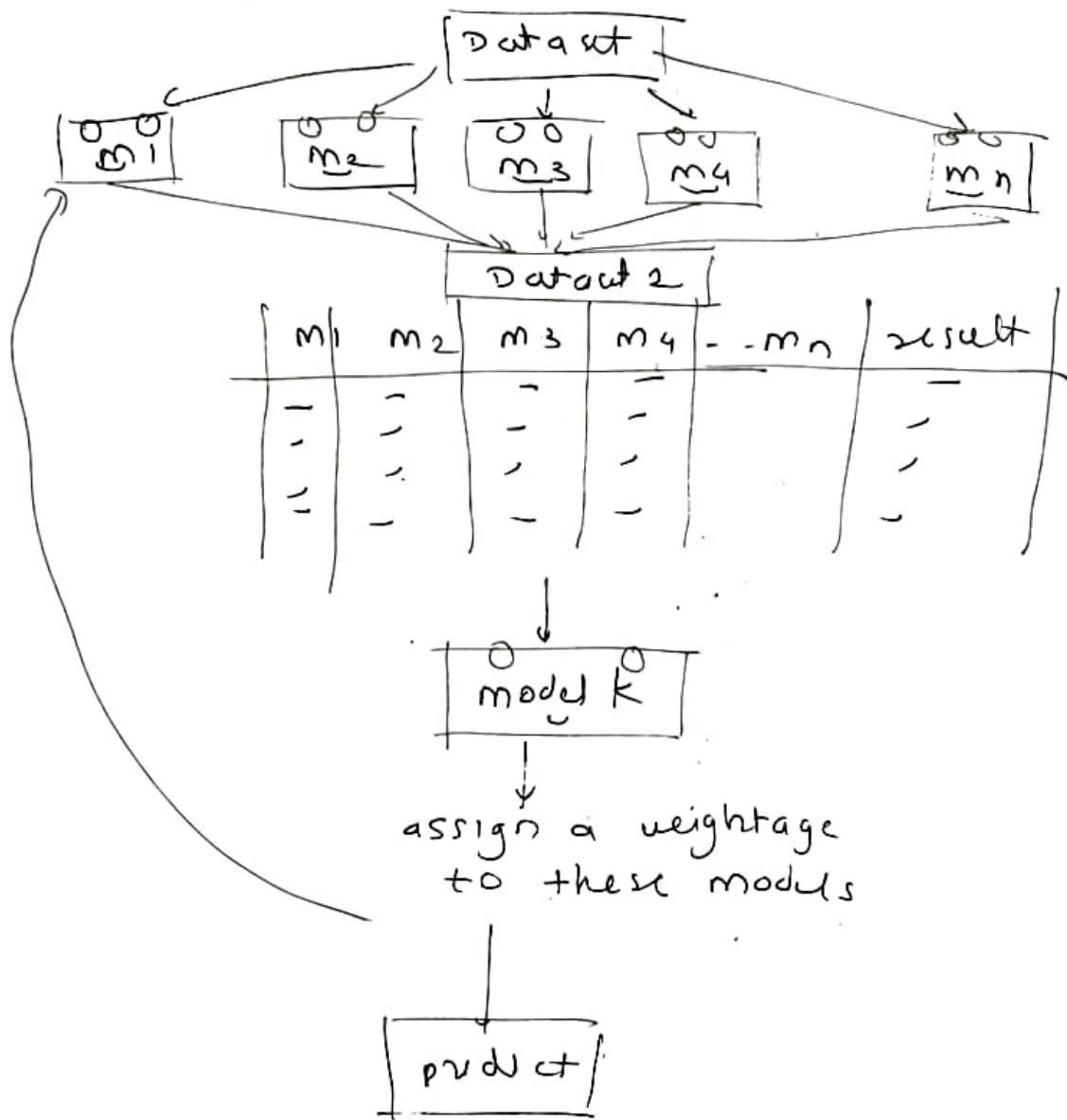
we have two majors, train diff models, or
train diff data. That's why difference.

Light Intro.

Voting ensemble - Base models different algorithms
give inputs.

- + classification,
go with majority
- + regression,
Go with mean.

Stacking -



+ train slice diff models,

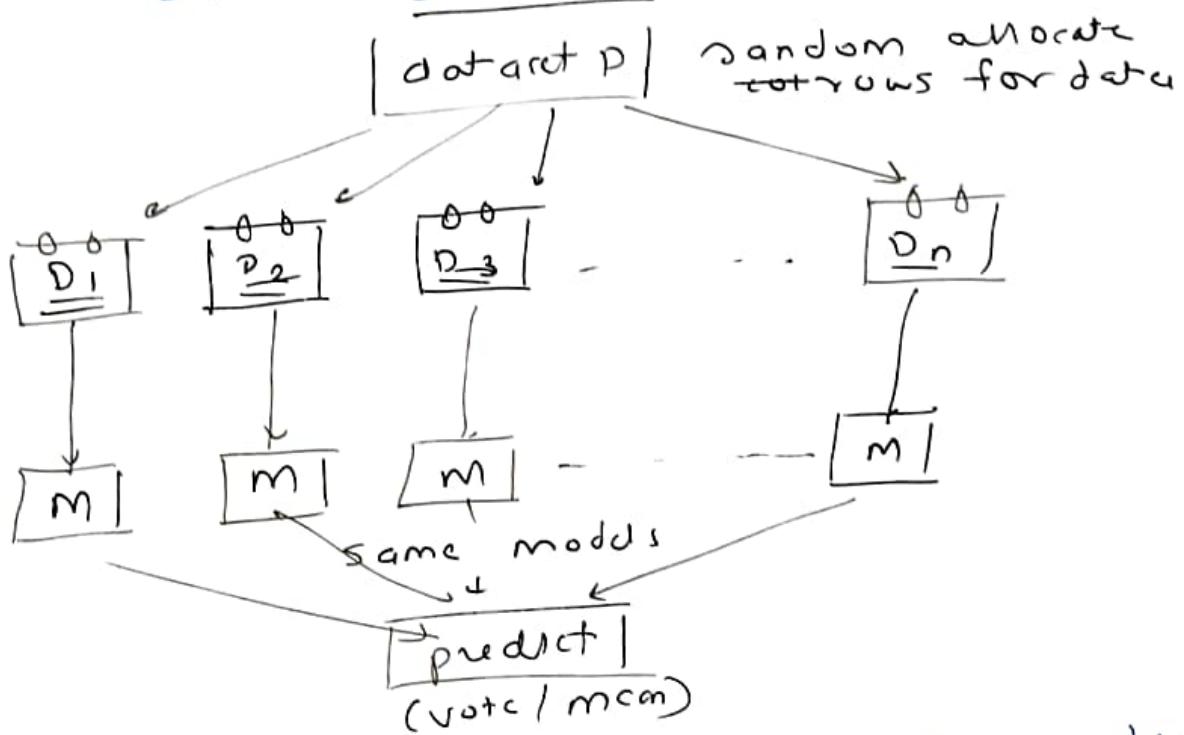
make a dataset, their model & output &
actual results.

+ train a master model, & decide weights.
and predict.

Bagging

Bootstrapped aggregation

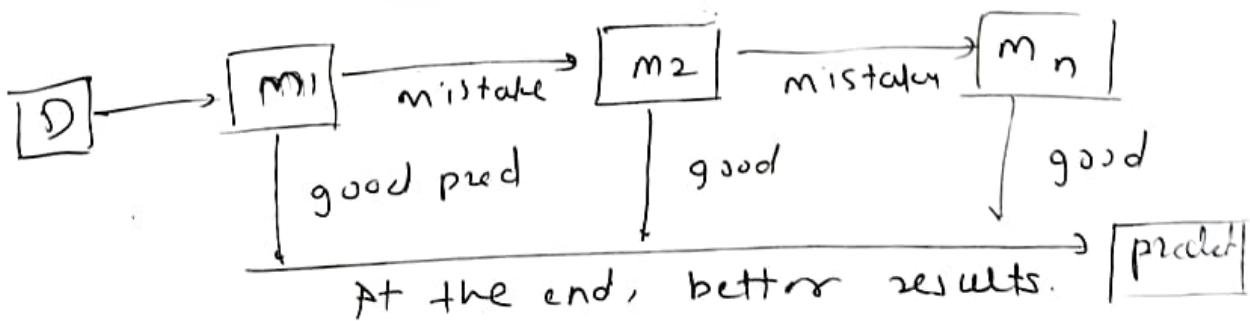
Same algorithm, diff data Input



random forest is a subset of ensemble training - (bagging) (diff trees)

Boosting

Boosting results in a series!



But still why it works?

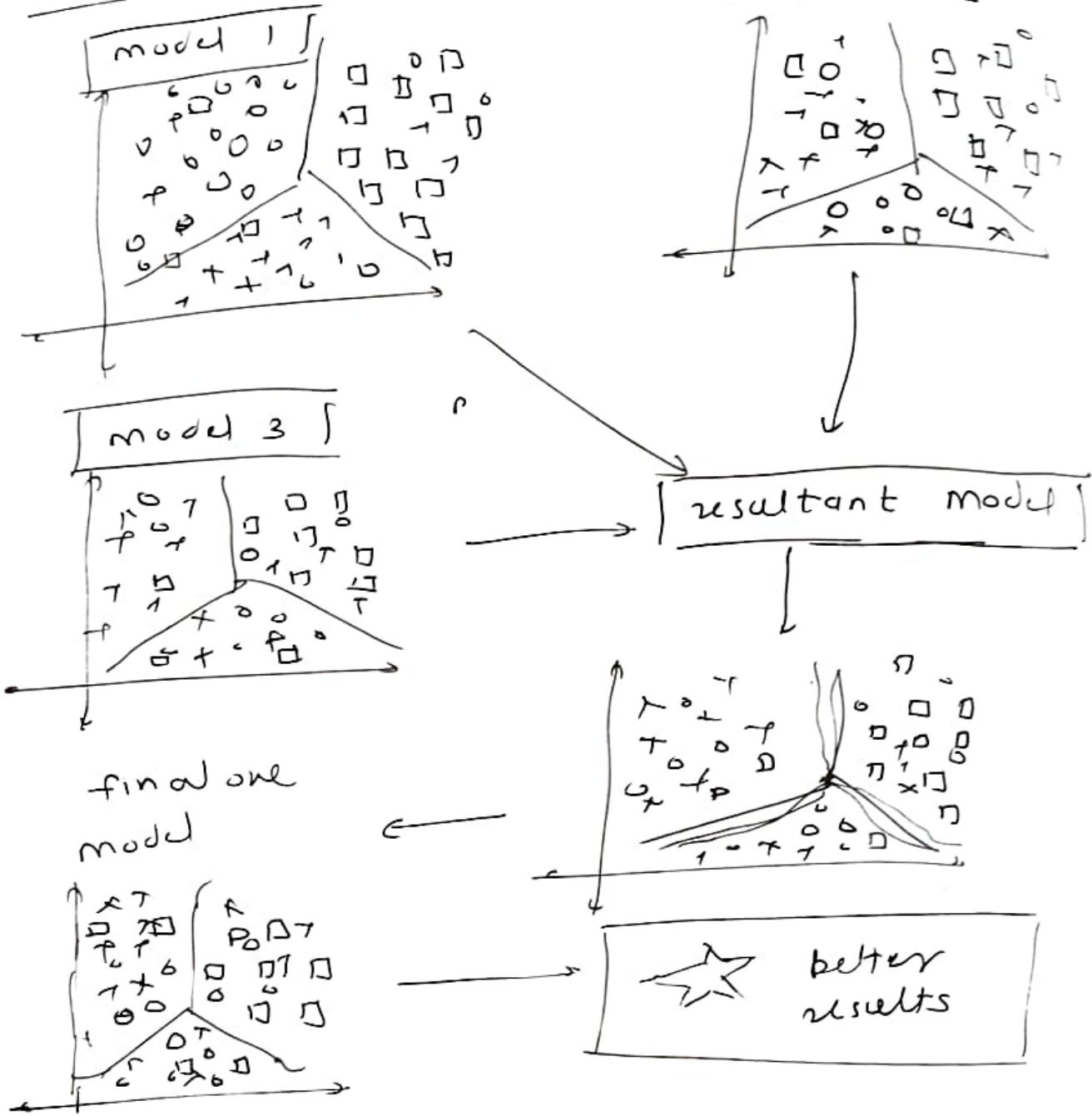
Let us see in the figure,

there are two different reasons.

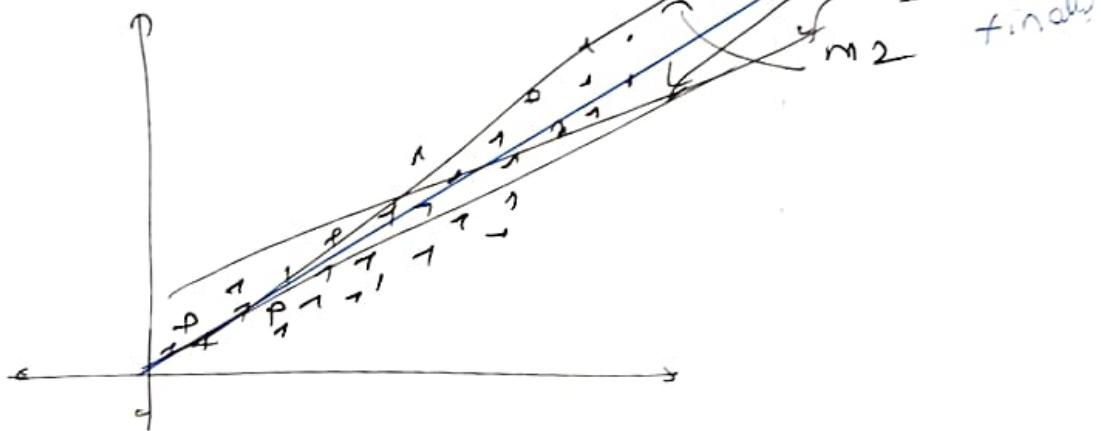
⇒ Classification

⇒ regression

Classification



Regression



High computation

Advantage

- Improvement in performance
- 2) Reduce bias and variance
It helps into achieve low (bias + var)
- 3) Robustness

When to use ensemble learning?

Always.

non reason to don't use it!

1. Voting ensemble

Please read the techniques, from intro,
how to work.

But why it work?

How it is possible that it outperform
from the combo of all algorithm?

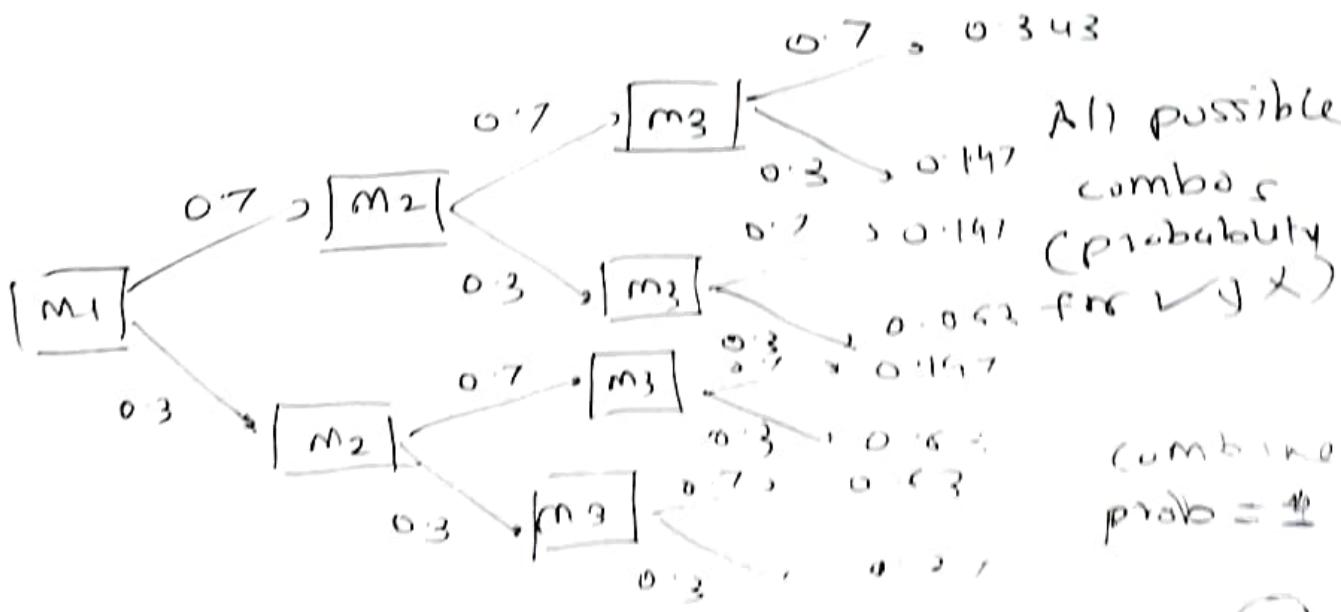
Assumption

models should be independent.

for all models, min accuracy should be
51% at least, else it will be worst
one ever.

example:

we have 3 models, each accency 0.7,



so, in the figure, the out of 3, the right prob at least 2 is, there are 4 chances where accuracy (votes) are at least 2 times same.

If we add, we get entire models right probability of prediction.

$$34.3 + 14.7 + 14.7 + 14.7 = \underline{0.78}$$

so by this way, we prove that a combination work.

But we have a rule that, accuracy at least for an individual model should be $\geq 51\%$ at least.

Let us prove this, by assuming it is 50%.

Now, instead of 0.7, assume 0.3 place.

Now, by combining these 0.3 models, we will get overall less accuracy.

So, $100 - 78\%$. (from true right).

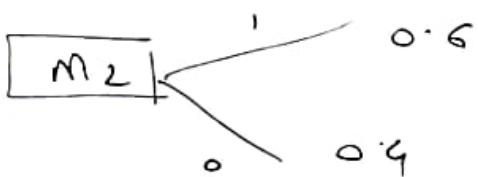
$$= 0.22 \quad \text{overall}$$

hard voting / soft voting



hard voting

$$\text{majority} = 1 \text{ so } 1.$$



soft voting

$$1 = \frac{0.8 + 0.6}{2} = 0.7 \text{ so } 1.$$

$$0 = \frac{0.2 + 0.4}{2} = 0.4$$

generally soft voting better.

Regression

see part of Intro for the regression.

Bagging

bootstrapping + aggregation

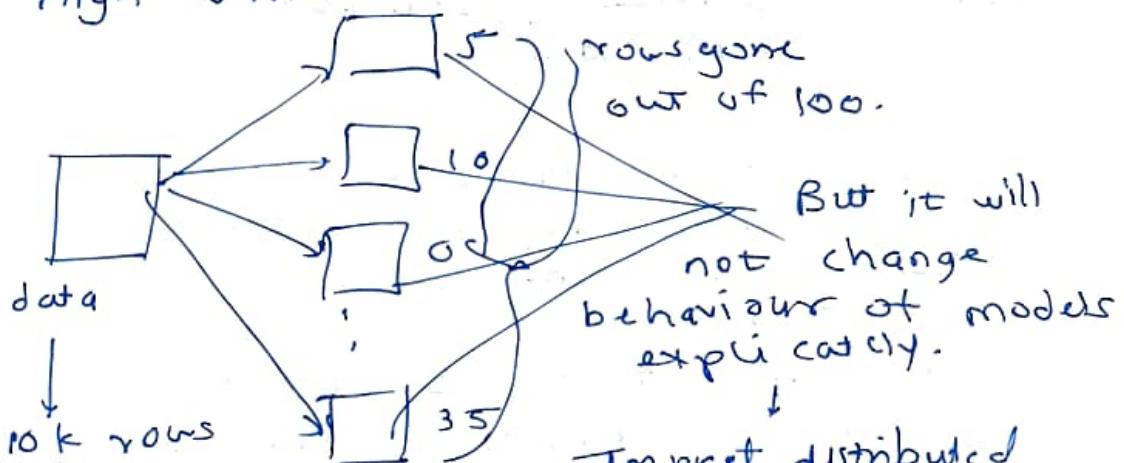
see intro part for intuition.

But why bagging improves score?

we can get } LB, MW } combu.
} HB, LN }

How?

we take models as a base for low bias
and high variance



We replace 100 rows, which are different from last data.

When to use?

Always give a try for it.

so we have random forest.

Types of bagging

can do row + columns

pasting - without replacement row sampling

Random Subspace - bagging with replacement

random patches - rows + column Sampling

No high variance affected.

OOB score

out of bag score

since we select at random, some rows go again and again, & some don't go. so they are called out of bag.

statistically it is proven that only 63% of rows go in. rest 37% never even go.

Bagging generally perform better than pasting

good results come around 25% to 50%. random patches & subspace could use on high dim data.

