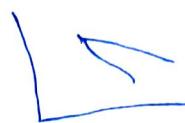


Maths is fun.

Slope

Eq of a line

$$y = mx + c$$



The major diff

$$\text{Slope} \rightarrow \frac{\text{change in } y}{\text{change in } x} = \frac{y_2 - y_1}{x_2 - x_1}$$

gradient:

↳ how steep the line is

$$\frac{m}{y_2 - y_1} = \frac{38}{6} \approx 7 \text{ not define}$$

if  $y = x^2$  → the how to calculate the slope of this in line is easy

→ because these not same perpendicular

in line close had constant rate

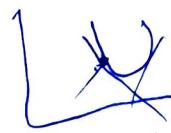
$$\text{can write: } \frac{y_2 - y_1}{x_2 - x_1} \times$$

in this case slope is very  
↓ lots

Steepness varies (so slope of 2 point not same)

so slope w.r.t 1 point that is called → instantan

instantaneous rate of change w.r.t y (called)



tangent/slope

is a line that touches a curve at one point

then how calculate

$$X \left[ \frac{y_2 - y_1}{x_2 - x_1} \right] \rightarrow \left\{ \frac{dy}{dx} = f'(x) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h} \right\} \rightarrow \text{first principle}$$

slope of tangent

this is called differentiation  
in range to 0 not 0

instantaneous rate of change w.r.t x  
in y w.r.t x

$$\frac{f(a+h) - f(a)}{h} \quad h \rightarrow 0 \quad y = a + h$$

Rules

$$\textcircled{1} \quad \frac{d}{dx} n^x = n^x \ln n$$

negative

gradient is slope  
+ve or -ve  
so line is upward called  
positive

$$\textcircled{2} \quad \frac{d}{dx} n^n = n^n \cdot n^{n-1}$$

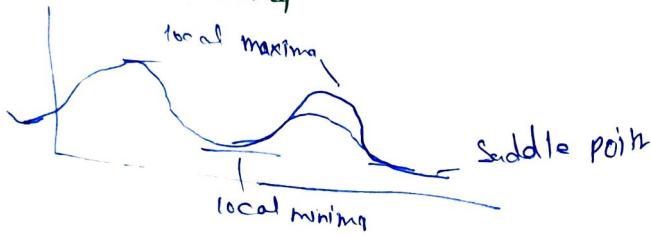
(2)

\* chain rule

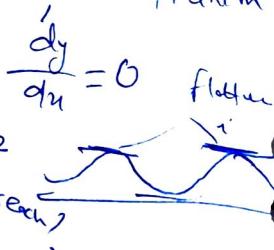
$$\frac{d}{dx} f(g(x)) = f'(g(x)) \cdot g'(x)$$

$$\begin{aligned} & \text{eg } (\sin(x^2)) \\ & \text{f'(x)} + g'(x) \\ & f'(g(x)) \end{aligned}$$

Use Case → in gradient descent (for ML optimization)  
# Maxima/minima



• Step = 0 (flattening)  
In case of  
maxima / minimum



Physics

① A trajectory of ball is followed by  $y = 4 + 10t - 5t^2$ . What is max. height that ball can reach?

We know that max. height & min. depth. Slope is 0. We know so we find the

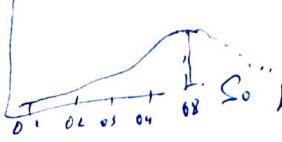
$$\begin{aligned} & \text{So find slope} \\ & \frac{dy}{dt} = \frac{d}{dt}(4 + 10t - 5t^2) = 0 + 10 - 10t \\ & = 10 - 10t \end{aligned}$$

So = Slope at max height

$$10 - 10t = 0 \Rightarrow t = \frac{10}{10} = 1 \text{ sec}$$

at 0.83 sec, it will reach max ht.

Why you say max it may be min



$$\begin{aligned} & \text{So height is} \\ & = 4 + 10t - 5t^2 \\ & = 4 + 10 \times 0.83 - 5 \times (0.83)^2 \end{aligned}$$

$$\Rightarrow 8.166 \text{ m at } 0.83 \text{ sec}$$

do second order diff in this

Matty (3) → Stats  
→ tool kit  
→ EDA

## ① Second order derivative

→ 10-12t

→ When after slope is 0 at x and the second order derivative at x is :-

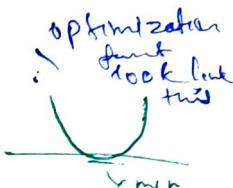
→ less than 0, it is local minimum

→ greater than 0, it local minima

→ equal to 0, you cannot say anything

$$\frac{d}{dt} (10-12t) = -12 \rightarrow +ve$$

$\frac{dy}{dx} < 0 \rightarrow$  it means  $\frac{dy}{dx} = 0$  is of local maxima.