Random forest

can solve any ML problem

good fit in any project

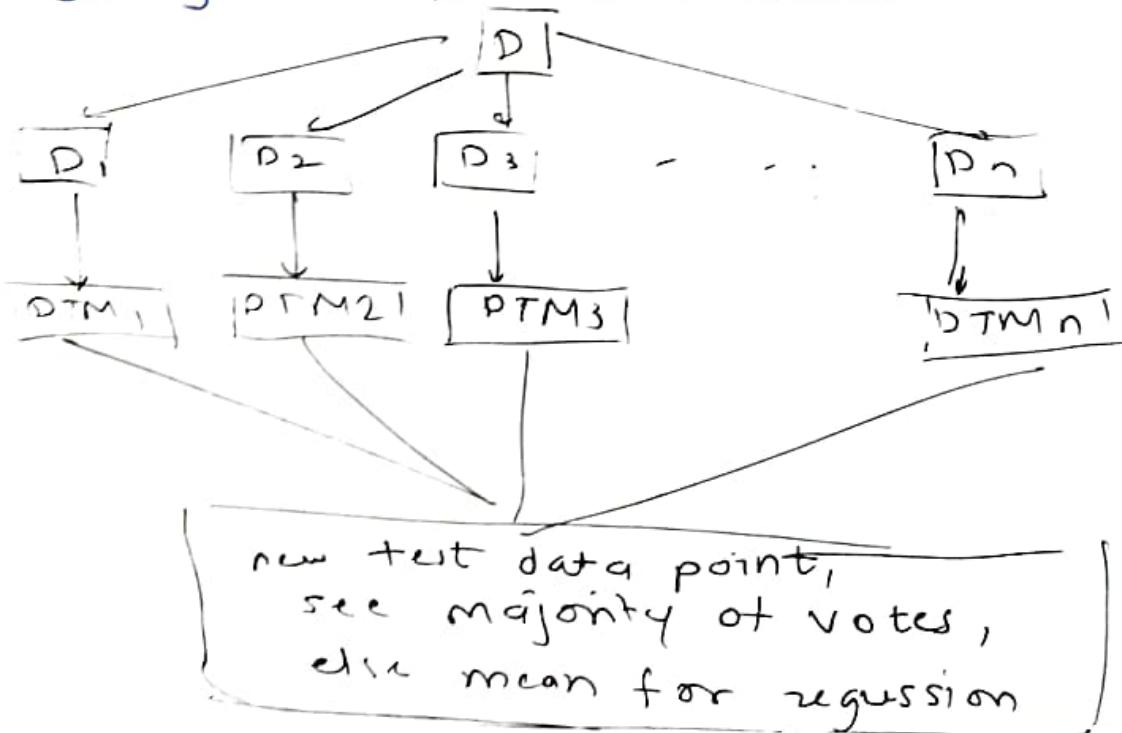
but if we dont change settings.

collection of trees of a bagging technique.

A group of trees. so forest.

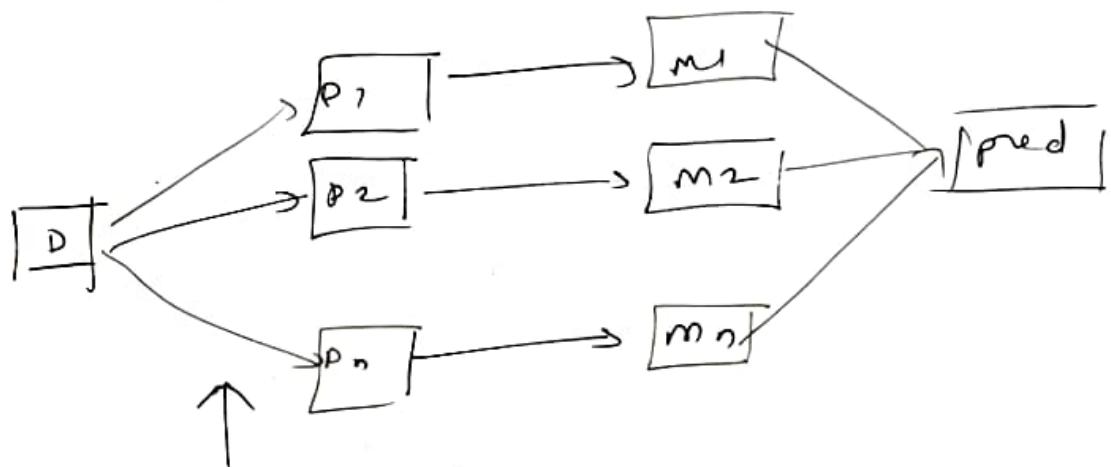
random mean for random sampling via bagging -

bagging technique with trees.



low bias high var	High bias low var
fully grown due tree	LR
SVM	
KNN	

random forest helps us to maintain, like it convert low bias high var algorithm to low bias low var algorithm.



Let say, we changed 1% of our data by outliers & noisy points.
It is not possible to get all noisy points to a single DT. Impact distributed.
so, no much variance

Hyperparameters

number of estimators ~~too~~

Trees to be made

max - features

number of columns considered

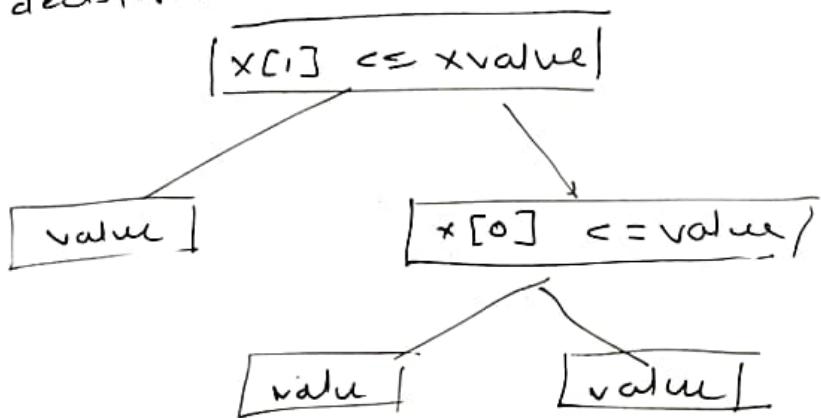
Bootstrap

will you select rows / cols random

max samples

rows

feature importance
in decision tree



So, here, we need to calculate feature importance.

So, for a column, we need to find how many nodes were made due to that column.

So,

$$\text{Feature Imp} = \frac{\text{No of nodes made due to that column}}{\text{Total number of Impurity (all nodes)}}$$

How to be node importance

$$n_i = \frac{N-t}{N} \left[\text{impurity} - \left(\frac{N-t-r}{N-t} \cdot \text{right-node_imp} + \frac{N-t-L}{N-t} \cdot \text{left-node_imp} \right) \right]$$

$N-t$ = no of rows to that node

N = total no of rows

impurity,

$N-t-r$ = no of rows in right node

right impurity = right node impurity

$N-t-L$ = no of rows in left node

left impurity = left node impurity

So, decision tree will calculate importance of every feature.

In random forest

for n no of models, we have n number of feature importances for every feature. If random forest then will average it.

AdaBoost

Weak learners - A model with less accuracy, around 50%.

Decision stumps - A type of weak learners which max depth is 1 as a decision tree.

+1 and -1 class in binary classification instead of 1 and 0.

A split on any axis with minimum entropy.

AdaBoost

stage wise } add all one by one
additive method } odd multiple
 } weak learners

+	+	+	-
+	-	-	-
+	+	-	-

+	+	+	-
+	-	-	-
+	+	-	-

miscalclassified. So, tell next to classify them more properly

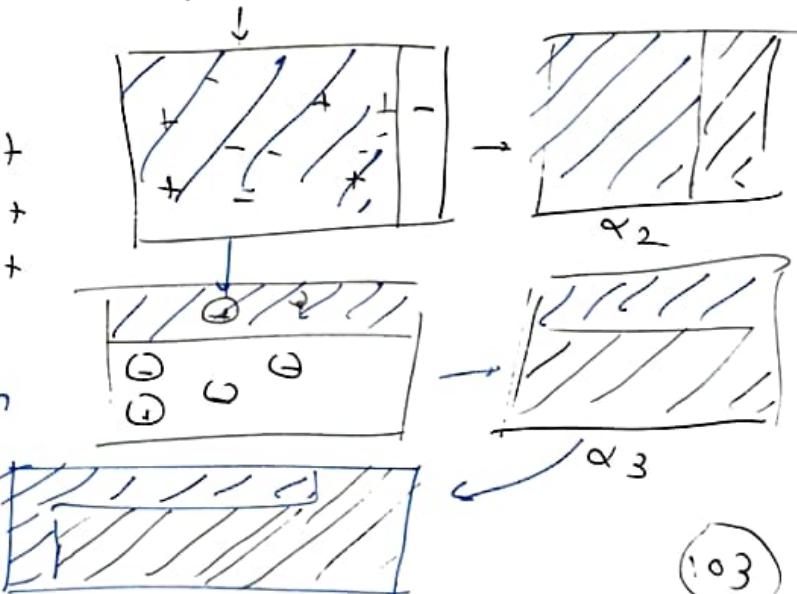
so eqn is

$$h(x) = \text{sign}(\alpha_1 * h_1(x) + \alpha_2 * h_2(x) + \alpha_3 * h_3(x) +)$$

Hx = AdaBoost eqn

α_i = weights assign

h_i = Indiv eqn at modus.



Core mathematics
A classification problem

$$n = 5$$

x_1	x_2	y	wt	Y-pred
3	7	1	0.2	1
2	9	0	0.2	0
1	4	1	0.2	0
9	3	0	0.2	0
5	7	0	0.2	0

In starting we will assign equal weight of 1/5.

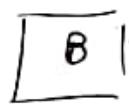
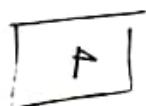
Step 1) Decide a 1-depth decision stump
cut/split on any column, name it
as model 1.

α_1 = for model 1 (weight)

α depend on model error rate

$$\text{error rate } \propto \frac{1}{\alpha}$$

How to calculate α ?



error
rate

0%

50%

100%

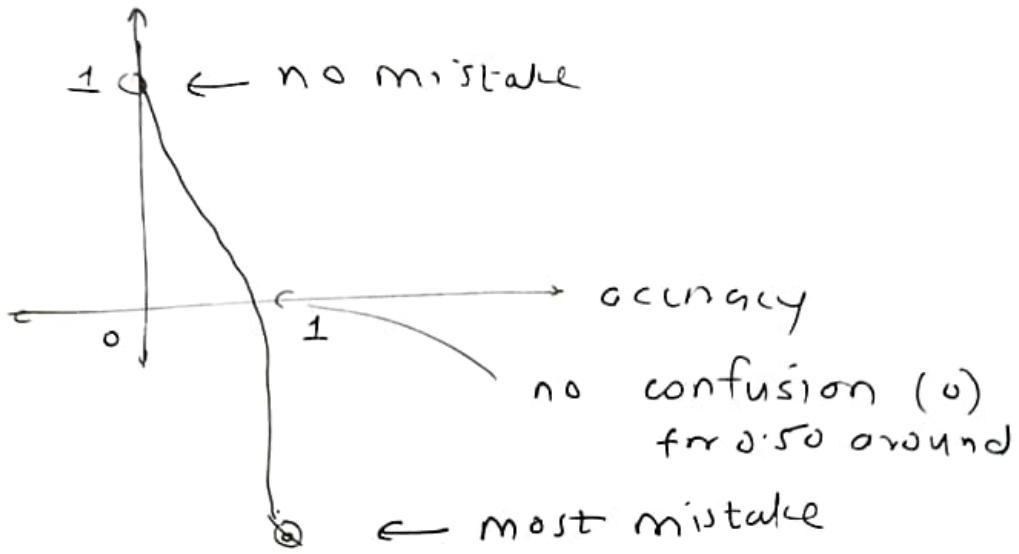
which is most reliable?

model A

except model B & C

model C

because if we inverse results,
we are accurate



This function exist, & that is,

$$\alpha = \frac{1}{2} \log_e \left(\frac{1 - \text{err}}{\text{err}} \right)$$

so, when model M1 predicted

$$\text{error was } \frac{2}{5} = 0.4$$

$$\alpha = \frac{1}{2} \log_e \left(\frac{1 - 0.4}{0.4} \right)$$

$$= \frac{1}{2} \log \left(\frac{0.6}{0.4} \right) = \frac{1}{2} \log \left(\frac{3}{2} \right) \\ = 0.20$$

Now, first prediction successful

Now we need to point to next model
that to need to classify that wrong
points properly.

But how algo will know about
that missclassified points?

We will increase weightage of misclass.
points and decrease weightage of classif.

But in which quantity to increase &
decrease?

for misclassified points

$$\text{new-wt} = \text{current-wt} \times e^{\alpha_i}$$

for classified points

$$\text{new-wt} = \text{current-wt} \times e^{-\alpha_i}$$

(105)

org-wt	new-wt		
	1		
0.2	0.16	0.166	0.166
0.2	0.24	0.25	0.166 - 0.710
0.2	0.24	0.25	0.410 - 0.610
0.2	0.16	0.166	0.610 - 0.12
0.2	0.16	0.166	0.162 - 0.710
0.2	0.16	5.100	
	0.96		

boosting misclassified points via upsampling
 Now generate 5 random number,
 and put them in rows range

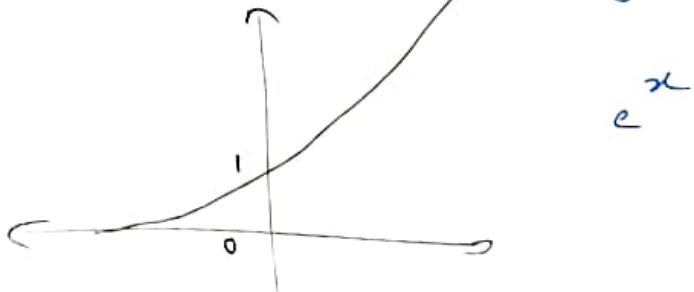
0.13	0.48	0.12	0.50	0.8
1	3	3	3	4

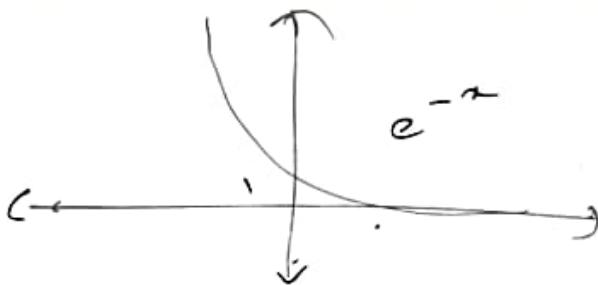
what happens here, that, the probability of "slating" that now, on random number generation, increase on having higher ranges, mean upsampling is occurred

now, I will again train model, and do stumps, again α_2 , & upsampling. do until we have n no of decision stumps.
 so $P = h_1(x)\alpha_1 + h_2(x)\alpha_2 + \dots + h_n(x)\alpha_n$

why e^x or e^{-x} term?

if a model has large α , mean less mistake for this value when model apply, it will a high value. A high trustable model get high value





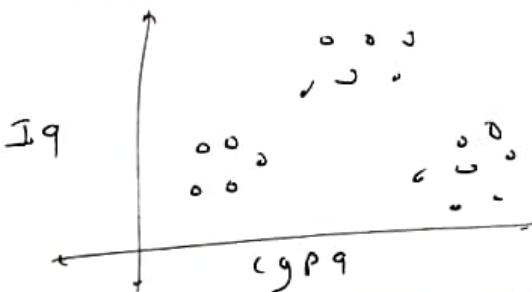
Reducing learning rate in adaboost decreased overfitting.

Bagging vs boosting

	<u>Bagging</u>	<u>Boosting</u>
Type of algorithm	Low Bias + high variance algorithm (fully grown DT)	High Bias + low variance algorithm (Decision stumps 1-length DT)
Learning	Parallel in bagging	Sequential in boosting
Weightage of base learners	same for all models	have some weightage to a model

K-mean clustering

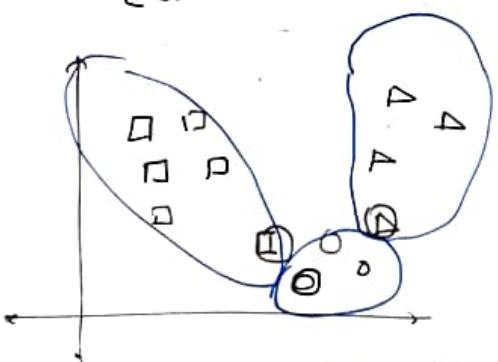
I have a data of our student, I wish to train them for exams, & I wish to classify them.



We wish to make cluster (groups). This is easy in 2D, But not in higher dimensions.

Steps working

- 1) Tell algo number of clusters. (k)
- 2) Initialize randomly, centroids = (k).
- 3) We will make clusters based on that points. So from the centroid, we will calculate every distance of every point, & will assign point to that centroid, if it is nearest one.



After round 1, we calculated the distance of classified points on least distance. Now we have these clusters.

Now for that cluster, calculate means with respect to the points and axes, and assign a new centroid.



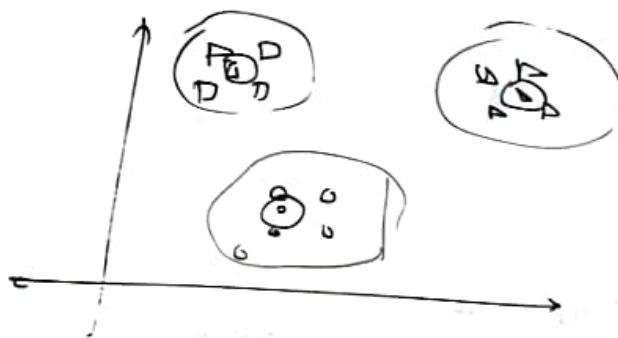
Original centroid
Moved centroid

run algorithm until:

our centroids goes in the groups.

In last, we still don't have achieved, as centroid so far from points.

Now, assign clusters again according to centroids.



still our centroid moved & cluster change, we will again move our centroids by mean.

If centroid moved, we need to repeat again.

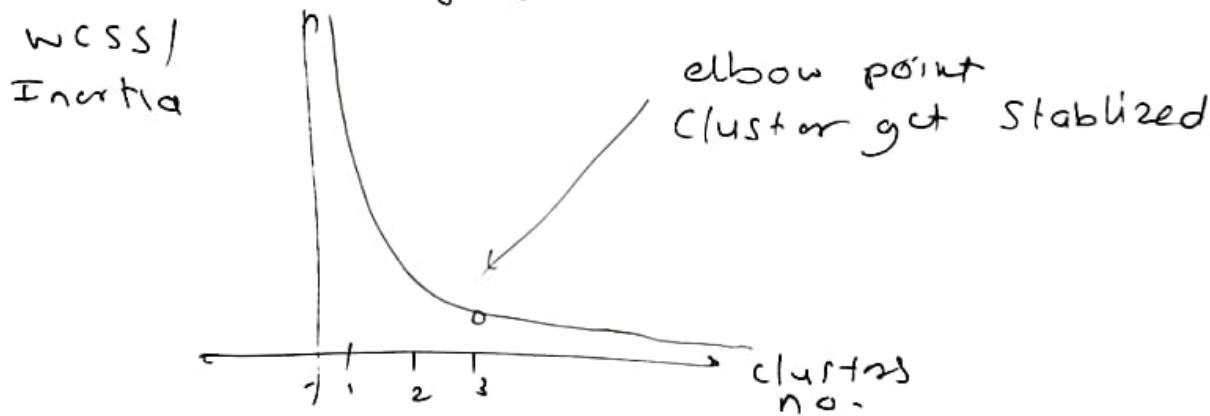
Assign cluster, calculate centroids, and fix centroid.

At the end, our centroids will not move.

How I will know How many clusters to take?

by Elbow method.

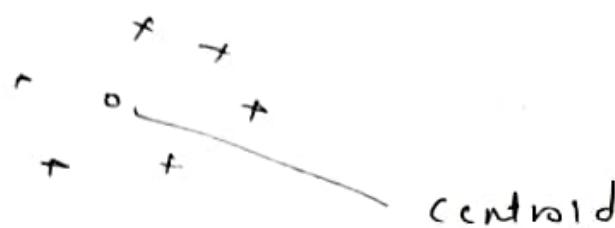
plot graph



WCSS - within cluster sum of square dist

WCSS₁ > WCSS₂ > WCSS₃ - - -

wcss for one cluster



Calculate, square, add distance b/w them
for two clusters

Calculate individual clusters, add wcss

Start with cluster 1, 2 ... 4 plot on the graph of elbow.

euclidean distance for calculating distance between centroid & points.

two points, $(x_1, y_1), (x_2, y_2)$

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

but higher dim,

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2 \dots}$$

but we can't calculate this way

we say,

$$= (x_2 - x_1)^2 + (y_2 - y_1)^2$$

$$= (x_2 - x_1) \frac{+}{(x_2 - x_1)} + (y_2 - y_1) (y_2 - y_1)$$

for the vector,

$$= [(x_2 - x_1)(y_2 - y_1)] [(x_2 - x_1) + (y_2 - y_1)]$$

we say $(x, y) = a$

$(x_2, y_2) = b$

$$= (b - a) = [(x_2 - x_1)(y_2 - y_1)]$$

so, thus way, Entire thing is

$n \cdot \sqrt{a \cdot a}$

$$n \cdot \text{dot}(a - b, a - b)$$

)

Gradient boosting

Again a boosting technique

cgpa	salary	pred_{M_1}	$\boxed{M_1} = \frac{3+4+6+6+3}{5} = 4.8$
90	8	4.8	res_1 in regression,
100	7	4.8	-1.8 It's mean of output.
140	6	4.8	-0.8
120	6	4.8	3.2
80	5	4.8	1.2
			-1.8

loss function = Actual - predicted
(pseudo residual)

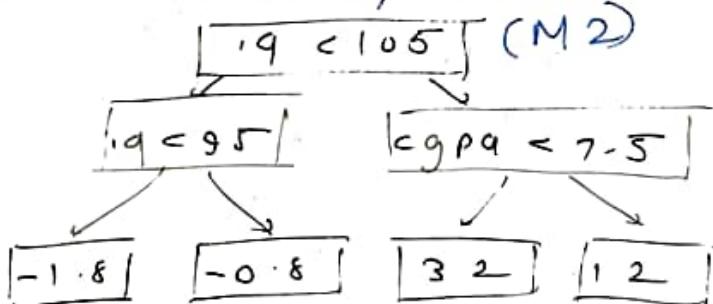
We want to tell me, the mistakes of M_1 , anyway.

$\boxed{M_2}$ a decision tree model (assumed)

Here come concept of gradient boosting for model 2,

Input - same OUTPUT - RES1.

We wish to let model 2 predict how many mistakes are done by model 1.



cgpa	salary	$\text{res}_1(M_1)$	$(\text{pred}_M - M_2)$
90	8	-1.8	-1.8
100	7	-0.8	-0.8
140	6	3.2	3.2
120	6	1.2	1.2
80	5	-0.8	-1.8

At this moment,
let assume we have 2 models.

Now, formula.

$$\text{pred_avr} = \overbrace{\text{out_M}_1 + \text{out_M}_2}^{\text{for student 3}},$$
$$= \frac{4.8}{(M_1)} + \frac{1.2}{(M_2)} = \text{Salary}$$

But, this is seriously an overfitting.
So, we use learning rate

$$\text{pred_avr} = \text{out_M}_1 + (\text{learning_rate}) \text{out_M}_2$$

We only take it to reduce overfitting,
and it is usually 0.1.

So, on prediction of
same student

$$4.8 + (0.1)(0.1) = 4.8 - 0.12$$
$$= \cancel{4.68} \quad 4.92$$

absolutely wrong,
but if we add more models, IT
will add and add, and go closer to
actual TO reduce overfitting.

So, $\text{res}_2 = \text{actual} - (m_1 + (0.1)m_2)$

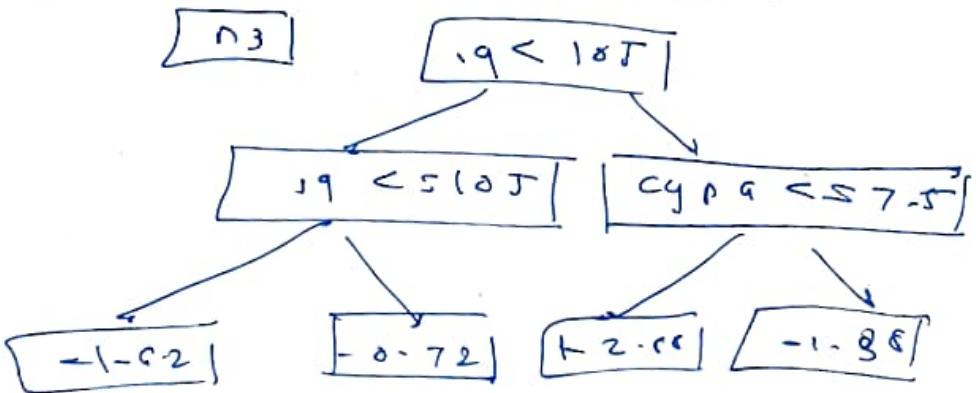
$$3 - (4.8 + 0.1 \times (-1.8)) = -1.62$$

So, Now, calculate res_2 ,

You notice AS you go to add models
further, you comes closer to your ideal
residual (0).

Now, let us train M_3 , as a DT.

Input same OUTPUT res_2 .



so finally,

$$M_{pred} = m_1 + (0.1)m_2 + (0.1)m_3$$

	cgpa	salary	pred1	res1	pred2	res2	pred3
19	5	3	4.8	-1.8	-1.8	-1.82	
20	7	4	4.8	-0.8	0.8	0.72	
100	6	8	4.8	3.2	3.2	2.88	
110	9	6	4.8	1.2	1.2	1.08	
120	5	3	4.8	-1.8	-1.8	-1.82	
BSO							

In gradboost, we usually take max leaf nodes between 8-32 for best results.

gradboost vs gradient-duct

	adaboost	grad-boost
max-leaf train concept	1 (stump) weight for different models.	(8-32) same learning rate for all models.

mathematics behind gradient-boosting algo.
Every Boosting algorithm use additive modeling.

Let's discuss what ML does exactly?

what are ML models?
 $y = f(x)$ - A function

look at dataset

uncorr \Rightarrow don't work

polynomial problem

Runge's phenomenon

Basically it disturb from line in higher degree

is this any direct func?

No

a comb o? Yes $5 \sin x + x$

breaking down func / phasing



13

Algorithm?

We need input
 $\{(x_i, y_i)\}$

We need a loss function, differentiable.

$$L(\hat{y}, f(x))$$

We take a math's loss fun'

$$L(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$$

\leftarrow + handle diff

No of iterations N

For i in range N \rightarrow Step 2.

$$f(x) = f_0(x) + f_1(x) + f_2(x) + \dots + f_n(x)$$

Step 1 find $f_0(x)$

$$f_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \hat{y}_i) \quad L = \frac{1}{2} (y_i - \hat{y}_i)^2$$

$$f_0(x) = \arg \min_{\gamma} \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y})^2$$

We need γ so the value of this expression should be minimum.

$$\frac{d}{d\gamma} f_0(x) = \frac{d}{d\gamma} \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y})^2$$

$$= \frac{1}{2} \sum_{i=1}^n \frac{d}{d\gamma} (y_i - \hat{y})^2$$

$$= \sum_{i=1}^n (y_i - \hat{y}) \frac{d}{d\gamma} (y_i - \hat{y})$$

$$= - \sum_{i=1}^n (y_i - \hat{y}) = 0$$

$$\sum_{i=1}^n (\hat{y} - y_i) = 0$$

So, let's make a dataset now

Ind	R&D	admin str	Profit	Marketing
0	165	137	192	572
1	101	92	144	250
3	29	127	91	201

so, +

$$\sum_{i=1}^3 (r - y_i) = 0$$

$$(r - 192) + (r - 144) + (r - 91) = 0$$

$$3r = 192 + 144 + 91$$

$$3r = 426$$

$$r = 142$$

kind of this is a mean!

step 2)

$$f(x) = \frac{f(x_0) + f(x_1) + f(x_2) + \dots + f(x_m)}{\text{m}} \quad \text{calc in loop}$$

for $i = 1 \text{ to } M$:

$$A) x_{im} = -\left[\frac{\partial L(y_i, f(x))}{\partial f(x_i)} \right] \quad f = f_{m-1}$$

i = row number

m = decision tree number

+ve first, $m = 1$

$$x_{i1} = -\left[\frac{\partial L(y_1, f(x_1))}{\partial f(x_1)} \right] \quad f = f_0$$

can do

pseudo residual

for

first row, x_{11}

first DT

~~first DT~~ x_{21}

~~first DT~~ x_{31}

+1st DT x_{31} ,

3rd row

less f_{4+}

No need of summation,
as already calculating
for row.

$$x_{11} = -\left[\frac{\partial L(y_i, f_i)}{\partial y_i} \right] \quad f = f_1$$

$$-\left[\frac{\partial}{\partial y_i} \left(\frac{1}{2} (y_i - \hat{y})^2 \right) \right] \quad f = f_0$$

$$= [(y_i - \hat{y}_1)] \quad f = f_0$$

$$= [(y_i - f_0(x_i))] + s \quad f_0 \quad (115)$$

$$x_{11} = y_i - f_0(x_i) \quad \text{---} \quad ①$$

at this the case,
so in

$$\begin{aligned}r_{11} &= 4_1 - f_0(x_1) = 192 - 142 = 50 \\r_{21} &= 4_2 - f_0(x_2) = 144 - 142 = 2 \\r_{31} &= 4_3 - f_0(x_3) = 91 - 142 = -51\end{aligned}$$

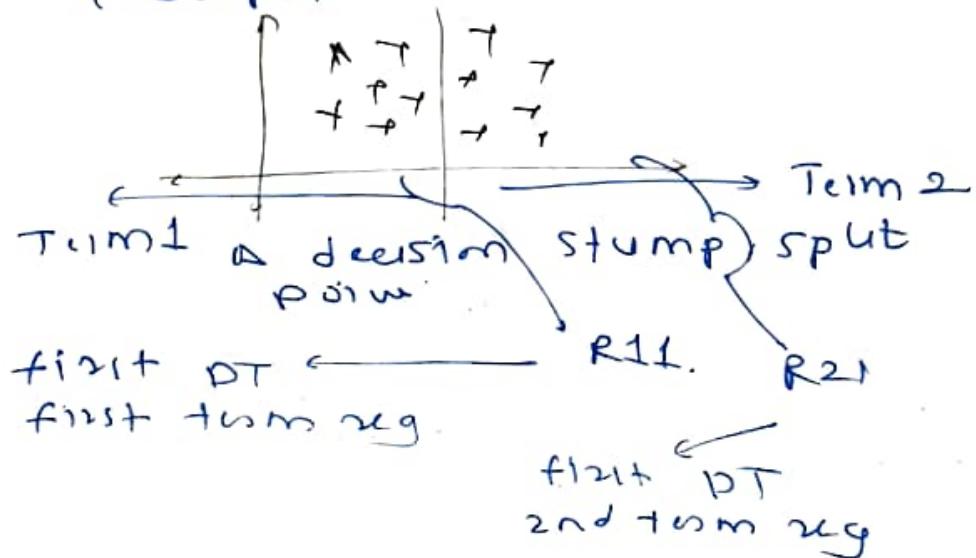
B) fit a regression tree, to targets x_m , using terminal regions.

$$R_jm, j = 1, 2, 3, \dots, m$$

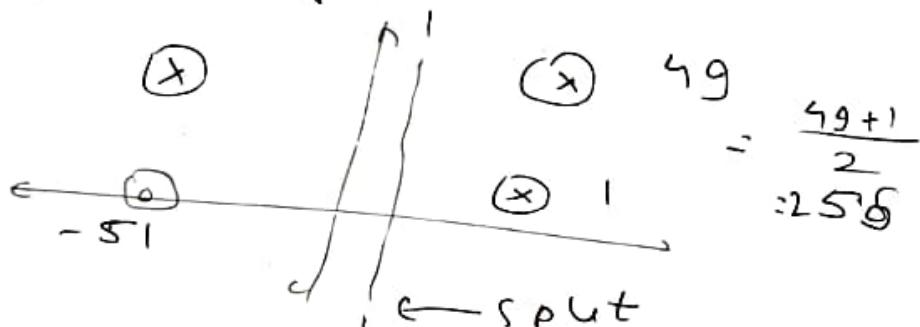
In short, fit a regression tree as data x input & y as output

What is a terminal region?
Answer in DT working.

in S4, A dataset, Input, & Output.



In decision tree, there is a 'value' named value, actually it is output.



⑥ calculate output for every terminal region.

$$r_{jm} = \arg \min_r \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + r)$$

$$r_{j1} = \begin{cases} r_{j1} \\ r_{j2} \end{cases}$$

$$y_{i1} = \arg \min_r \sum_{x_i \in R_{j1}} L(y_i, f_{m-1}(x_i) + r)$$

so, basically means, consider rows fall in that region/leaf only.

$$r_{i1} = \arg \min_r \left(\frac{1}{2} (y_i - (f_0(x_i) + r))^2 \right)$$

$$\frac{dL}{dr} = \frac{1}{2} \cdot 2(y_i - (f_0(x_i) + r))$$

$$\frac{d}{dr} (y_i - f_0(x_i) - r) = 0$$

$$= -(y_i - f_0(x_i) - r) = 0$$

$$= y_i - f_0(x_i) - r = 0 \quad \text{--- (1)}$$

$$r_{i1} = g_1 - l_{g2} - r = 0$$

$$r = g_1 - l_{g2} = -s_1$$

$$r_{i2} = \arg \min_r \left(\sum_{x_i \in R_{i2}} L(y_i, f_0(x_i) + r) \right)$$

$$= \arg \min_r \left(\sum_{i=1}^2 (y_i - f_0(x_i) + r)^2 \right)$$

$$= \sum_{i=1}^2 (y_i - f_0(x_i) - r) = 0$$

$$= \sum_{i=1}^2 (y_i - f_0(x_i) - r) = 0$$

$$= y_1 - f_0(x_1) - r + y_2 - f_0(x_2) - r = 0$$

$$= l_{g2} - l_{g2} + r - l_{g1} - l_{g2} - r = 26$$

(117)

You may note that, this is exactly same output as a leaf of decision tree.

But - This was due to loss function, ordinary square.

Like, we can consider output of DT, as output of OLS.

⑦ update

$$f_{m+1}(x) = f_{m-1}(x) +$$

$$\frac{\sum_{j=1}^M \text{lim}_i l_i(x \in R_{ij})}{\text{DT output}}$$

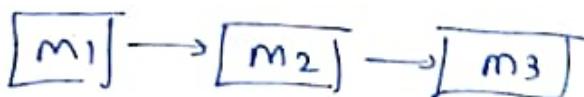
$$f_1(x) = f_0(x) + \text{DT}_1$$

$$f_2(x) = f_1(x) + \text{DT}_2$$

$$f_3(x) = f_2(x) + \text{DT}_3$$

Step 3), output = $\hat{f}(x) = f_m(x)$

Gradient boosting in classification
for examples let us consider 3 models



$$f_0(x) \quad f_1(x) \quad f_2(x)$$

phase 1) build M_1 - a very simple model

$$\log(\text{odds}) = \frac{\text{no}(1)}{\text{no}(0)}$$

This is our first model.

We need to calculate residuals - so, we do need something

But we can't use that fun directly.

Instead, it is giving us probability.

To handle log(odds) and probability, we need to calculate ... from the

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

so for the dataset,

cgrp	Iq	pval + log(odd)	prob(pval)	err
6.82	76	0	0.5108	0.625 - 0.00
6.34	125	1	0.5108	0.625 - 0.375
5.39	99	1	0.5108	0.625 - 0.375
5.50	106	1	0.5108	0.625 - 0.375
6.39	148	0	0.5108	0.625 - 0.125
9.13	158	1	0.5108	0.625 - 0.375
7.17	147	1	0.5108	0.625 - 0.375
7.72	72	0	0.5108	0.625 - 0.125

so, now we will set up a threshold of 0.5, like if $0.5 < P \Rightarrow 1$, else 0.

Here by M1, all are placed.

phase 2 M2

calculate errors of M1.

$$\text{pseudo-residual} = \hat{y} - f_0(x)$$

so, zero col is as follows in table

Now, for model 2, train a DT with same input but pval as output.

So train a model decision tree regressor

Hence, In the decision tree, we will get a value we used to put it directly in regression.

But now, this is probability. To predict,

$$f_0(u) + f_1(u), \text{ i.e}$$

$$\log(\text{odds}) + f_1(u)$$

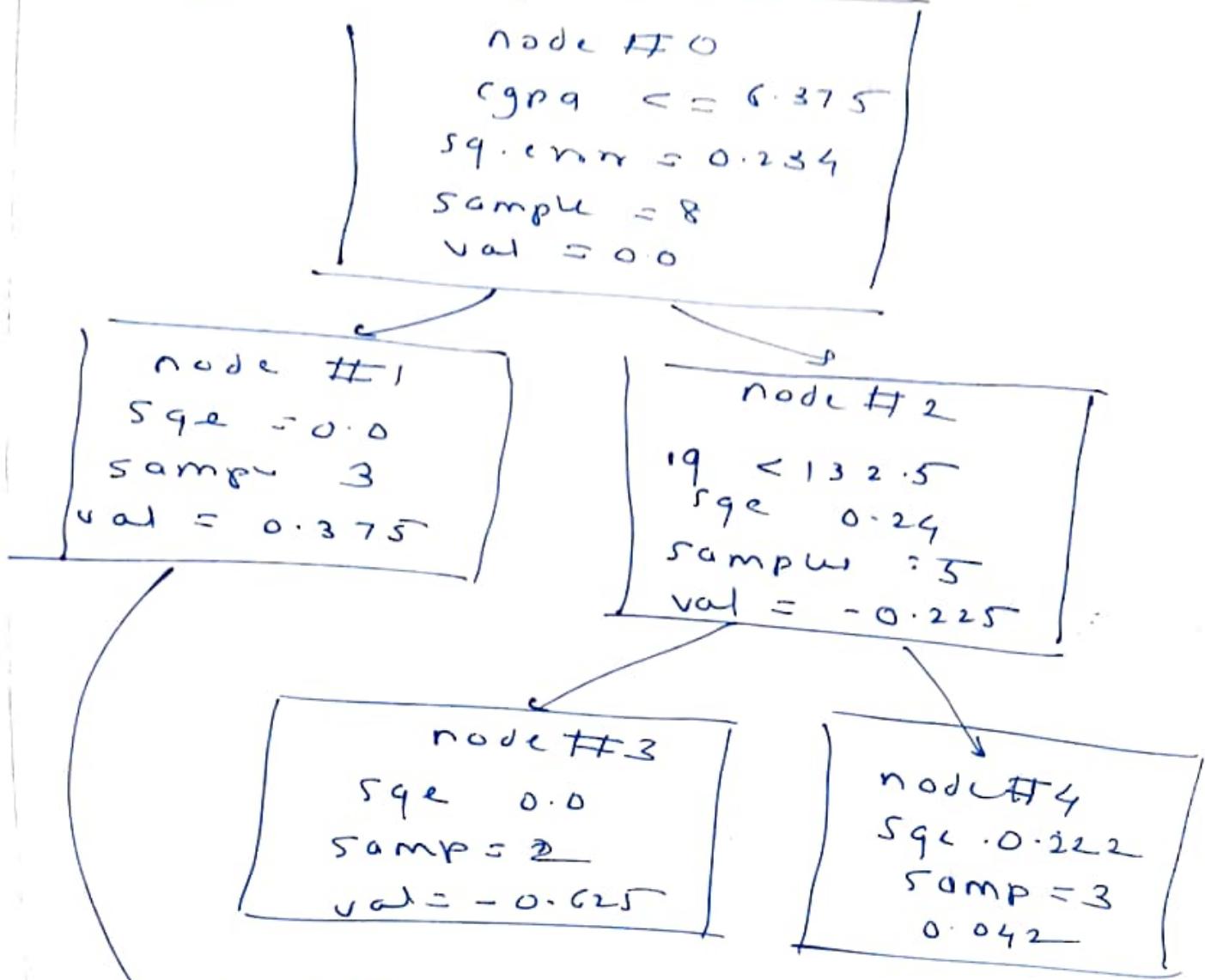
we cannot add probability directly.

So, we need to convert these probs to log(odds).

for the purpose ($\text{prob} \Rightarrow \log(\text{odds})$),
we use the formula for every leaf
node

$$\left[\begin{array}{l} \leq \text{residual} \\ \leq (\text{prev prob}) * (1 - \text{prev prob}) \end{array} \right]$$

so, for right now, my DT is



find for node1, find all points fall in
node1.

sklearn tells us that node
for node 1, rows are [2, 3, 4].