18 April

→ Maths → behind ML → intuition ① ✓  
implementation ② ✓  
→ Maths ③ little bit

→ intro ML

Terminology

Pre-process/feature Engineering

→ Matrices

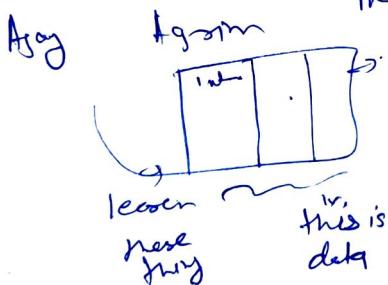
Show examples / motivation of ML

→ Cataloguing ads example

4 Alarms  
→ Automation  
→ Weather forecast  
→ GPS, Board, LLM (Stem-MDL → NLP)

this is process

6



mall example

Age, promotion, discount, offers → Promote

↳ FB, tinder, bin, Snap → social tab

→ no student is bad  
Teacher's View  
④ how I teach  
6 ask question  
② Practice/Practice  
& Practice

↑ of People  
Not a Niche  
(list no)  
methodology  
→ Strengthen

②

What is Machine Learning?

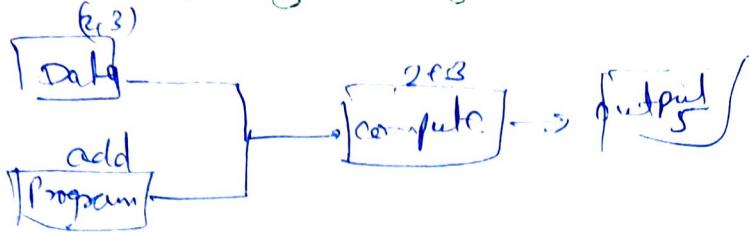
→ Machine learning patterns from the data and tries to replicate the same in future.

Arthur Samuel  
1959

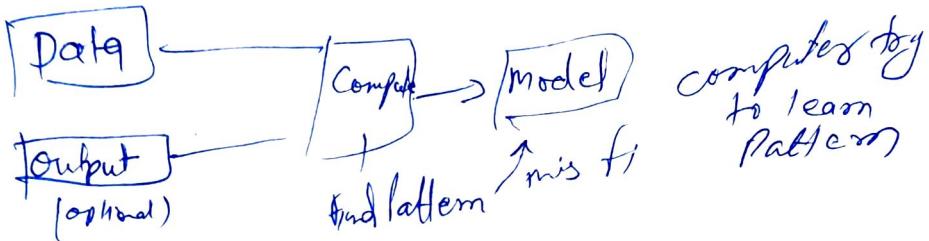
→ The subfield of CS that gives complete the ability to learn without explicitly programmed

→ Tradition Programming Paradigm

## ⇒ Traditional Programming Paradigm



## ⇒ ML Programming Paradigm



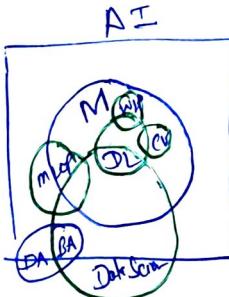
Tom Mitchell  
(1997)  
Geoffrey Hinton  
Andrew Ng

1989 → 1997  
1997  
2011  
2016  
2017  
2017

Linear Regression  
Logistic Regression } old

AI vs ML vs DL vs DS

→ Introduction to ML



→ types of → SL

Historical data

Supervised

Predict Price of House

feature targets

labels

Regression Classification

y is continuous

e.g. Rice

y is discrete (finite)  
spam/ham

SSL

No Supervision

→ y is not present

No Historical data

Then can you find trends in that data

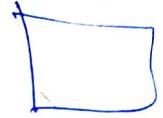
Semi-supervised learning

Combination

→ So either you drop labels & do SL

e.g. Netflix recommendation

Reinforcement

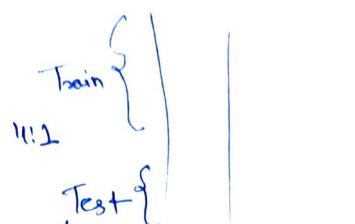


Self driving car

(5)