Now, let us calculate
 $\log(\text{odds})$ for node 3,
rows are 1 and 8.

for node #3,
rows are 1 and 8.

$$\begin{aligned}\sum \text{residual} &= -0.625 - 0.625 \\ &= -1.250\end{aligned}$$

$$\begin{aligned}\sum(\text{puv prob}) &= \text{fr}(1) + \text{fr}(r) \\ &= 0.625\end{aligned}$$

so,

$$\begin{aligned}&\frac{-1.250}{(0.625)(1-0.625) + (0.625)(1-0.625)} \\ &= \frac{(-2)(0.625)}{(0.625 \times 0.375) + (0.625 \times 0.375)} \\ &= \frac{(-2)(0.625)}{2(0.625 \times 0.375)} \\ &= \frac{-1}{0.375} = -2.66\end{aligned}$$

so, log(odd fr node #3) = -2.66

for node #1
rows are, (2, 3, 4).

$$\begin{aligned}&\text{so, } \frac{\sum(\text{residual})}{\sum[(\text{puv-prob})(1-\text{puv-prob})]} \\ &= \frac{0.375 + 0.375 + 0.375}{3[(0.625)(1-0.625)]} \\ &= \frac{0.375}{(0.625)(1-0.625)} \\ &= \frac{1}{0.625}\end{aligned}$$

log(oddss(node #3)) = 1.666

for node #4

$$\begin{aligned}&= \frac{-0.625 + 0.375 + 0.375}{3 \times 0.625 + (1-0.625)} \\ &= \frac{0.125}{3 \times 0.625} = \underline{\underline{0.066}} \quad 0.18\end{aligned}$$

(121)

$$\log(\text{odds}(\text{node } \#1)) = -0.1 - 5$$

$$\log(\text{odds}(\text{node } \#3)) = - - 2.65$$

$$\log(\text{odds}(\text{node } \#4)) = 0.18$$

Now,

$$\text{pred_until_now} = M_1 + M_2$$

$$= \log(\text{odds}(M_1)) + \log(\text{odds}(M_2))$$

$$\text{for leaf 3 } (1, 6) = 0.51 + (-2.65) = -2.14$$

$$\text{for leaf 1 } (2, 3, 4) = 0.51 + (1.8) = 2.11$$

$$\text{for leaf 4 } (5, 6, 7) = 0.51 + (0.18) = 0.69$$

Now, we will convert this again to probab, because we need to calculate recall.

$$\text{recall} = 4 - [f_0(x) + f_1(x)]$$

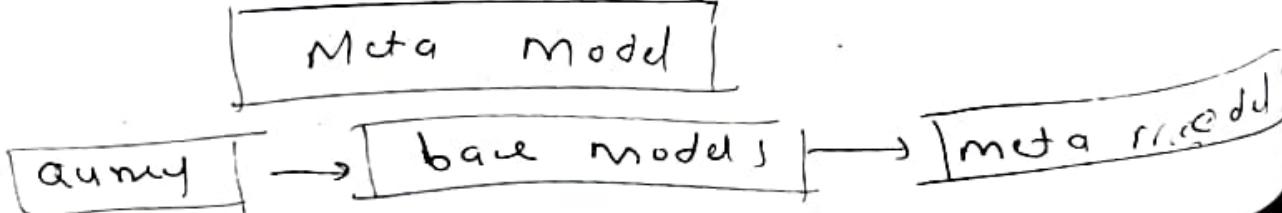
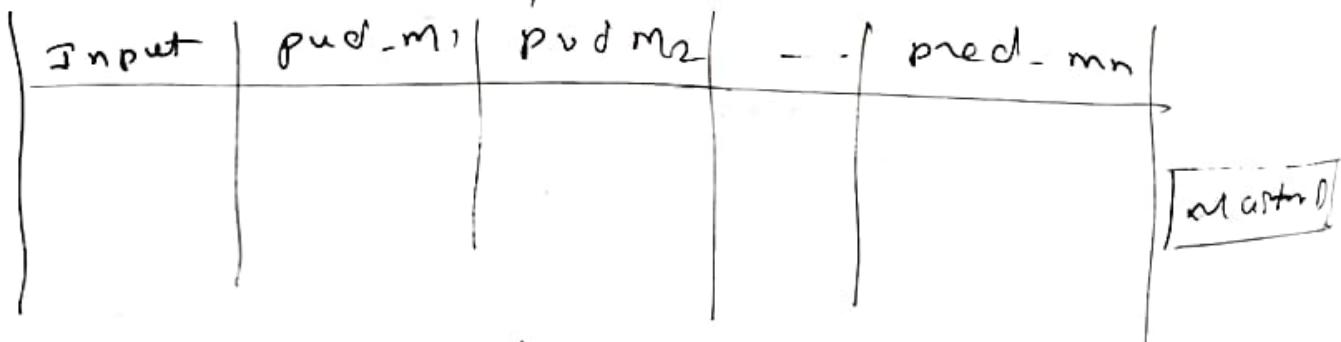
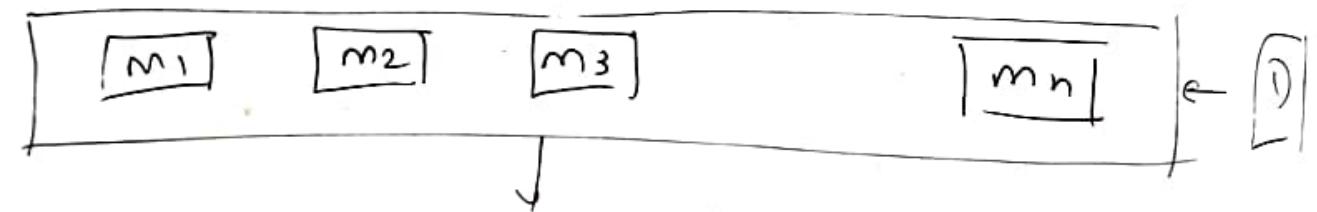
Now, again train DT2 & M3,
as same input but res2.

If these jumps are high, use learning rate.

$$\text{i.e., } \text{rec} = 4 - [f_0(x) - (0.1) f_1(x)]$$

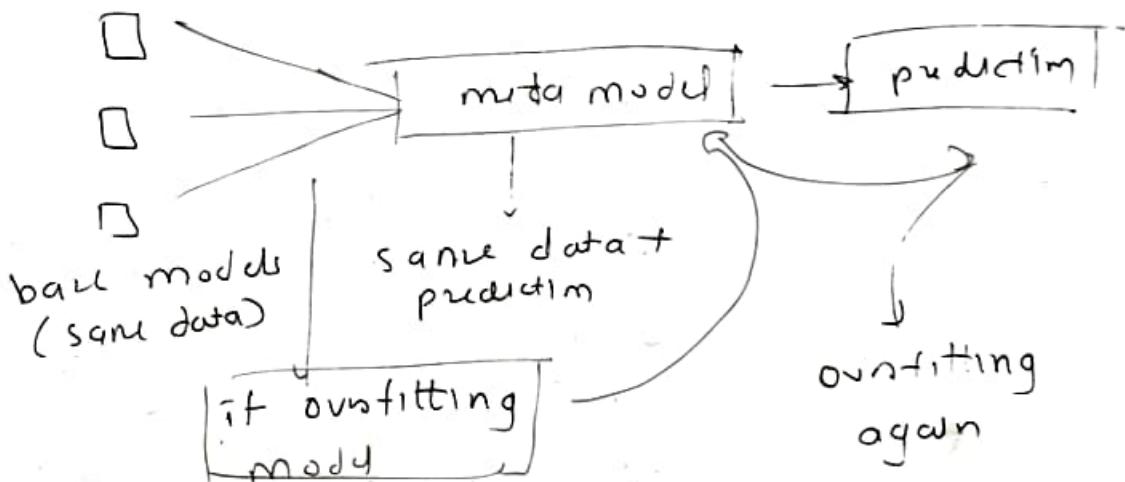
$$\text{pred} = \frac{f_0(x)}{\text{odd log}} + f_1(x) + f_2(x)$$

Stacking



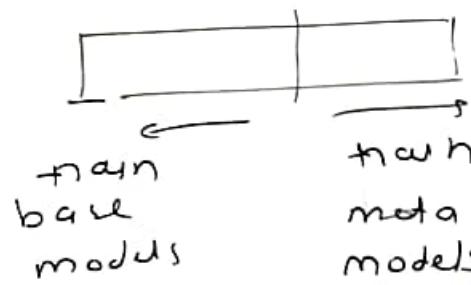
Difference stacking v/s bagging/boosting

	stacking	bagging / boosting
base models	can be difficult	can be same. only
output use	Indirect	direct
problem	△	

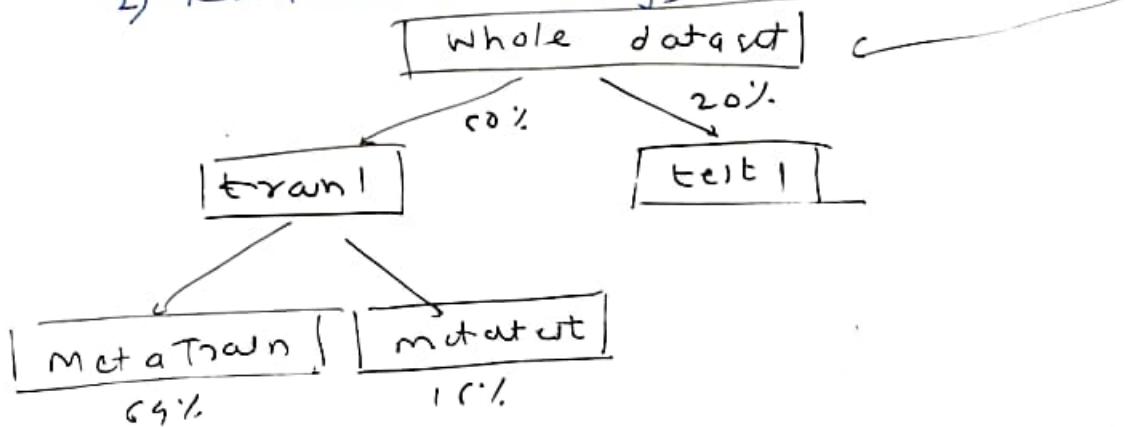


Solutions ✓

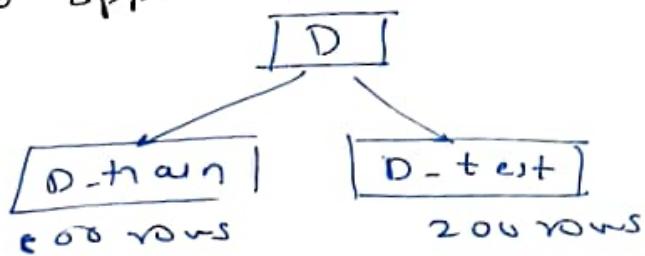
1) Hold out method (Blending)



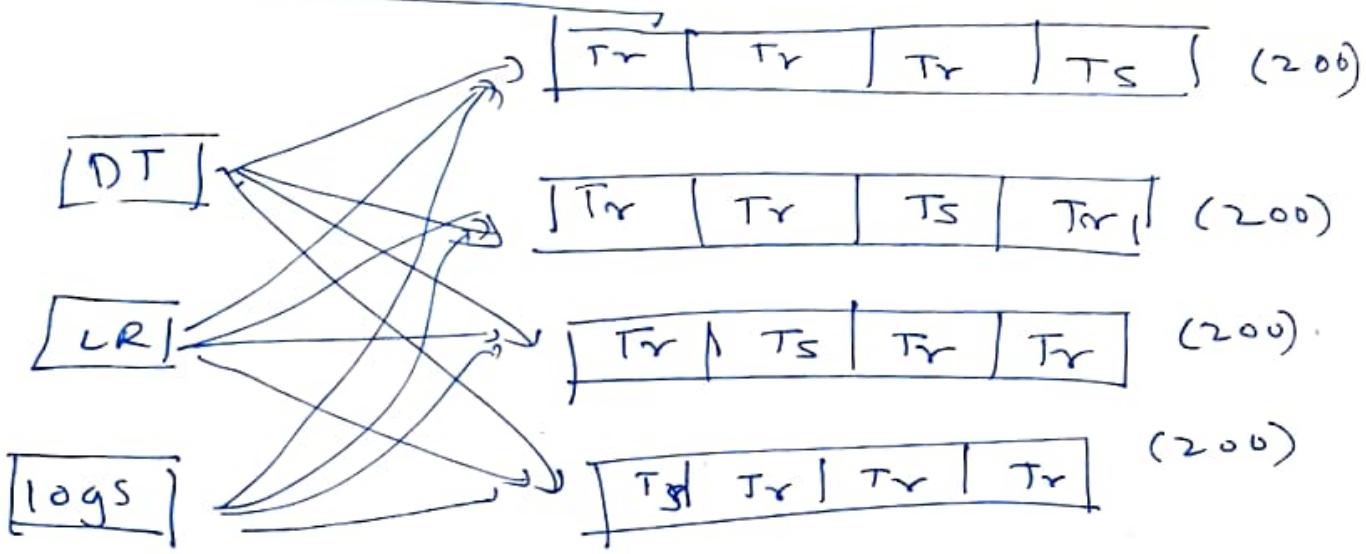
2) k-fold (stacking) → (further)



K-fold approach

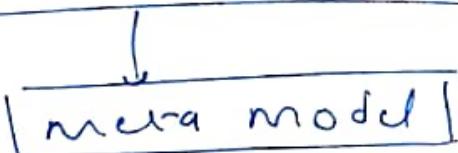


Preferred values for $k = 5/10$
or $k = 4$



	DT	LR	logs	actual
200				
+				
200	-	-	-	-
+				
200				
+				
200				
<hr/>	<hr/>	<hr/>	<hr/>	<hr/>
400				

So, 12 models trained.

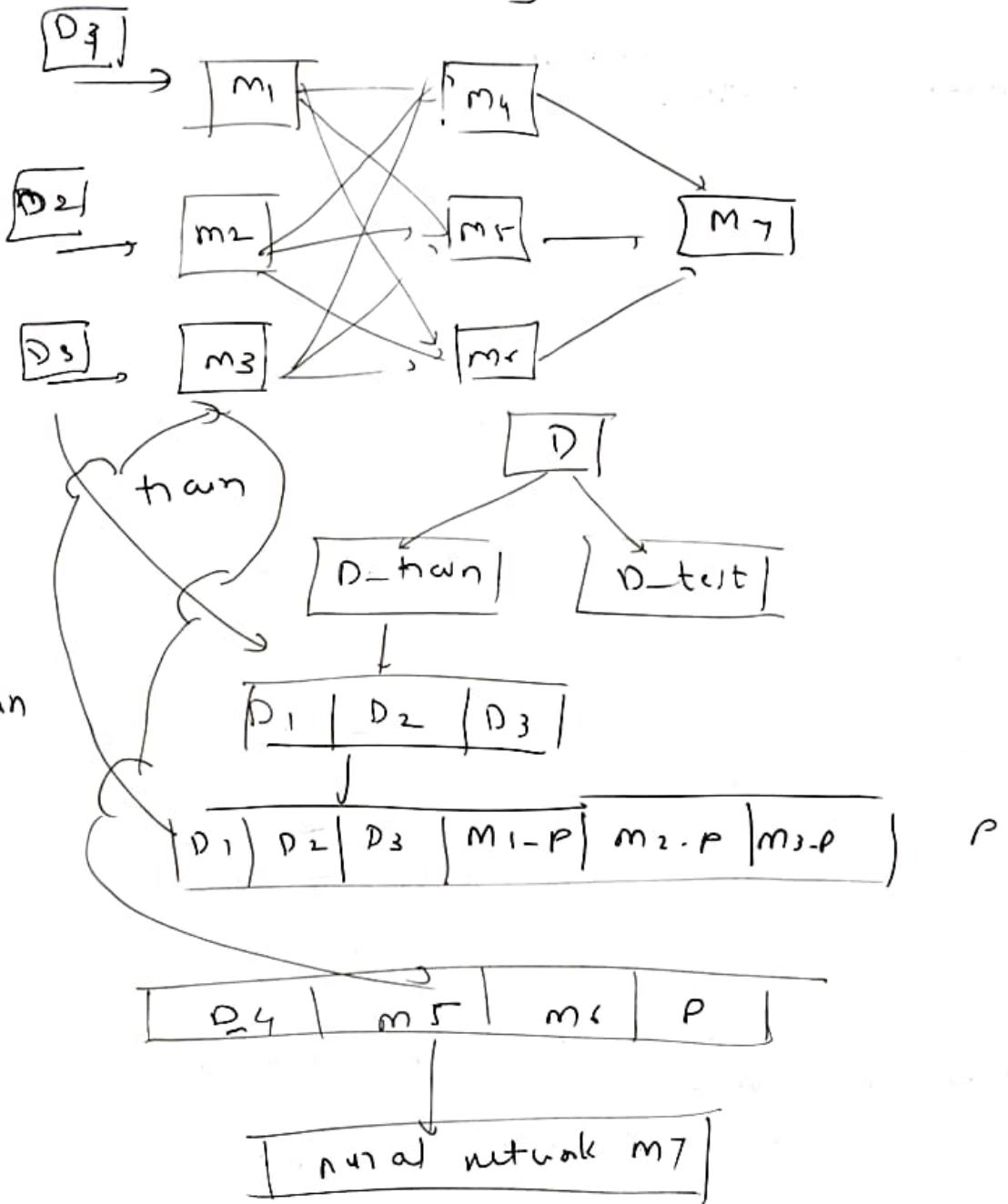


so, But which base model to choose?
 lift whole d-train,
 train LR / DT and Log
 (gettrain)

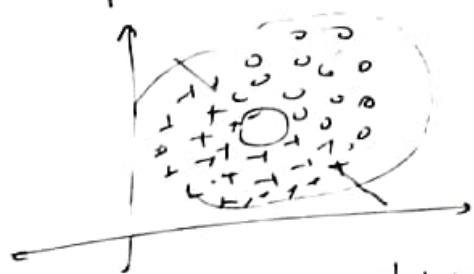
so,



multi-layer stacking



~~Agglomerative hierarchical clustering~~
problem is in centroid working



Types of hierarchical clustering

- ① Agglomerative clustering
Divisive clustering

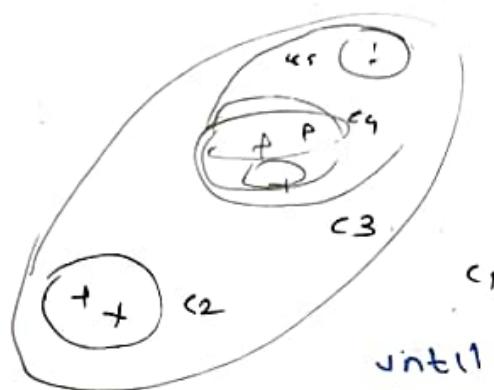


Initially, make assume that each point is an individual cluster.

Then, start merging them, and keep a record.

② Divisive clustering

exact opposite to agglomerative



until every point, becomes a cluster.

But, merging is easy than breaking.

Algorithm

- 1) Initialize proximity matrix

	P_1	P_2	P_3	\dots	P_n
P_1	0	d_1	d_2	d_3	
P_2	d_1	0			
P_3	d_2		0		
P_n	d_3				0

- 2) make each point as a cluster.

- 3) Inside loop

merge → the two closest clusters

update proximity matrix.

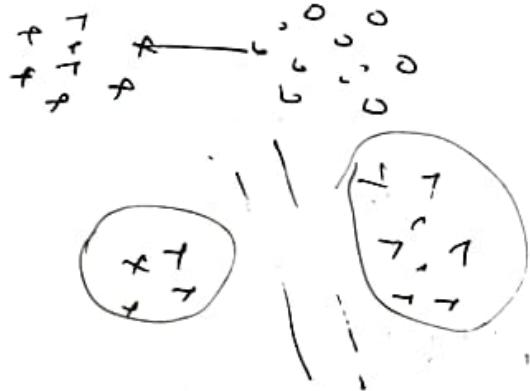
	P_1	P_2	P_3	P_4
P_1	0	-	-	-
P_2	0	-	-	-
P_3	-	-	0	-
P_4	-	-	-	-

calculate
here
again

To find the distances can vary - it depends on what method we use to calculate distances between two clusters.

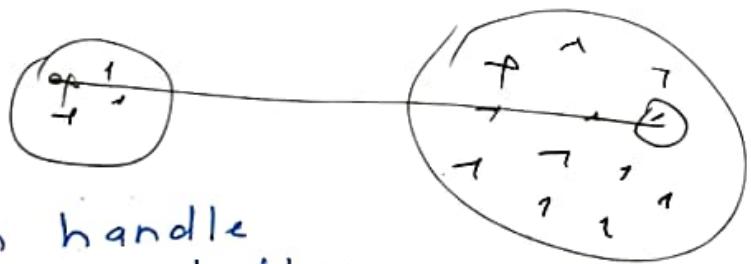
- Types of agglomerative clustering
- 1) min (single link)
 - 2) max (complete link)
 - 3) average
 - 4) Ward

1) min (single link)



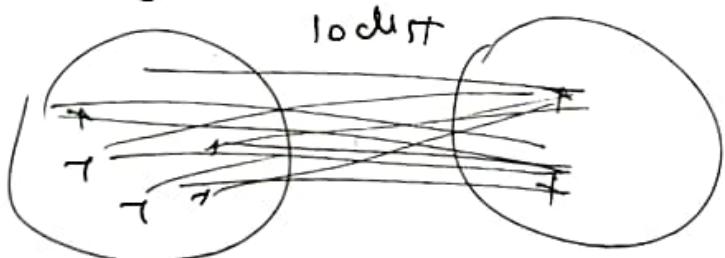
good in logic, can cluster if there is a gap.
But, fails if outliers.

2) complete link



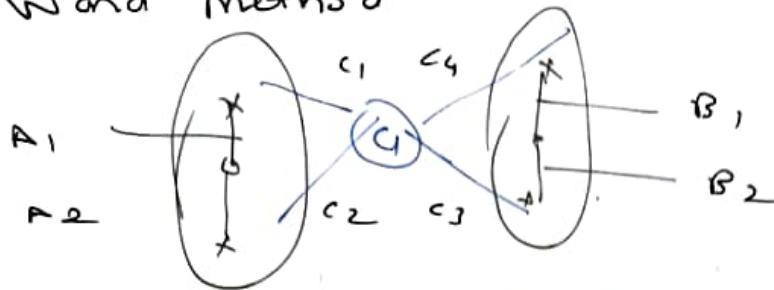
It can handle noise / outliers better.

3) group average



$$\frac{(\text{dist})_{10}}{2} = \text{gap / distance bet^n clusters.}$$

7) Ward method



- 1) plot a centroid.
- 2) calculate distance between points and centroid.

So, Distance is,

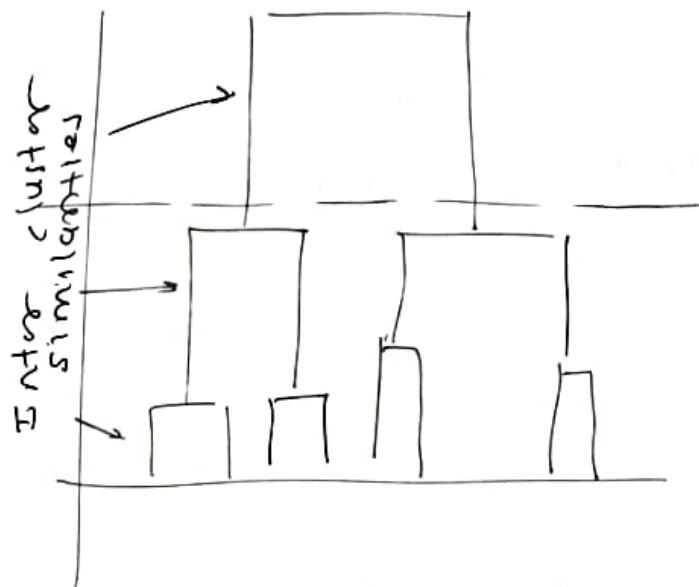
$$d = c_1^2 + c_2^2 + c_3^2 + c_4^2 - (\underline{A_1^2 + A_2^2 + B_1^2 + B_2^2})$$

To minimize the variance

Default in sklearn.

How to decide no of clusters?

plot a dendrogram, and see longest vertical lines.



Hypoparameters (Agglomerative clustering)

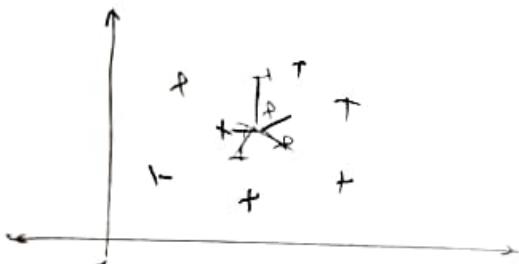
- 1) n-clusters
- 2) affinity, default - 'euclidean'
- 3) linkage = ward, (different types of agglomerative clustering)
- 4) Distance-threshold.

- benefit
- ① widely applicable
 - ② dendrogram
- disadvantage
- ③ cannot apply on large datasets

KNN - K nearest neighbours

You are like your k-neighbors.

see a datapoint,
calculate distances,
sort them, in ascending
take majority for first k.



But how to decide k?

not any fixed answer.

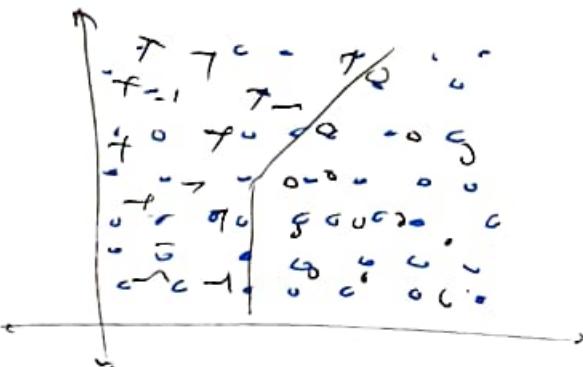
① heuristic approach

② experiment

\sqrt{n} → avoid even

sput data, train diff KNN models for
diff k .

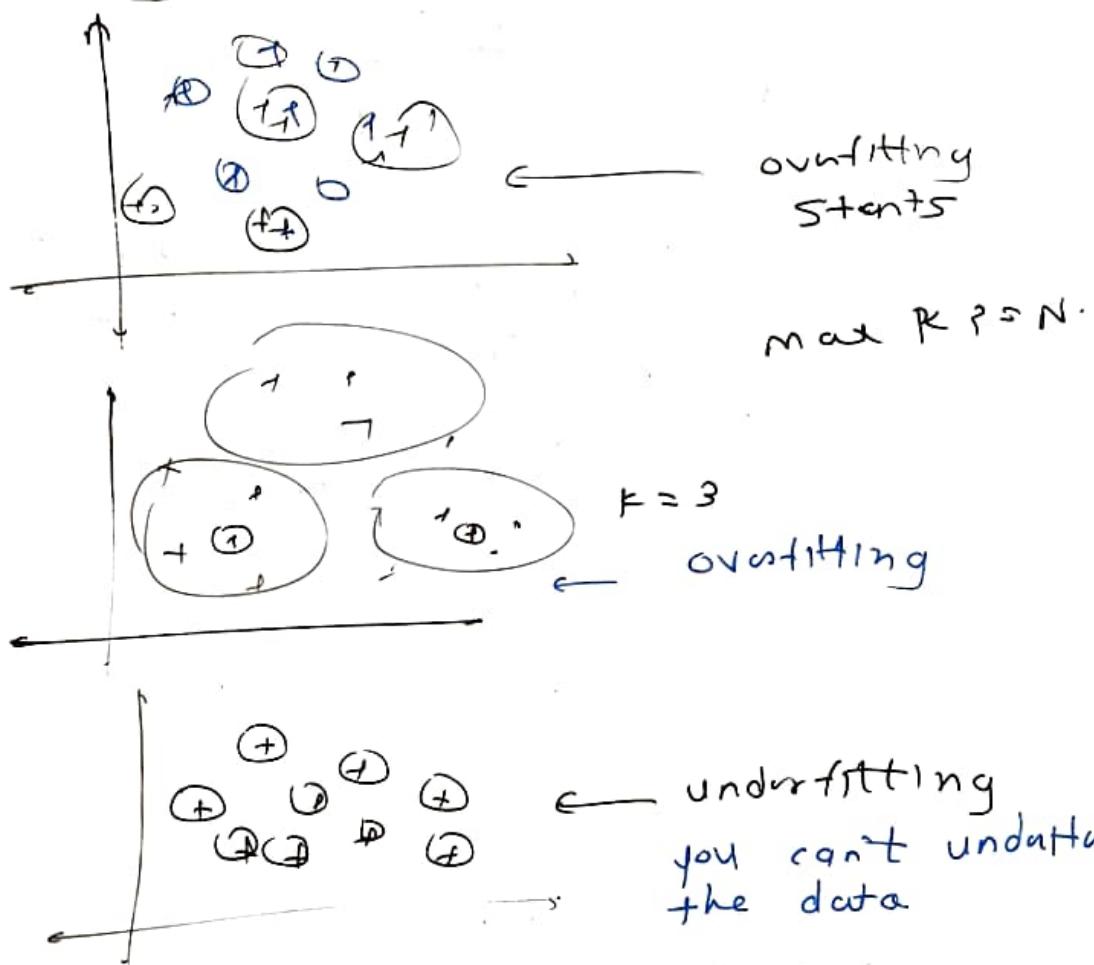
Decision Surface / Boundary



How to build?

- 1) plot training data
- 2) calculate range for axes
- 3) for ranges, generate a numpy matrix.
- 4) Train KNN by these points & send all points in it.
- 5) Encode points by the KNN response
- 6) Show points as a pixel on an image.

Overfitting & Underfitting



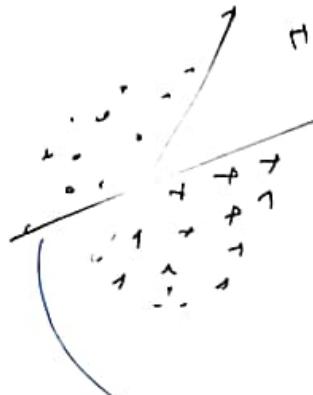
Limitations of KNN

- 1) Large datasets
- 2) Lazy learner,
most of work of in predictions.
In training, we do just store
- 3) In higher dimensions,
cannot see it.

- 4) Imbalanced data
- 5) A black box model - No feature imp

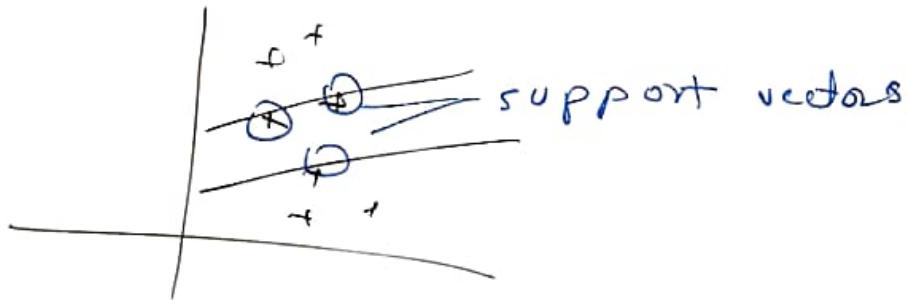
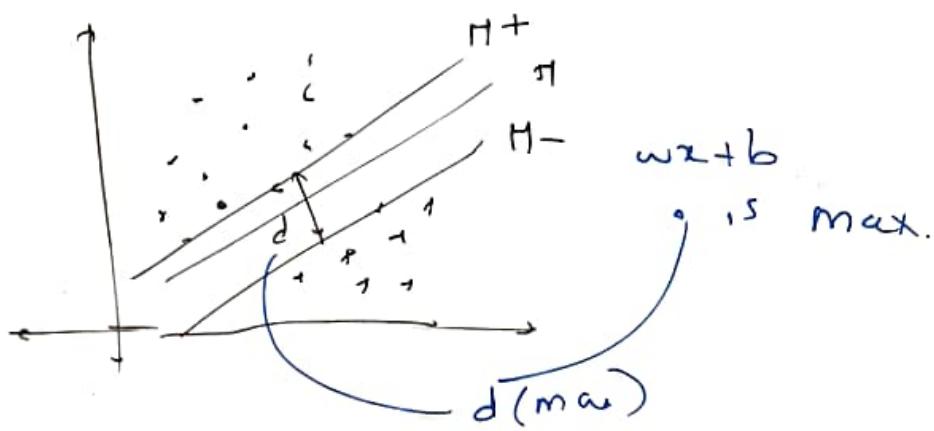
(13)

SVM - support vector machine

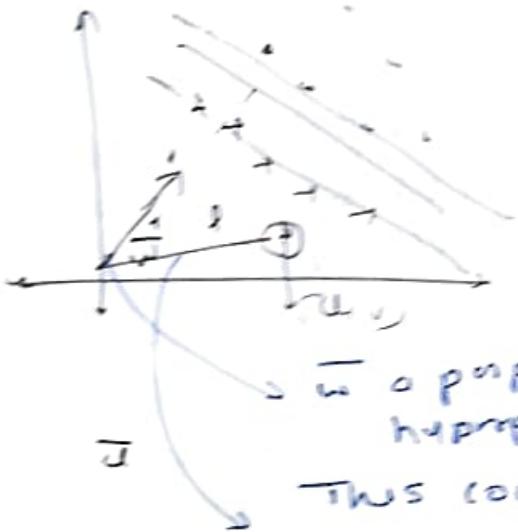


$\|w\| \leftarrow$ but why?
↓ due to high distance
g less chance of overfitting

margin maximizing hyperplane?
how?



- 1) robust to outliers
- 2) non linear data handles
- 3) classification + regression



\hat{w} is perpendicular to all of
hyperplanes

This could also be a vector.

Now, let us define our decision rule
project u due to w .

$\hat{w} \cdot \hat{J} \geq c$ if next (upper) to vector
then positive
if before (lower) to vector
then negative

$$\hat{w} \cdot \hat{J} - c \geq 0$$

$$w \cdot u - c = b$$

$$\hat{w} \cdot \hat{J} - b = 0$$

$$w \cdot u \text{ for } u = (2, 3)$$

$$2x + 3y = 3$$

$$2(2) + 3(3) - 3 \geq 0 \text{ (upper)}$$

$$u = (2, -4)$$

$$= 2(2) + 3(-4) - 3 = 0$$

$$= 4 - 12 - 3 = 0$$

$$= -11 < 0 \text{ (lower)}$$



$$\hat{w} \cdot \hat{J} - b \geq 0$$

$$\hat{w} \cdot \hat{J} - b = 0$$

(ii)

To find the distance d , we need to find the both lines equation.

As for now, let us assume upper vectors equation is $w^T x + b = 1$ and lower vectors equation is $w^T x + b = -1$.

But why? Why this magnitude should be 1 and -1?

Because our (M) should exactly between (H^+) and (H^-) . So equal magnitude $(1, -1)$.

Now, why (1) only? Actually it doesn't matter, can be anything.

If we multiply that $x_i y$ coeff with val > 1 , it will compact, shrink, & if val ≤ 1 , it will expand.

But to calculate loss, ie that boundaries,

The point, either upper or lower for the hyperplane should satisfy the equation.

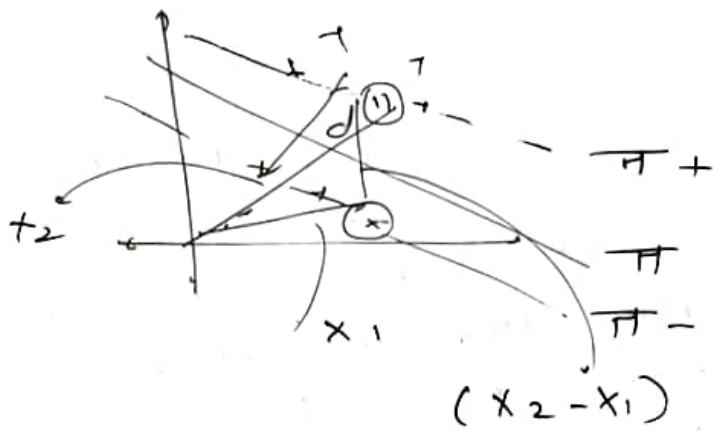
$$y_{i+1} \rightarrow \begin{cases} \bar{w}^T x + b \geq 1 \\ y_{i+1} \end{cases}$$

$$\hookrightarrow \begin{cases} (\bar{w}^T + b)(y_i) \geq 1 \\ (y_i)(\bar{w}^T x + b) \geq 1 \end{cases}$$

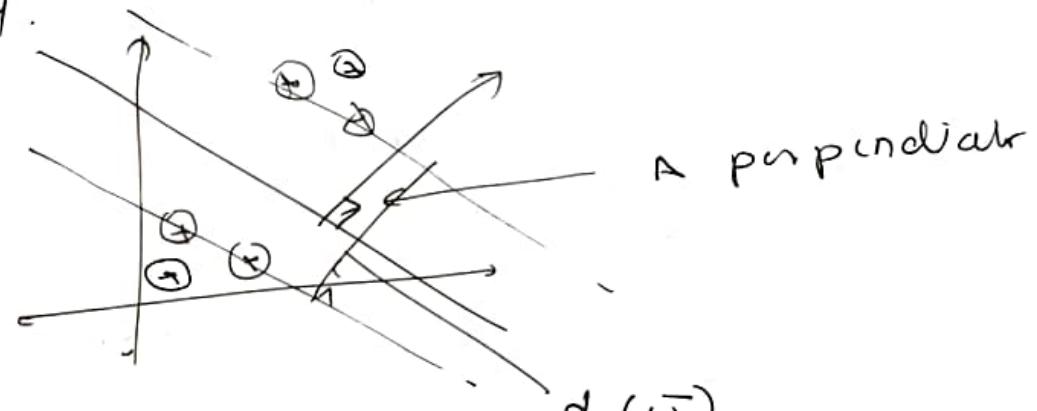
y_i - neighbour

$$\left(\begin{array}{l} y_i(\bar{w}^T x + b) \geq 1 \\ (= \text{for support vws}) \\ y_i(\bar{w}^T x + b) = 1. \end{array} \right)$$

Calculate distance until this is true



But we cannot calculate distance thus way.



$$\begin{aligned} &= \frac{\bar{w}}{\|\bar{w}\|} \\ &\approx d = \left(\frac{\bar{w}}{\|\bar{w}\|} \right) (\bar{x}_2 - \bar{x}_1) \\ &= \frac{\bar{x}_2 \cdot \bar{w} - \bar{x}_1 \cdot \bar{w}}{\|\bar{w}\|} \end{aligned}$$

$$x_1 \cdot (\bar{w} \cdot \bar{x}_1 + b) = 1$$

$$(1) (\bar{w} \cdot \bar{x}_2 + b) = 1$$

$$\bar{w} \cdot \bar{x}_1 = (1 - b)$$

and so

$$(-1) (\bar{w} \cdot \bar{x}_1 + b) \leq 1$$

$$-b - 1 \leq -b - 1$$

$$\rightarrow \frac{(1 - b) - (-b - 1)}{\|\bar{w}\|}$$

$$\therefore \frac{1 + b - b + 1}{\|\bar{w}\|} \approx$$

$$\boxed{d = \frac{2}{\|\bar{w}\|}}$$

135

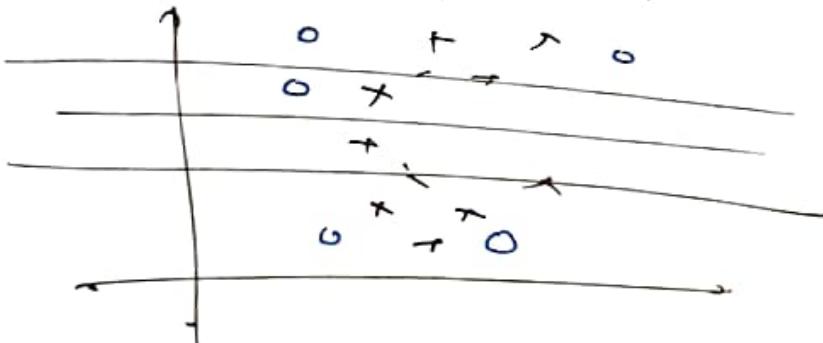
optimization function

$$\underset{(\omega^*, b^*)}{\operatorname{argmax}} \left(\frac{2}{\|\omega\|} \text{ such that } v_i(\bar{\omega}x_i + b) \geq 1 \right)$$

so, all this kind of clear linear separable data, was a hard margin SVM.
But it is possible in real world, we get non-linear data. For that, we need soft margin SVM.

soft margin SVM

make space for outliers!



$$\max(f(x)) = \min\left(\frac{1}{f(x)}\right)$$

$$= \frac{\|\omega\|}{2} + \sum_{i=1}^n \beta_i$$

for all every correctly classified point, β should be equal to 0. (0) error points.

$\beta \rightarrow$ distance from their respective hyperplanes.

Now, I wish to minimize their sum.

$$\begin{aligned} \text{sum error} &= \text{margin} + \text{classification error} \\ &\quad + \text{error} \\ &= \frac{\|\omega\|}{2} + \sum_{i=1}^n \beta_i \end{aligned}$$

if C high: \Leftrightarrow Don't worry about width
worry about misclassif.

algo should not misclass any point

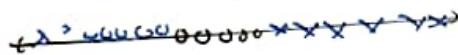
If C is very small, ignore classification error, focus on margin error, & vice versa for higher C .

In logistic also

logistic loss + $\lambda||w||^2$

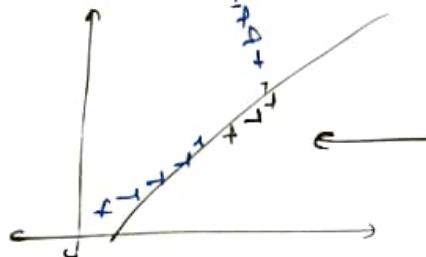
$$C \propto \frac{1}{\lambda}$$

Kernel trick in SVM



We cannot classify this, as it is a non linear data

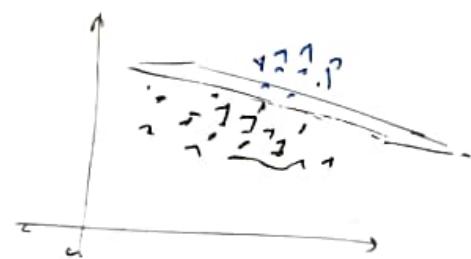
But by some mathematical eqns, we can convert it to higher dimensions



kernel
RBF
polynomial
sigmoid



$$\rightarrow f(x) = e^{-x^2}$$



Naive Bayes
conditional probability

A, B

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{given } P(B) \neq 0$$

Let us consider a dice's problem.

two dice

Sample space $S = \{(1,1), (1,2), \dots, (1,7),$
 $(2,1), (2,2), \dots, (2,7)$
 \dots
 \dots
 $(6,1), (6,2)\}$

	1	2	3	4	5	6
1	2	3	4	5	6	7
2	3	4	5	6	7	8
3	4	5	6	7	8	9
4	5	6	7	8	9	10
5	6	7	8	9	10	11
6	7	8	9	10	11	12

$$\textcircled{1} \quad P(D_1 \leq 5) = \frac{14}{36} = \frac{1}{6}$$

$$\textcircled{2} \quad P(D_1 + D_2 \leq 10) = \frac{33}{36} = \frac{11}{12}$$

\textcircled{3} $D_i \in S$ given $D_1 + D_2 \leq 10$



$$P(D_1 \leq 5 \mid D_1 + D_2 \leq 10) =$$

$$P(\cap B) = \frac{5}{33}$$

$$P(\cap B) = \frac{5}{36} / \frac{33}{36} = \frac{5}{33}$$

Independent events

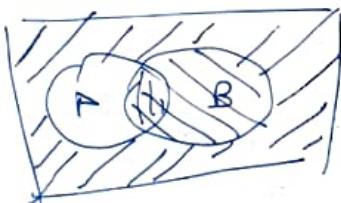
ex.