→ Test the model on some other data  
 take this model  
 we test

so Time/Money  
 take



### Train test split

→ what ration should decide?  
 Q No any ideal any? / Many auto off

big data → SD:SD → x why  
 small data → 70:30  
 80:20  
 90:10 → more data to train

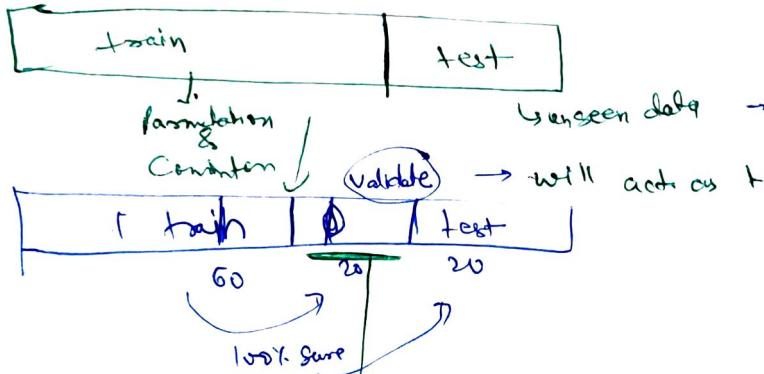
train > test  
 data

representation of unseen data

→ used to test model accuracy

D/M

↓  
 appears  
 for interview  
 rock interview  
 practice



### • data leakage

### • Why Validation is need

Hyperparameters tuning

Ex: Bateman  
 Virat kholi → Practice  
 Board retain for

Coach

R.D.

tuned

Hyperparameter tuning

Parameter → Mathematical relationship

you have used  
 test information to  
 the training of the model

Model will Perform well  
 on both

Study

Study

① NCERT  
 study: Sample  
 paper

② Tenay

③ Board Exam

11  
 test of Model

loss of  
 time/money  
 Model perform  
 fair

(use)

Data leakage  
 & five board Exam  
 papers are given  
 to Students

(P)

• Costas  
Validation

## ML framework

- f is set of rules
- ↳ a predefined process
- ↳ set of tools
- ↳ guidelines
- ↳ Standardized process

• Area of focus

many  
ways  
Price  
20%

↳ many → variety of mode → Understanding the parts  $\rightarrow$  Mathematical relationships

20%

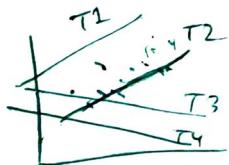
$y = f(x)$   
Input  
Output  
↳ Predictive functions

→ Model training → Given the data estimate the prediction function to minimize error

to  
ml edges

evaluation  
matrix

② If you learn from Training data  
then apply it to new prediction



## Overfitting & Underfitting

train  $\rightarrow$  Model  $\xrightarrow{\text{train}}$  Accuracy 71% 95%

test  $\rightarrow$  Model  $\xrightarrow{\text{test}}$  Accuracy is low 60%

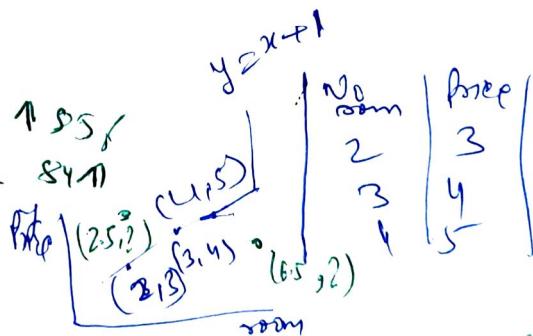
→ Model perform well on train data but worse on test data

Underfit  $\rightarrow$  Train  $\rightarrow$  accuracy  $\downarrow$  80%  
test  $\rightarrow$  accuracy  $\downarrow$  25%

## Best Generalization Model

Train  $\rightarrow$  Acc. 100%  
test  $\rightarrow$  Acc. 84%

overfitting  
↳ memorize the data



why?  $\rightarrow$  will you be able to predict  
new data  $\rightarrow$  nice  
they don't fall in the pattern

7

## Generalize Model

How to identify ??  
 Model is Underfitting  $\rightarrow$  not well while training

- overfitting  $\rightarrow$   $t_{\text{train}}$  is high  $\uparrow$   $t_{\text{test}}$   $\downarrow$   
 diff more than 5%

Generalized model  $\rightarrow$   $t_{\text{train}} \uparrow$   $t_{\text{test}} \uparrow$   
 $t_{\text{train}} \approx t_{\text{test}}$  diff  $\approx 5\%$  range

Preassumption

### Bias Variance

\* training error is also known as bias

High training error means high Bias

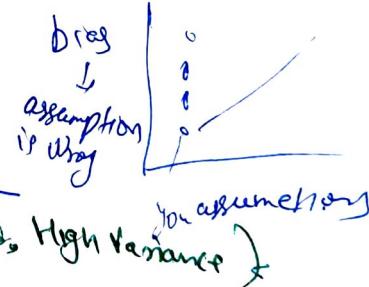
\* testing error is also known as Variance.