Two dice rolled,
if die 1 comes 3,
what is prob for
die 2 to come 6?
They are independent,
no matter to other.

$$P(A \cap B) = P(A) + P(B)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \textcircled{1}$$

thus way, $P(B)$ is
already happen.
we are looking
for existence of
 A in that.



$$\begin{aligned} P(A) &= \frac{P(A \cap B)}{P(B)} \\ &= \frac{n(A \cap B)}{\cancel{n(S)} \times \frac{n(B)}{\cancel{n(S)}}} \\ &\underset{P(B)}{=} \frac{n(A \cap B)}{n(B)} \quad \textcircled{2} \end{aligned}$$

$$P(A) = P(A|B)$$

B doesn't matter.
No contribution

Mutually exclusive event

for A, B

$$P(A \cap B) = 0$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B) = 0$$

example,

you are driving car.
what is probability
of turning right
if we took left
already?

like 2 dice, In
first dice, come 3,
what is prob to
come 6 in first
itself?

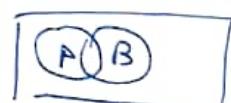
Bayes theorem

likelihood	\rightarrow	prior
$P(A B)$	\times	$P(A)$
↑		\leftarrow
posterior		evidence

$$P(A|B) = \frac{P(A|B) \cdot P(A)}{P(B)} \quad \textcircled{1}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \textcircled{1}$$

$$A \cap B = B \cap A$$



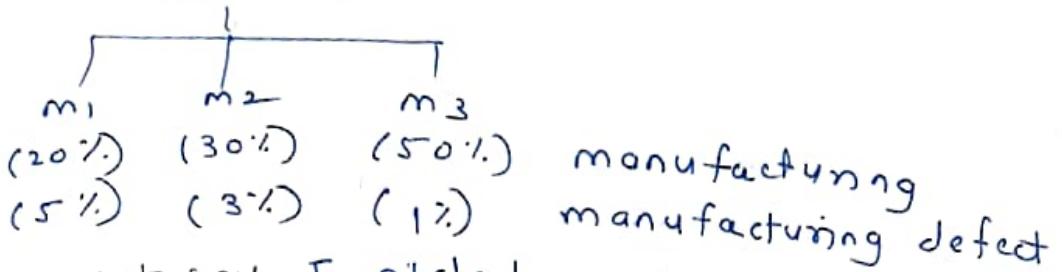
$$P(B|A) = \frac{P(B \cap A)}{P(A)} = \frac{P(A \cap B)}{P(A)}$$

$$P(A \cap B) = P(B|A) \cdot P(A) \quad \textcircled{2}$$

$$P(A \cap B) = P(B|A) \cdot P(A)$$

example by bayes theorem

factory



Now, let say, I picked an object, which was defective from random bag. Now, tell probability it was from m_3 .

$$P(m_1) = \frac{1}{5} \quad P(m_2) = \frac{3}{10} \quad P(m_3) = \frac{1}{2}$$

$$P(D|m_1) = \frac{1}{20} \quad P(D|m_2) = \frac{3}{100} \quad P(D|m_3) = \frac{1}{100}$$

$$P(m_3|D) = \text{prob of } m_3, \text{ given it is defective}$$
$$= P(D|m_3) \times \frac{P(m_3)}{P(D)}$$

$$\left\{ \begin{array}{l} P(D|m_3) = \frac{1}{100} \\ P(m_3) = \frac{1}{2} \end{array} \right.$$

$$P(D) = P(D|m_1) + P(D|m_2)$$
$$+ P(D|m_3)$$

obj being def & come from m_1

$$P(D) = P(D|m_1)(P(m_1)) +$$
$$P(D|m_2)(P(m_2))$$
$$+ P(D|m_3)(P(m_3))$$

$$= \left(\frac{1}{20}\right)\left(\frac{1}{5}\right) + \left(\frac{3}{10}\right)\left(\frac{3}{100}\right) + \left(\frac{1}{2}\right)\left(\frac{1}{100}\right)$$

$$= \frac{\frac{1}{100}}{1000} + \frac{\frac{9}{1000}}{1000} + \frac{\frac{1}{200}}{1000}$$

$$= \frac{10 + 9 + 5}{1000} = \frac{24}{1000}$$

Now,

$$= \frac{\left(\frac{1}{100}\right)\left(\frac{1}{2}\right)}{\left(\frac{24}{1000}\right)} = \left(\frac{1}{200}\right) / \left(\frac{24}{1000}\right) = \frac{\frac{5}{1000}}{\frac{24}{1000}} = \frac{5}{24}$$

Sample example and mathematics

CSK

toss	venue	outlook	result
won	mumbai	Overcast	won
lost	chennai	Sunny	W
w	Kolkata	S	N
w	C	S	W
w	M	S	Lost
L	M	S	L
w	C	O	L
w	K	O	L
w	M	S	W

Now, predict

win loss

by datapoint

{ lost.

Mumbai,

sunny }

Prob(lost \cap mumbai \cap sunny) = ?
by Naive bayes.

In more clarity

$P(\text{winning} | (\text{lost} \cap \text{mumbai} \cap \text{sunny}))$

$P(\text{lost}) | (\text{lost} \cap \text{mumbai} \cap \text{sunny})$

find probab Given condition
will find prob, < greater one will answer.

$$P(A \cap B) = \frac{P(B|A) P(A)}{P(B)}$$

$$P(W | \text{lost}, \text{mumbai}, \text{sunny})$$

$$= \frac{P(\text{lost}, \text{mumbai}, \text{sunny} | W) P(W)}{P(\text{lost}, \text{mumbai}, \text{sunny})}$$

$$P(L | \text{lost}, \text{mumbai}, \text{sunny})$$

$$= \frac{P(\text{lost}, \text{mumbai}, \text{sunny} | L) P(L)}{P(\text{lost}, \text{mumbai}, \text{sunny})}$$

At the end, we will divide greater one, so we don't need denominator, which is same in both cases.

So remove that

$$P(w | \text{lost}, \text{mumbai}, \text{sunny}) =$$

$$P(\text{lost}, \text{mumbai}, \text{sunny} | w) P(w)$$

$$P(L | \text{lost}, \text{mumbai}, \text{sunny}) =$$

$$P(\text{lost}, \text{mumbai}, \text{sunny} | L) P(L) \quad \text{--- (1)}$$

"Naive Bayes", because it applies a very simplistic assumption on bayes theorem.

$$P(L | \text{lost}, \text{mumbai}, \text{sunny}) =$$

$$P(\text{lost} | w) P(\text{mumbai} | w) P(\text{sunny} | w)$$

Just did, because [--- (1)] is so specific
seems to be occur very unlikely,

$$\text{win} = \left(\frac{1}{8} \right) \left(\frac{2}{5} \right) \left(\frac{4}{5} \right) \times \left(\frac{5}{8} \right) = \frac{1}{25} = 0.04$$

$$\text{loss} = \left(\frac{1}{3} \right) \left(\frac{1}{3} \right) \left(\frac{1}{3} \right) \left(\frac{5}{8} \right) = \frac{1}{72} = 0.013$$

Core mathematics

$$x = \{x_1, x_2, \dots, x_n\}$$

$$c. = \{c_1, c_2, \dots, c_n\}$$

$$P(c_k | x) = \frac{P(x | c_k) P(c_k)}{P(x)} \leftarrow \text{no condnb}$$

$$= P(x | c_k) P(c_k)$$

$$P(c_k | x) = P(x | c_k)$$

$$P(A \cap \bar{B}) = \frac{P(A \cap B)}{P(B)}$$

$$P(A \cap \bar{B}) P(B) = P(A \cap \bar{B})$$

$$P(C_k | \bar{x}) = P(x_n | k) \quad \text{---} ①$$

$$P(x_k | x) = P(x, c_k)$$

$$P(\frac{x_1, x_2, x_3, \dots, x_n, c}{A \cap B}) \text{ Now assume}$$

$$P(A \mid B)$$

$$\underbrace{p(x_1 | x_2, x_3, \dots, x_n)}_{a} p(x_1, x_2, x_3, \dots, x_n | c_k)$$

$$P(c_k | \omega) = a \cdot P(x_2, x_3, \dots, x_n | c_k)$$

2001

$$P\left(\frac{x_2, x_3, \dots, x_n < k}{A \quad B}\right)$$

$$= a \cdot \frac{P(x_2 | x_3 x_4 \dots x_n \in k)}{b} P(x_3 x_4 \dots x_n \in l)$$

$$= a \cdot b \cdot P(x_3, x_4, \dots, x_{n+1}, k)$$

chain rule for conditional probability

↓ break down

↓ break down

lost + com =

$$\frac{P(x_n | c_k)}{P(c_k)} \varphi(c_k)$$

$$= p(x_1 | x_2 x_3 x_4 \dots x_{n-k})$$

$$P(x_2 | x_3 x_4 x_5 \dots x_n)_{(k)},$$

$$P(x_3 | x_4 x_5 x_6 \dots x_n(k))$$

$$P(\underline{x_n} \mid x_n(k))$$

this only depends
on Ck. — Assumption

$$P(C_k)$$

↓
Conditional independence

143

We are taking the assumption because, finding that very specific event probability is very very less or even zero.

Independent event

$$P(A|B) = P(A)$$

So further,

$$P(A|B, C) = P(A|C)$$

A is independent on B conditionally

So further simplifying eqn ②

$$= P(x_1|c_k) P(x_2|c_k) \dots P(x_n|c_k) P(c_k)$$

$$\therefore P(c_k|x) = P(c_k) \prod_{i=1}^n P(x_i|c_k)$$

↓
actually

$$\therefore P(c_k|x) \propto P(c_k) \prod_{i=1}^n P(x_i|c_k) \quad \text{--- } ③$$

↓

$$= P(c_k|x) = \frac{1}{Z} \left(P(c_k) \prod_{i=1}^n P(x_i|c_k) \right)$$

$$Z = P(x)$$

Ignore coz it will come same

$$\hat{y} = \arg \max_{k \in \{1, 2, \dots, K\}} \left\{ P(c_k) \prod_{i=1}^n P(x_i|c_k) \right\}$$

↔

maximum a posteriori rule
MAP

Handle numerical data

Height weight gender

172 150 M

find {

180 170 M

$$H = 185, w = 170, G \}$$

165 140 M

190 200 M

$$P(M | H = 165, w = 170)$$

139 100 F

$$= P(H = 165 | M) P(w = 170 | M)$$

$P(M)$

145 120 F

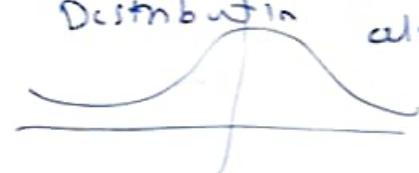
Assumption:

160 190 F

Height is a gaussian

172 150 F

Distribution also weight



find out mean and G.

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$

Gaussian Normal Bayes

use histplot, and find best one fit for
your case

use Gaussian NV,

Binomial NV,

Poisson NV

XGBoost - A rockstar guy

To understand, first, let revise ML history

(70-80) → simple algorithm,
that can only work on a
specific kind of data

(80-90) → some complex
SVM, random forest
powerful
but lead to overfitting
and scalability is not much

Introducing XgBoost

(power, performance, Scalability,
speed).

fact: It is a library!

Tinogi chaini made XgBoost by SE concepts
and gradient boosting.

History

early days
(2010)
grad boost

flexibility

loss funⁿ

performance

kaggle

robust — performance
big data

→ kaggle → open source

power

+
kaggle
work

→ features

optimization

multiple features

kaggle

XgBoost

EXTREME!!

flexibility

speed

performance

1) flexibility

- 1) cross platform (windows, unix)
- 2) multiple language support
- 3) Integration with other libs and tools
- 4) customizable with any ML problem.

2) Speed

1) parallel processing

It is a sequential boosting trick.
To build a model individually,
we use parallel processing.
like In a dataset, when multiple
features, & we will In that
column assume a column, so,
In case of DT, we will found
best splits, by gini and entropy,
and this computation will
individual for any column so
parallel computing.

2) Optimized data structure

others interpret row-wise
It (xgb) is columnwise

3) Cache awareness

On numerical features, Xgboost
calculate bins, we need it regularly,
so we store it in cache memory.
and like this, many things.

4) Out of core computing

splitting data / making chunks
that cannot fit in computers
main memory.

5) distributed computing



perform computing parallel on
different nodes (machines)

- 1) split data
- 2) vocal each data on a single
node
- 3) one master node performs
aggressim & communication
(dask, kubamotes)

6) GPU support

3) Performance

① regularization leaning objective,
use in Uncor modus, to reduce
overfitting

② handling missing values

③ sparsity aware split finding

④ efficient split finding (weighted quantile
search + Approximate tree learning)

⑤ Tree pruning (cut down depth of tree)

regularization term in loss fun? itself.

xgb detects sparsity (very large 0 value)
and handles it. But how?

f

9 - 4.5

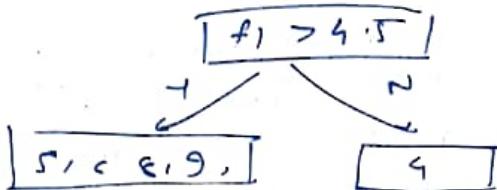
5 - 5.5

6 - 5.5

7 - 7

8 - 8.5

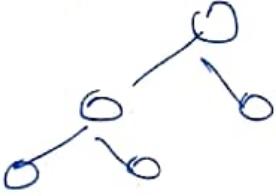
9 - 9



Now, send missing value
in both node & calculate
gain., and go with max
gain node

→ Trying to make sense of depth
of missing value

we train a lot of trees, likely,



f_1 → exact greedy
search, and
splitting criteria

But, it is slow

Instead for the ← approach

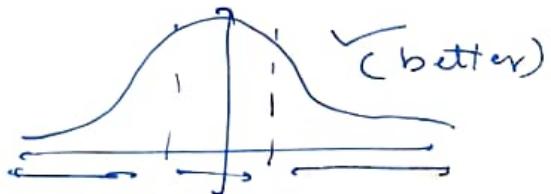
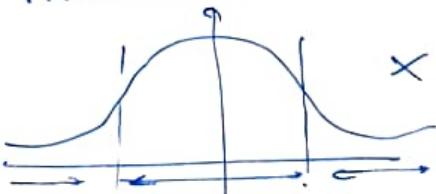
numerical column,

we divide it in bins.

(continuous → discrete),

this is approximate tree learning.
(Histogram based) (weightage
quantile sketch)

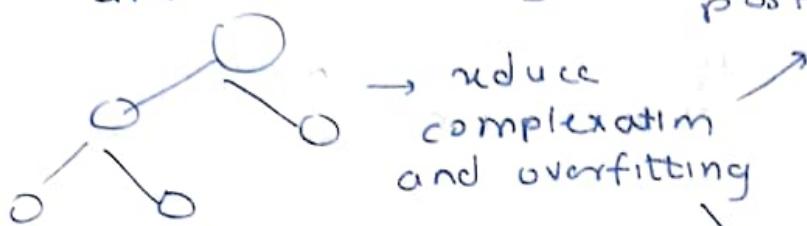
→ we use quantile binning instead of uniform,
to maintain distribution, & decide bins.



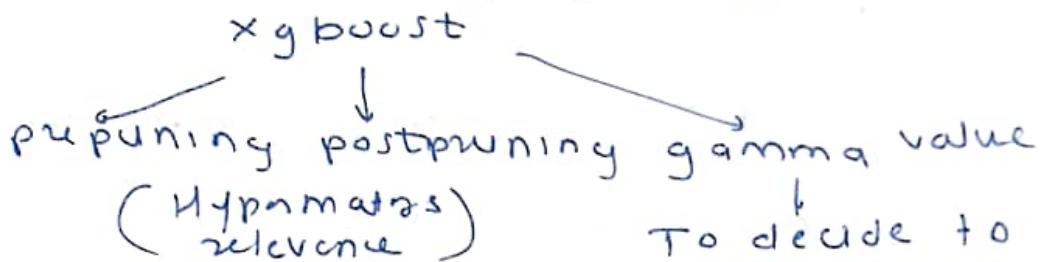
7) tree pruning
To cut down tree for reducing complexity

and overfitting.

post pruning



pre-pruning



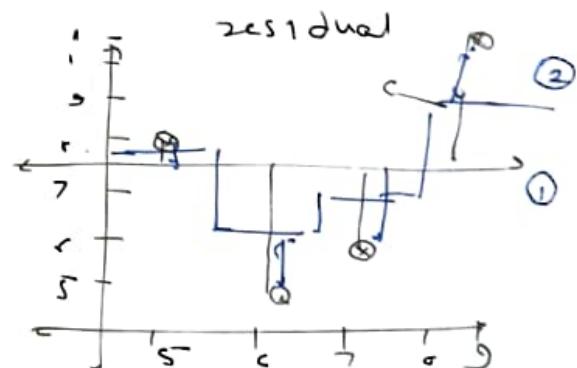
Outro, related to gradient boosting

1) LightGBM, Light GBM

2) catBoost → (Categorical)

for regression

cgpa	package
6.7	45
9	11
7.5	6
5	8



like gradient boosting, start first model as mean. calculate residuals. Now,

train a DT to predict residuals.

calculate ① and, then predict DT,
build model ②. Predict errors of M2.

train M3, same input, & rest of M2
& it will start mg, go to accuracy

& GB will work same as GB

But one of the major difference is
they train DT differently. unlike
gini & entropy (vanilla).

So, thus way, DT train differently.

Algorithm

PD Take mean as base estimator
calculate prediction & error.

cgra	package	model 1	res 1	model 2	res 2
6.7	4.5	7.3	-2.6	6.49	-2.19
9	11	7.3	3.7	8.41	2.59
7.5	6	7.3	-1.3	6.49	-0.49
5	8	7.3	0.7	7.51	0.49

Now train a DT on cgra and res 1.
before that,

form a leaf node for residual

$$LN = [-2.6, 3.7, -1.3, 0.7]$$

calculate similarity score²

$$SSC = \frac{(\text{sum of residuals})^2}{(\text{no of res}) + 1}$$

for right now, regularization parameter $\lambda = 0$

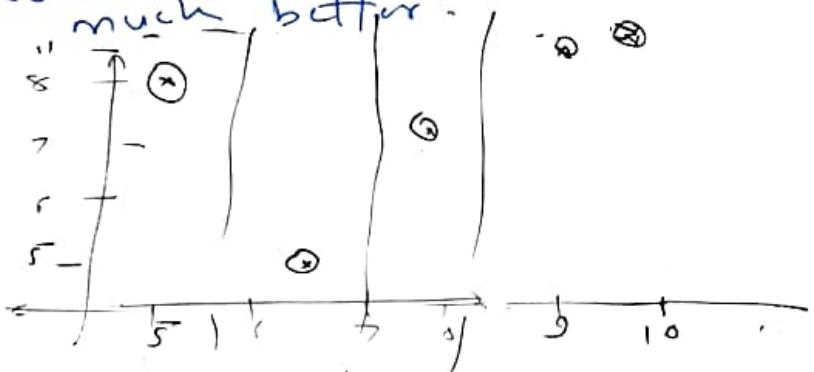
$$SSC = \frac{(-2.6 + 3.7 - 1.3 + 0.7)^2}{9 + 0}$$

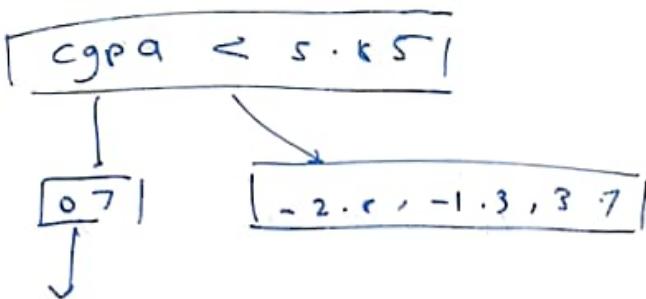
$$= \underline{(0.3)^2} = 0.02$$

we will now try to split residuals and see if we can increase similarity score
sort them & calc mean,

$$\begin{array}{cccccc} 5.0 & 6.7 & 7.5 & 5.0 \\ 5.45 & 7.1 & 6.25 \end{array}$$

Now see which one can improve similarity score "much better."