High testing error means high Variance.

Underfitting:

$\left. \begin{array}{l} \text{High} \\ \text{Train} \rightarrow \text{Accuracy} \downarrow \text{of High training error} \rightarrow \text{High Bias} \\ \text{Test} \rightarrow \text{Acc} \uparrow \rightarrow \text{High testing error} \rightarrow \text{High Variance} \end{array} \right\}$

$\left. \begin{array}{l} \text{Underfitting} \rightarrow \text{High Bias} \\ \text{High Variance} \end{array} \right\}$



Train Validate  
for rebates  
to imp

use Validation  
data not test data

Overfitting  $\rightarrow$

$\rightarrow$  (low bias, high variance)

### Variance

$\rightarrow$  why Variance is analogous to High testing error?

$\rightarrow$  Test data varied from the Pattern,  
 that why high testing error



Train  $\rightarrow$  trained  $\rightarrow$  Accuracy  $\uparrow$   $\rightarrow$  mean low bias (low  $t_{\text{train}}$ )

Test  $\rightarrow$  tested  $\rightarrow$  (error)

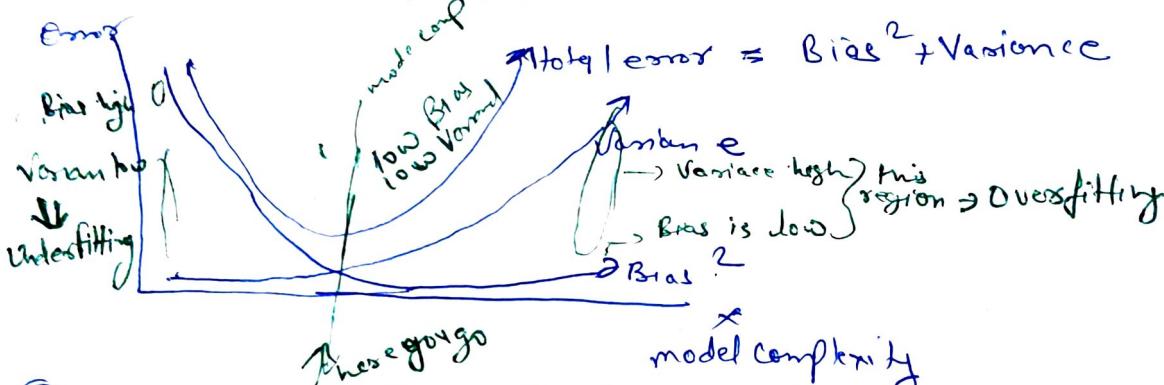
Accuracy (high test error)  $\rightarrow$  High Variance

(Generalised Model)  $\rightarrow$  train acc $\uparrow$  (low bias)  $\rightarrow$  low error in training  
 (low bias, low variance) test acc $\uparrow$  (low variance)  $\rightarrow$  low error in test

## Bias-Variance Trade off

it describe the relationship b/w a model's complexity and the accuracy of its predictions and how well it can make predictions on previously unseen data that were not used to train the model.

Variance  
 spread out  
 bias  $\rightarrow$  how close to target you are



① where you want to go (know)

↳ means you used more than required features.

$\rightarrow$  in order to reduce the variance produce high bias.

Simply inject the bias (means drop unnecessary feature)  
 ↳ increase bias, Select only relevant features

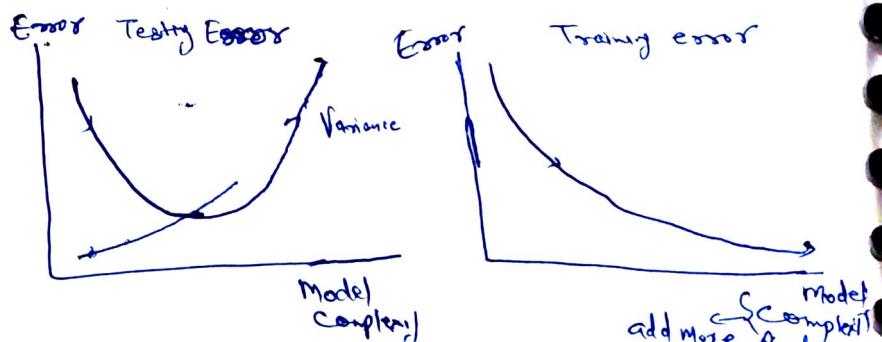
\* In order to reduce Variance (reduce the test error)  
 ↳ So introduce some bias in training of model by selecting relevant features or high bias