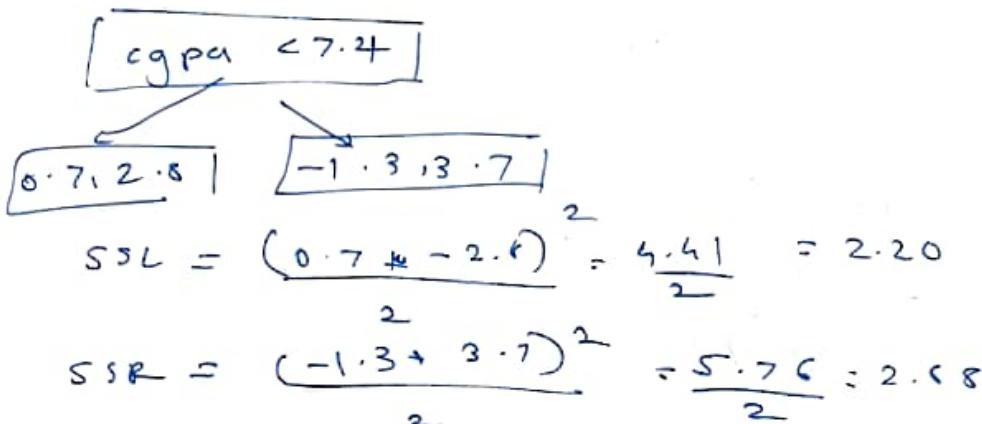




$$SS_{L1} = \frac{(0.7)^2}{1+0} = 0.49$$

$$SS_R = \frac{(-2.4 - 1.3 + 3.7)^2}{3+0} = \frac{0.16}{3} = 0.05$$

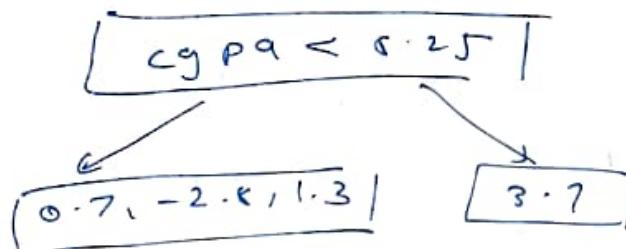
$$\begin{aligned} \text{gain} &= (SS_L + SS_R) - SS_{\text{root}} \\ &= (0.49 + 0.05) - 0.02 \\ &= 0.52 \end{aligned}$$



$$SS_L = \frac{(0.7 + 2.8)^2}{2} = \frac{4.41}{2} = 2.20$$

$$SS_R = \frac{(-1.3 + 3.7)^2}{2} = \frac{5.76}{2} = 2.88$$

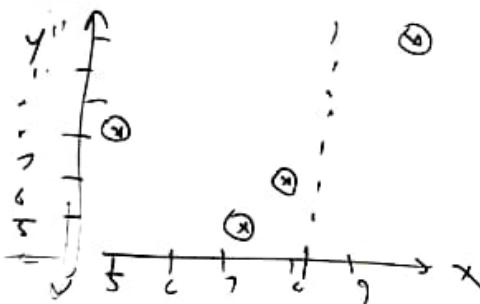
$$\begin{aligned} \text{gain} &= (SS_L + SS_R) - SS_{\text{root}} \\ &= (2.20 + 2.88) - 0.02 \\ &= 5.06 \end{aligned}$$



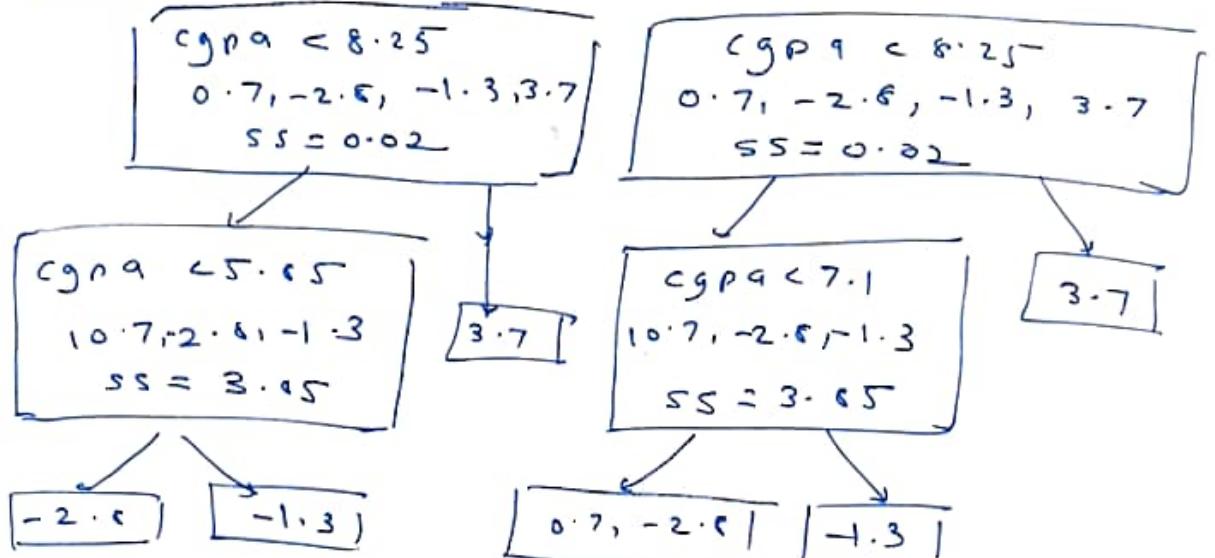
$$SS_L = 3.45$$

$$SS_R = 13.49$$

$$\text{gain} = 17.52 \rightarrow \text{choose this}$$



Now we can have two possibilities on set ones.



$$SS_L = \frac{(0.7)^2}{1} = 0.49$$

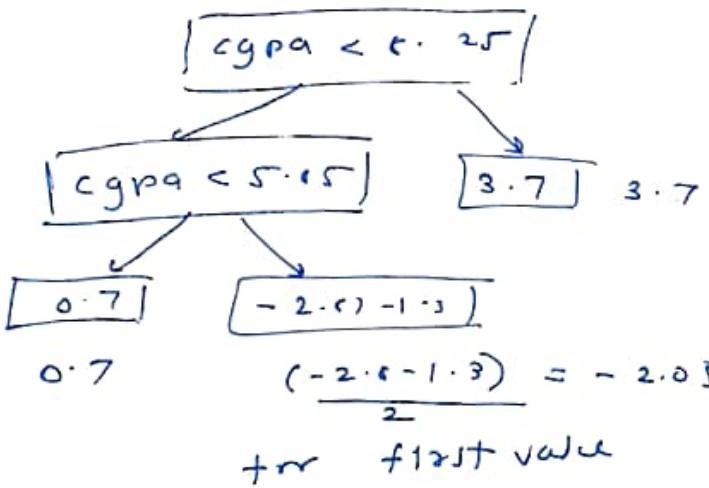
$$SS_R = 2.40$$

$$\text{gain} = 5.04$$

$$SS_L = 2.2$$

$$SS_R = 1.19$$

$$\text{gain} = 0.04$$



$$0.7 \quad \frac{(-2.6 - 1.3)}{2} = -2.05$$

for first value

package ? cyna 6.7 leaf

$$(7.3) + (0.7)(-2.05) = 6.79$$

\uparrow
(regularization)

In multiclass multiple features

so consider all, do all possible splits, calculate gain & SS and consider max gain.

In categorical

consider all unique In categorical, and consider both as a potential split point, calculate gain & similarity with max gain.

multi-class categorical variables

is identity unique & consider everyone as a potential split point, calculate GSS & maximum gain.

Two approach:

1) one single node + rest in a group

{small}, {big, long, tiny}

2) non empty non overlapping subset

{small, big}, {long, tiny}

Now for every possible split
calculate similarity scores, gain
and go with max gain.

By the ~~west~~ way,

there is not only one single algorithm,
the one we were trying out every time,
this is exact greedy algorithm.

in small dataset

But it is not only way,

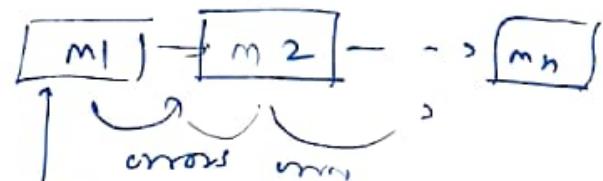
we can use approximation
in large datasets.

Xgboost in classification

We need to predict a student's placement.

cgpa	placement
5.70	0
6.25	1
7.10	0
8.15	1
9.00	1

Overview



DTs we train are not
vanilla like DT. they are
based on similarity score
rather than gini/entropy

so let's start problem solving.

phase 1) base model

$$\log(\text{odds}) = \left(\frac{P}{1-P} \right) = \log\left(\frac{P}{1-P}\right) = \log\left(\frac{3/5}{2/5}\right)$$
$$= \log\left(\frac{3}{2}\right) = \log(1.5)$$
$$= 0.405$$

But, we need to convert this to probability

$$P = \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}} = \frac{e^{0.405}}{1 + e^{0.405}} = 0.60$$

with output to phase 1 model, we need to calculate pseudo residual. so, it is placed-prob.

CGPA	placed model	prob	resid
5.70	0	0.405	-0.6
6.25	1	0.405	0.6
7.10	0	0.405	-0.6
8.15	1	0.405	0.6
9.0	1	0.405	0.6

phase 2) DT

Sort data on basis of input
for a leaf node

$$\{-0.6, 0.4, -0.6, 0.4, 0.6\}$$

$$SS = \sum (\text{residual})^2$$

$$\sum \text{prev_prbi} (1 - \text{prev_prbi}) + 1$$

$$= (-0.6 + 0.4 + 0.4 - 0.6 + 0.6)^2$$

$$5(0.6[0.4])^2 + 1 = 0$$

now, split on splitting criteria, goal is to
maximize similarity score, on basis of CGPA
splitting criteria is simple

$$\begin{array}{cccccc} 5.70 & 6.25 & 7.10 & 8.15 & 9.00 \\ | & | & | & | & | \\ 5.97 & 6.68 & 7.52 & 8.67 & \end{array}$$

calculate gain max on splitters.

sputting cuttoff 5.97

$$[\text{CGPA} < 5.97]$$

$$\begin{array}{c} -0.6 \\ \text{SSL} \\ \frac{(-0.6)^2}{(0.6)(0.4)} = 1.5 \end{array} \quad \begin{array}{c} 0.4, 0.6, 0.4, 0.4 \\ \text{SSR} \\ \frac{(0.4+0.4-0.6+0.4)}{(4 \times 0.6 \times 0.4)} = 0.37 \end{array}$$

gain

$$= (\text{SSL} - \text{SSR})$$

$$= \text{SSL} - \text{SSR}$$

$$= (1.5 - 0.37) - 0 \\ = 1.13$$

sputting cuttoff 6.17

$$\text{CGPA} < 6.17$$

$$\begin{array}{c} -0.6, 0.4 \\ \text{SSL} \\ \frac{(-0.2)^2}{2(0.6)(0.4)} = 0.16 \\ \frac{0.09}{2+0.24} = 0.04 \\ \frac{0.04}{0.24} = 0.16 \\ = 0.16 \end{array} \quad \begin{array}{c} 0.09 \\ 3+0.24 \\ = 0.05 \\ = 0.13 \end{array}$$

so DT is

on 7.25

$$\text{SSL} = 0.68$$

$$\text{SSR} = 1.33$$

$$\text{gain} = 2.22$$

for 8.575

$$\text{SSL} = 0.17$$

$$\text{SSR} = 0.67$$

$$\text{gain} = 0.69$$

$$[\text{CGPA} < 6.17]$$

$$\begin{array}{c} -0.6, 0.4, -0.6 \\ \text{SSL} \end{array} \quad \begin{array}{c} 0.4, 0.4 \\ \text{SSR} \end{array}$$

output for a particular leaf for classification

output

$$\frac{\sum(x_{i,j})}{\sum(p_{i,j}, p_{j,i}) (1 - p_{i,j}) + \lambda}$$

$$\begin{array}{c} -0.6 + 0.4 - 0.6 \\ 3(0.6 \times 0.4) \\ -0.6 \\ = \frac{-0.6}{3 \times 0.24} = -1.11 \end{array}$$

$$\begin{array}{c} (0.4 + 0.4) \\ (2 \times 0.24) \\ = 1.66 \end{array}$$

155

so basically this model says,

$$\boxed{\text{cga} < 7.66}$$

as an output
DTI

phase 2 model prediction

$$\frac{(m_1) + (0.3)(m_2)}{\text{log odds}}$$

DTI

$$(0.905 + (0.3)(\text{DTI}))$$

(2)

$$= \underline{0.905}$$

$$0.905 + (0.3)(-1.11) = 0.072$$
$$0.905 + (0.3) + 1.66 = 0.903 \rightarrow \text{log of odds format.}$$

so we need it in probabilities
we need it in probabilities

so we will do,

prob = $e^{\text{log odds}}$ for the residuals.

$$\text{prob} = \frac{e^{0.072}}{1 + e^{0.072}} = 0.518$$
$$= \frac{e^{0.905}}{1 + e^{0.905}} = 0.712$$

mathematics behind xgboost

Regression

$$\text{SSC} = \frac{\sum (\text{residuals})^2}{N+\lambda}$$

$$\text{OUT} = \frac{\sum (\text{residuals})}{N+\lambda}$$

Classification

$$\frac{\sum (\text{residuals})^2}{\sum (P(1-P)) + \lambda}$$

$$\frac{\sum (\text{residuals})}{\sum (P(1-P)) + \lambda}$$

for a combination chain of models,

$$[m_1] \rightarrow [m_2] \rightarrow [m_3] \rightarrow \dots \rightarrow [m_n]$$

this way it is a chain of functions,
so, let us say,

$$\hat{y} = f_1(x_1) + f_2(x_2) + \dots + f_n(x_n)$$

$$\hat{y} = \sum_{i=1}^n f_i(x_i)$$

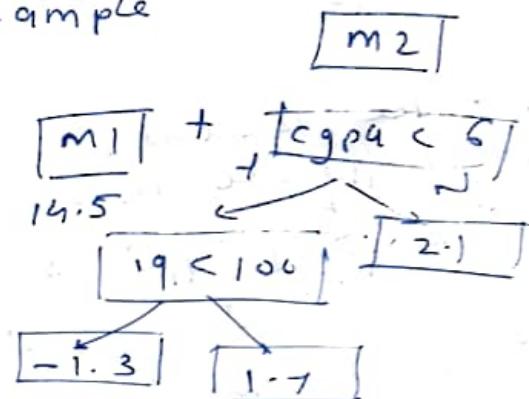
so, let us say our model is a XGB tree.

Now, for a output leaf weight (node)
the formula is,

$$\text{out} = \frac{\text{sum of resid for node leaf}}{N+1}$$

let us take an example

cgpq		19	percentage
6.5		71	13
5		91	10
3.1		120	16



$$\hat{y} = m_1 + m_2$$

$(f_1(x_1)) + (f_2(x_2)) \rightarrow$ calculate 3 different
values for output.
(for 3 leaf nodes)

$$\hat{y} = \sum_{i=1}^2 (f_i(x_i))$$

so, we go closer to reality, i.e., y_i .

so, the relation between y_i and \hat{y}_i , and
this is called loss function, we need
to minimize it.

In gradient boosting, there is no fixed
loss function,

$L[(y, \hat{y})]$ for gradient boosting.

we can put mse, log, etc. but it should
be differentiable

In xgb, there is a parameter λ .

according to RGB paper,

$$L = \sum_{i=1}^n L(y_i, \hat{y}_i) + \gamma_2(f_2(w))$$

↓ RP ← only on our
last one DT

regularization → $\gamma^T + \frac{1}{2} \lambda \|w\|^2$
parameter ↑ $w \in$ weight of node
no at leaf nodes in DT

so basically, we need such values of w_1, w_2, w_3 such that it should minimize the loss function.

derivation

$$L = \sum_{i=1}^n L(y_i, \hat{y}_i) + \gamma_2(f_k)$$

$$\frac{\text{loss} + \text{fun}}{\text{objective fun}}$$

let us take 2 cases

case 1) only one model (mean / log odds)

phase 1) $L^{(1)} = \sum_{i=1}^n L(y_i, \hat{y}) + 0 \rightarrow$ coz no DT

case 2) $\frac{m_1}{f_1(x_i)}$ ← so no consideration

phase 2) $\frac{m_1}{f_1(x_i)} \quad \frac{m_2}{f_2(x_i)}$ DT

so, $f_1 = f_1(x_i) + f_2(x_i)$

$$L = \sum_{i=1}^n [y_i, f_1(x_i) + f_2(x_i)] + \gamma_2(f_2(x_i))$$

case 3)

phase 3) $\frac{m_1}{f_1(x_i)} \quad \frac{m_2}{f_2(x_i)}$ DT $\frac{1}{f_3(x_i)}$ DT

(3) $L = \sum_{i=1}^n (y_i, f_1(x_i) + f_2(x_i) + f_3(x_i)) + \gamma_2(f_3(x_i))$

why only $f_3(x_i)$ in regularization in f_3 ?
because for f_2 we already handled
 f_2 regularization in phase 2.

stage t)

(59)

m_1

m_2

m_t

m_m

D_T

D_T

D_T

$(f_1 u_i)$

$(f_2 u_i)$

$(f_t u_i)$

last mod4

$$L^{(t)} = \sum_{i=1}^n L(y_i, f_1(u_i) + f_2(u_i) + f_3(u_i) + \dots + f_t(u_i))$$

$$L^{(t)} = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(u_i)) + r_2(f_t)$$

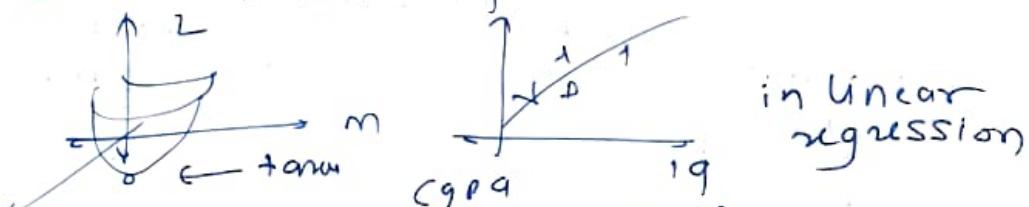
at any stage t,

$$L^{(t)} = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(u_i)) + r_2(f_t)$$

so our goal is to find such a value of $f_t(u_i)$ that should minimise all objective function.

But this optimisation problem is not as simple as linear regression,

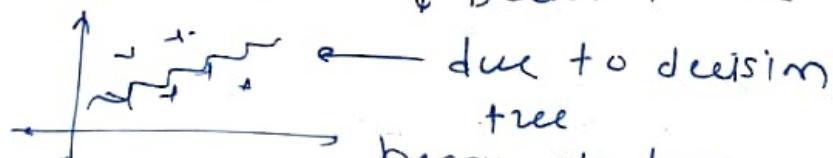
use in linear regression



b in Decision tree, use for now.

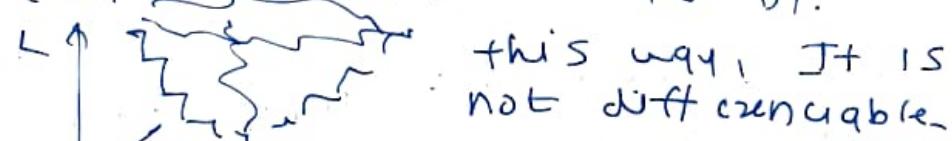
$$\hat{y}_i = \underbrace{\hat{y}_i^{(t-1)}}_{\text{constant}} + f_t(u)$$

↓ decision tree



because it has sudden jumps in values.

so, basically, if we imagine it in a 3D, it will not smooth as linear regression. It can be like this due to DT.



But for us, in

optimization, it is our requirement that this fun' should differentiable so we unable to find slope.

but what, if I try to replace this equation with a simple smooth approximation (closer to curve), we make it smooth and differentiable?
we can do it by Taylor series.

Taylor Series

A powerful mathematical approximating technique you can give it a complex equation and then it will convert to simple by a polynomial equation.

As we add more polynomial terms, x^2, x^3, \dots , our approximation getting better & accurate. It uses derivatives (first, second, third), to convert a complex function to a simple polynomial.

To approximate a funcⁿ around point a ,

$$f(x) \approx f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 - \dots - \frac{f^{(n)}(a)}{n!}(x-a)^n$$

for example

$$f(x) = e^x \quad a=0$$

$$f(a) = e^a = 1$$

$$f'(a)(x-a) = e^a(x-a) = x$$

$$f''(a) \frac{(x-a)^2}{2!} = e^a \frac{x^2}{2} = \frac{x^2}{2}$$

$$f'''(a) \frac{(x-a)^3}{3!} = \frac{x^3}{3!}$$

approx e^x using Taylor series

$$e^x = 1 + x + \frac{x^2}{2} + \frac{x^3}{6} \dots \frac{x^n}{n!}$$

this way we need to apply Taylor series on our function.

$$f(x) \approx f(a) + f'(a)(x-a) + \frac{f''(a)(x-a)^2}{2!}$$

$$\approx \left[\sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) \right] + \gamma (f_t(x))$$

$$x = \hat{y}_i^{(t-1)} + f_t(x)$$

a = changing point.

Basically, we will take a as $\hat{y}_i^{(t-1)}$.
This after will make the change.

$$f(a) = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(x_i))$$

$$f(a) = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)}) \quad \text{gradient}$$

$$f'(a)(x-a) = \sum_{i=1}^n \frac{\partial L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}} \leftarrow f_t(x)$$

$$f''(a)(x-a)^2 = \sum_{i=1}^n \frac{\partial^2 L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)} \partial x_i} \leftarrow \text{hessian} \quad f_t''(x)$$

now, the equation will become,

$$L = \sum_{i=1}^n \frac{L(y_i, \hat{y}_i^{(t-1)})}{10} + g_1 f_t(x) + \frac{1}{2} h_1 f_t^2(x) + \gamma (f_t(x))$$

$$L = \sum_{i=1}^n [g_1 f_t(x) + \frac{1}{2} h_1 f_t^2(x)] + \gamma (f_t(x))$$

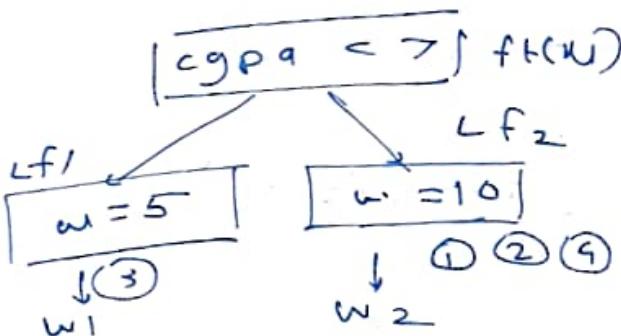
$$L = \sum_{i=1}^n \frac{[g_1 f_t(x) + \frac{1}{2} h_1 f_t^2(x)]}{1-n(\text{no of nodes})} + \gamma T + \frac{1}{2} \sum_{j=1}^T w_j^2$$

basically these two summaries are different.

we are going to do

$$w_j \rightarrow f_t(x_i) \left| \begin{array}{c} \sum_{i=1}^n \rightarrow \sum_{j=1}^T \sum_{i \neq j} \end{array} \right.$$

so, let us take a simple example again.



cgpa	package
7.1	4.5
8.2	6.9
6.5	5.1
9.1	10

$$\begin{aligned}
 & \left(g_1 f_t(x_1) + \frac{1}{2} h_1 f_t^2(x_1) \right) + \left(g_2 f_t(x_2) + h_2 f_t^2(x_2) \right) + \\
 & \quad \cdots \left(g_4 f_t(x_4) + \frac{1}{2} h_4 f_t^2(x_4) \right)
 \end{aligned}$$

what being $i \in I_j \rightarrow$ instance set of leaf j

for Lf_1 , we have only one row.

for Lf_2 we have 3 rows.

$j = 1$

$$\begin{aligned}
 & g_3 f_t(x_3) + \frac{1}{2} h_3 f_t^2(x_3) + g_1 f_t(x_1) + \frac{1}{2} h_1 f_t^2(x_1) \\
 \xrightarrow{\text{same}} & + g_2 f_t(x_2) + \frac{1}{2} h_2 f_t^2(x_2) \\
 & + g_4 f_t(x_4) + \frac{1}{2} h_4 f_t^2(x_4)
 \end{aligned}$$

$w_j \rightarrow$ output of tree for a particular leaf node

$$\begin{aligned}
 L^{ft} &= \sum_{j=1}^t \left[\sum_{i \in I_j} g_i w_j + \frac{1}{2} \sum_{i \in I_j} h_i w_j^2 \right] \\
 &\quad + \gamma T - \frac{1}{2} \lambda \sum_{j=1}^T w_j^2
 \end{aligned}$$

$$L^{ft} = \sum_{j=1}^t \left[\sum_{i \in I_j} g_i w_j + \frac{1}{2} \sum_{i \in I_j} h_i w_j^2 + \frac{1}{2} \lambda w_j^2 \right] + \gamma T$$

$$L^{ft} = \sum_{j=1}^T \left[\sum_{i \in I_j} g_i w_j + \frac{1}{2} \left(\sum_{i \in I_j} h_i + t \right) w_j^2 \right] + \gamma T$$

now we are going to differentiate it.

so now, for the differentiation,

$$L^{(t)} = \sum_{j=1}^T \left[\sum_{i \in I_j} g_i w_i + \frac{1}{2} \left(\sum_{i \in I_j} h_i + 1 \right) w_i^2 \right] + \gamma T$$

$$\frac{\partial L}{\partial w_j} = \sum_{i \in I_j} g_i^j + \left(\sum_{i \notin I_j} h_i + \lambda \right) w_j = 0$$

$$w_j = - \frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$$

Now, we

$$g = \frac{\partial L(y_i, \hat{y}_i^{<t-1>})}{\partial \hat{y}_i^{<t-1>}}, h = \frac{\partial^2 L(y_i - \hat{y}_i^{<t-1>})}{\partial y_i^{<t-1> 2}}$$

$t \rightarrow \text{mean} \rightarrow$

on iteration, we say, on model training,
we are on nth model

$$\boxed{\quad} \quad \boxed{\quad} \quad \boxed{\quad} \quad \boxed{\quad} \quad \boxed{(t+1)} \quad t$$

$\xrightarrow{\text{on this step loss}}$

In case of regression, MSE

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

$$g_i = \frac{\partial L}{\partial \dot{q}_i} = - \sum_{j=1}^n (q_j - \dot{q}_j) = \sum_{j=1}^n (\dot{q}_j - q_j)$$

so, basically

$$\sum_{i \in I_j} \frac{(y_i - \hat{y}_i(t))}{\text{sum of residuals}} \leftarrow \begin{matrix} \text{residuals for} \\ \text{every point} \end{matrix}$$

$$h_i = \frac{\partial L}{\partial f_i} = -(\hat{f}_i - y_i)$$

$$\frac{\partial^2 L}{\partial y_1^2} = \frac{\partial}{\partial y_1} (\hat{y}_1 - y_1) = 1$$

so moving forward

$$= \frac{\sum_{i \in I_j} (y_i - \hat{y}_i)^2 / T}{\sum_{i=1}^T 1 + \lambda}$$

= no of residuals in a point

$$\text{output } w_j = \frac{\text{sum of residuals}}{\text{no of residuals} + 1} \quad \text{--- (1)}$$

classification problem

$$\frac{\log f(y)}{\log s_s} = \log \text{loss}$$

$$w_j = \frac{\text{sum of residuals}}{\sum_{i \in I_j} p_i(1-p_i) + \lambda} \quad \begin{array}{l} p_i \rightarrow \text{predicted} \\ \text{probability of} \\ \text{previous itn} \end{array} \quad \text{--- (2)}$$

so put w_j in loss function, Basically calculate minimum loss on minimum w_j .

$$L^{(t)}(q) = \sum_{j=1}^T \left[\left(\frac{\sum_{i \in I_j} g_i}{\sum_{i=1}^T g_i} \right) \left(\frac{\sum_{i=1}^T g_i}{\sum_{i \in I_j} h_{i+1}} + \frac{\sum_{i \in I_j} h_{i+1}}{\sum_{i \in I_j} h_{i+1}} \right) \right] + \gamma T$$

$$L^{(t)}(q) = \sum_{j=1}^T \left[\frac{- \left(\frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_{i+1}} \right)^2}{\sum_{i \in I_j} h_{i+1}} + \frac{1}{2} \frac{\left(\sum_{i \in I_j} g_i \right)^2}{\sum_{i \in I_j} h_{i+1}} + \lambda T \right]$$

$$= \sum_{j=1}^T \left[-\frac{1}{2} \frac{\left(\sum_{i \in I_j} g_i \right)^2}{\sum_{i \in I_j} h_i + \lambda} \right] + \lambda T$$

$$\boxed{= -\frac{1}{2} \sum_{j=1}^T \frac{\left(\sum_{i \in I_j} g_i \right)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma T}$$

SSC

minimum loss for lost model added so basically, this equation can be used as a scoring function, basically to measure the quality of last one tree we made.

similarity score

$$= -\frac{1}{2} \sum_{j=1}^T \frac{\left(\sum_{i \in I_j} g_i \right)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma T$$

and this is exactly.

regression

$$SSC = \frac{\left(\sum \text{residuals} \right)^2}{N + \lambda}$$

classification

$$SSC = \frac{\sum (P \cdot (1-P)) + \lambda}{N + \lambda}$$

How to build tree using similarity scores?
Let us say, we are on t th step in modus.



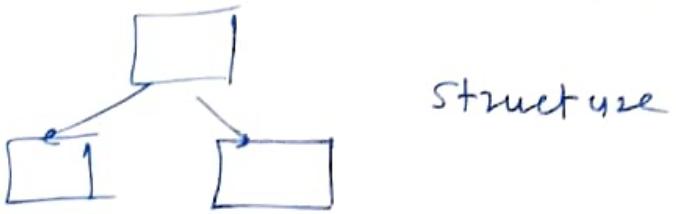
for a leaf node

$$L^{(1)}(q_p) = -\frac{1}{2} \frac{\left(\sum_{i \in I_p} g_i \right)^2}{\left(\sum_{i \in I_p} h_i + \lambda \right)} + \gamma T$$

for one leaf node.

(165)

Then we will do split for the leaf



We splitted it and made new nodes, because we wish to reduce the loss too more.

Then we will now calculate overall loss for tree structure

$$L^{(t)}(q_c) = \frac{1}{2} \left[\frac{\sum_{i \in I_L} (g_i)^2}{\sum_{i \in I_L} \text{hit}} \right] + \left[\frac{\sum_{i \in I_R} (g_i)^2}{\sum_{i \in I_R} \text{hit}} \right] + 2\gamma \quad \because (T=2)$$

so now I have both parent & child loss. I wish to know how much loss reduced after the split.

$$\begin{aligned} L^{(t)}(q_P) - L^{(t)}(q_C) &= \left\{ \left[\frac{-\frac{1}{2} \left(\sum_{i \in I_P} g_i \right)^2}{\sum_{i \in I_P} \text{hit}} + \gamma \right] \right. \\ &\quad \left. - \left(\frac{1}{2} \left[\left(\frac{\sum_{i \in I_L} g_i}{\sum_{i \in I_L} \text{hit}} \right)^2 + \left(\frac{\sum_{i \in I_R} g_i}{\sum_{i \in I_R} \text{hit}} \right)^2 \right] + 2\gamma \right) \right\} \\ &= \left[\frac{-\frac{1}{2} \left(\sum_{i \in I_P} g_i \right)^2}{\sum_{i \in I_P} \text{hit}} + \frac{+\frac{1}{2} \left(\sum_{i \in I_L} g_i \right)^2}{\sum_{i \in I_L} \text{hit}} + \frac{+\frac{1}{2} \left(\sum_{i \in I_R} g_i \right)^2}{\sum_{i \in I_R} \text{hit}} \right] - \gamma \end{aligned}$$

$$= \frac{1}{2} \left[\frac{\left(\sum_{i \in IL} g_i \right)^2}{\sum_{i \in IL} \text{hit}} + \frac{\left(\sum_{i \in IR} g_i \right)^2}{\sum_{i \in IR} \text{hit}} - \frac{\left(\sum_{i \in IP} g_i \right)^2}{\sum_{i \in IP} \text{hit}} \right]$$

This is exactly formula for when

we did

$(SSL + SSR) - SSP$ to calculate loss.

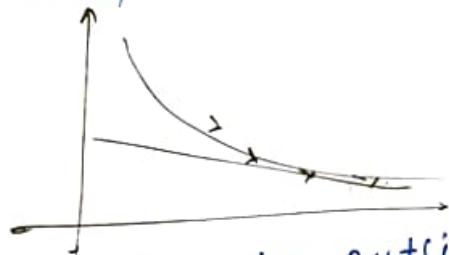
DBScan (Density Based spatial Clustering of Apparition with noise)

Why DBScan? Why to change

Some elbows for k-means

① need to tell how many clusters, should be done.

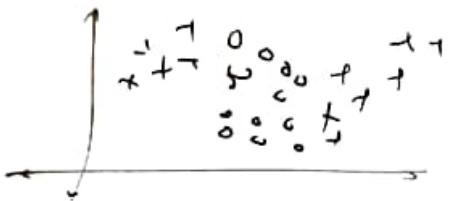
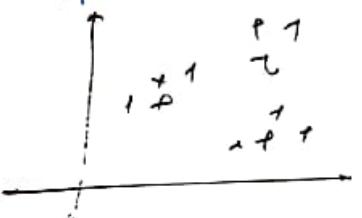
done. So we use elbow method. but many times it is difficult to tell.



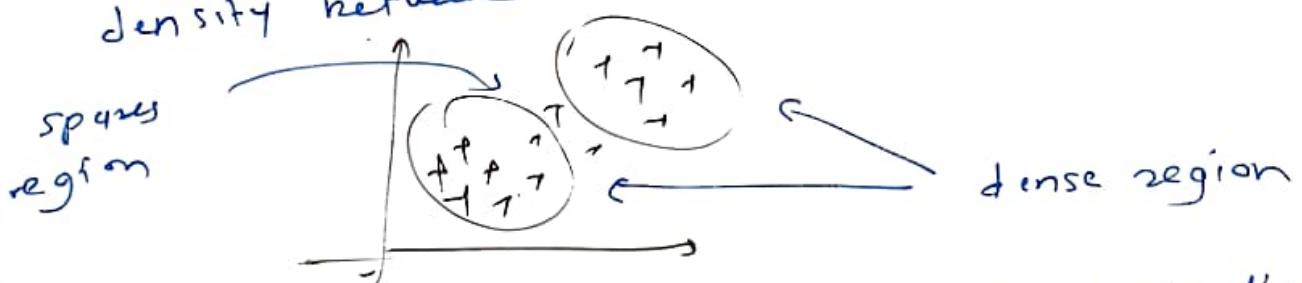
② very sensitive to outliers.

③ k-means anywhere can't form sphericals.

data points.

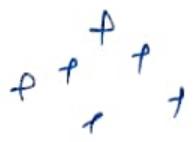


Density clustering
we try to cluster the data on the basis of density between two data points.



basically, DBScan see data, separate them out form sparse region, and separate dense region.

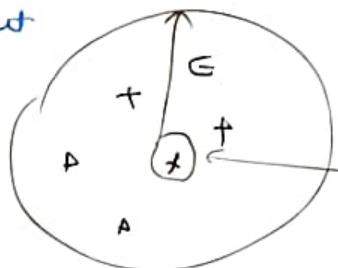
basically let's say,



we need to define a threshold, and it will work radius for the targeted point. and if some K number of ~~points~~ come under the radius of that specific threshold, we need to consider it dense if they are less, then sparse.

some concepts

1) core point

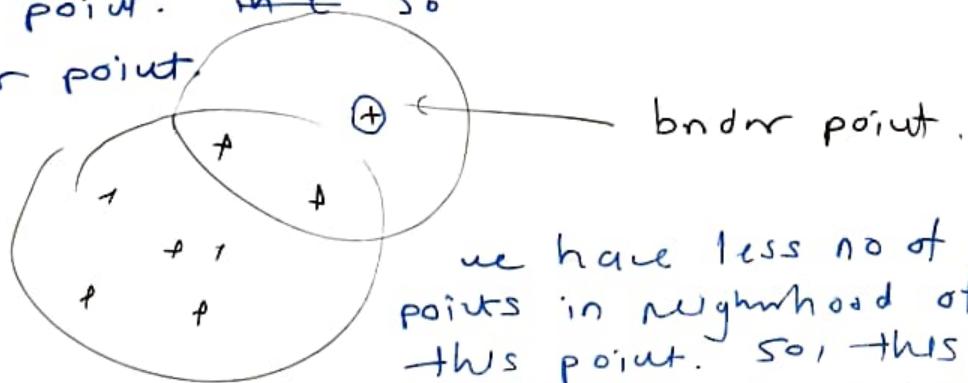


core points

MIN point = 4

in E neighborhood, we have greater or equal to datapoints, that is called as core point. ~~not E~~ so

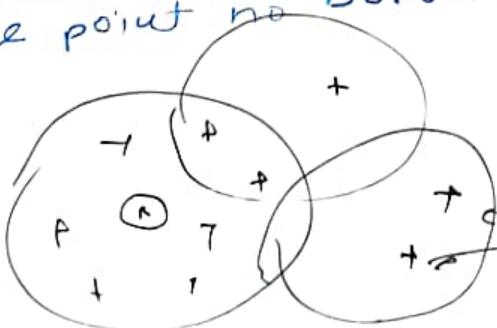
2) Border point.



we have less no of points in neighborhood of this point. So, this is a border point. BUT it have also a core point in E region.

3) noise point

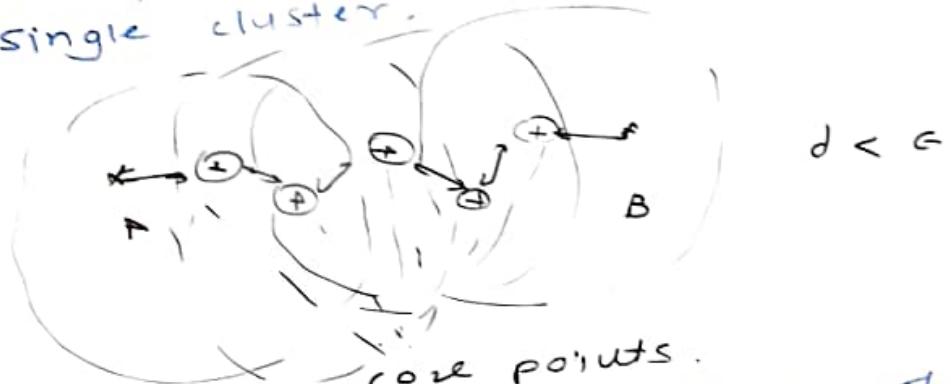
No core point no border point



noise points

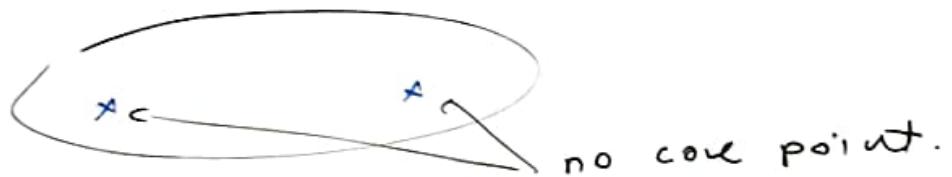
density connected points.

If A and B are two points, and if they are density connected, then we put both them in a single cluster.

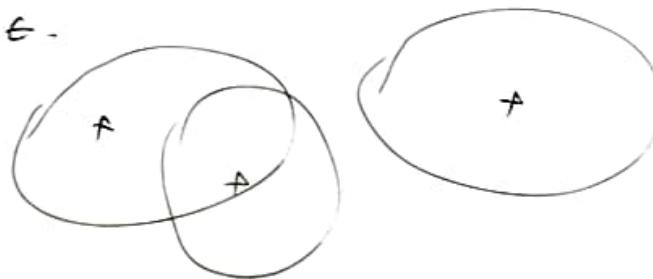


basically under the influence of core points. Basically connected indirectly.

This fails when the main cause of failure can be



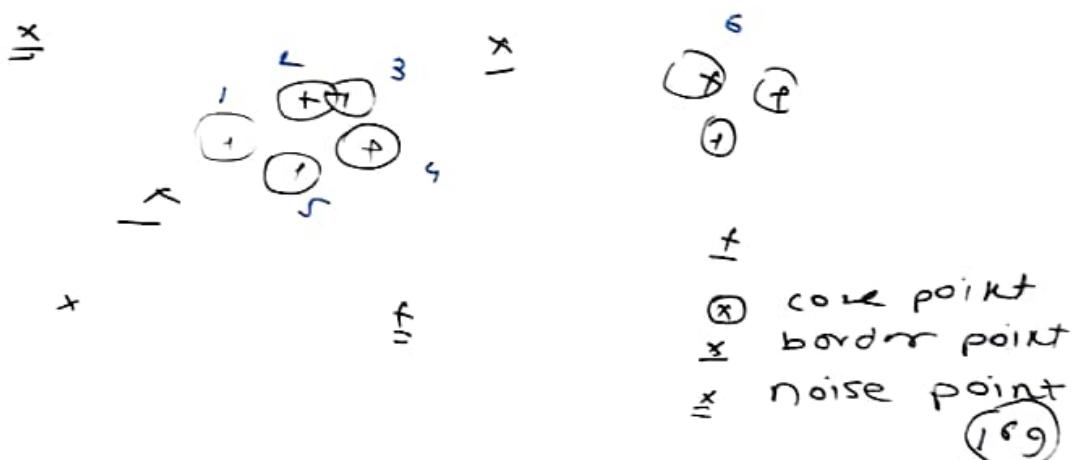
If distance between two core points is greater than ϵ .



Algorithm

Step 0) define ϵ and no points.

Step 1) classify core, border, noise



add all points that are unclustered and density connected to the current point to this cluster.

Let start with ① and check conditions one point's density connected?

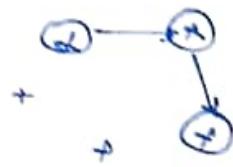
$$① \rightarrow ② \quad \checkmark$$

diagram

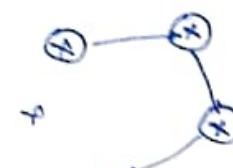
$$④ \rightarrow ⑦$$



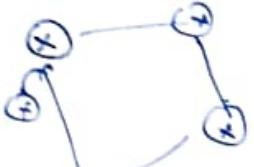
$$④ \rightarrow ⑤ \quad \checkmark$$



$$① \rightarrow ④ \quad \checkmark$$



$$① \rightarrow ⑤ \quad \checkmark$$



$$① \rightarrow ⑥ \quad \times$$



a cluster formed

do now for all unclustered points.

Step ③ for each unstructured border point, assign it to the nearest cluster point.

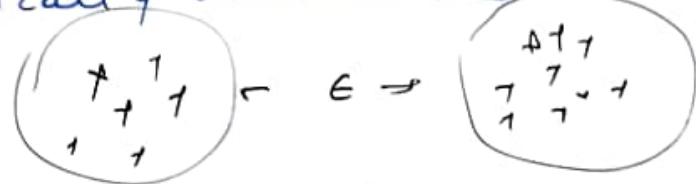
Step ④ leave all noise points as it is.

Advantages

- 1) Robust to outliers
- 2) no need to specify clusters.
- 3) can find arbitrary shaped clusters
- 4) 2 hyp to tune (m_n, ϵ)

disadvantage

- ① sensitive to hyperparameters
- ② difficulty with varying density clusters



- ③ cannot do prediction

* * * * *

completed playlist!

Credit goes to CampusX!



171