↳ Don't learn Noise / outliers

↳ Use some algorithms which is not overfitting

\* In order to reduce high bias  $\rightarrow$  Use some more datapoint, do more feature engineering or use some other algorithms

PCA in m.  
 ↳ Vectors



Model complexity  
 add more feature

(9)

Vectors

$$\begin{bmatrix} 5 \\ 4 \end{bmatrix}$$

Set of vectors

$$w^T x + b = 0$$

$$OA^2 = OB^2 + AB$$

$$\begin{bmatrix} 5 \\ 4 \end{bmatrix}$$

$\begin{bmatrix} u \\ v \end{bmatrix}$

Same my These can be multiple vectors in space of same magnitude & direction, the origin is different.

Operations on Vectors

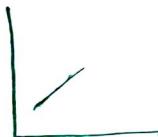
represent  $\mathbb{R}^2 \rightarrow 2D$  space  
Real numbers

$$i = \begin{bmatrix} 5 \\ -4 \end{bmatrix}; j = \begin{bmatrix} -4 \\ 5 \end{bmatrix}$$

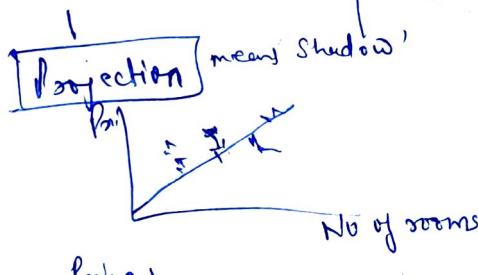
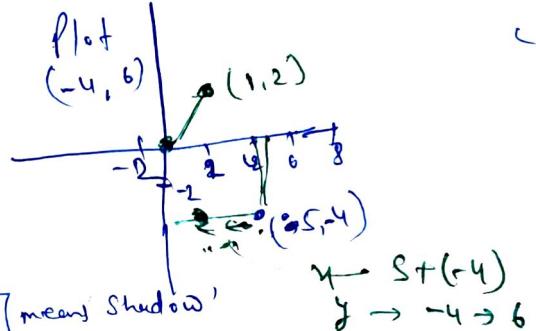
$$i, j \in \mathbb{R}^2$$

Add  $i + j = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \begin{bmatrix} 5 + (-4) \\ (-4) + 6 \end{bmatrix}$

vector is point



Projection of this line



Projection of a data point will be smaller than best fit line  
 that in same direction

Q) Why do Projection

$$U_0 = \begin{bmatrix} U_1 \\ U_2 \end{bmatrix} \& \text{8 magnitude}$$

$$2 + 8 = 10$$

$$\begin{bmatrix} w \\ x \\ - \\ - \\ - \\ - \\ - \\ - \\ - \\ - \end{bmatrix}$$

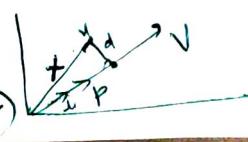
$$8 \text{ DataPoint} = 16 \text{ car}$$

$$2 \text{ cols} = 16 \text{ coordinate}$$

Can we all datapoint represent by

$$\text{How we chose } U = \frac{V}{\|V\|}$$

unit vector (magnitude is 1)



according to vector rule  
 $U = P + d$   
 $d = x - P \Rightarrow P = x - d$   
 $d = x - kU$

(16)

Then dot product will be 0  $a \cdot b = ab \cos 90^\circ$

$$8m - \frac{V}{M} = 4 \text{ m/s}$$

$\Rightarrow 8m$  East

$\Rightarrow 8m$

$\Rightarrow 1$  East Unit vector

$\Rightarrow 0$

$$\text{prod} = 0$$

$$kx_4 \times (x^2 - kx_4)$$

$$\Rightarrow kx_4 - k^2 x_4 \cdot x_4 = 0$$

$$= k(x_4 - kx_4 \cdot x_4)$$

$$x_4 - kx_4 \cdot x_4 > 0$$

$$k = 2x_4$$

Magnitude  $x_4 = 1$

$$u \cdot u = 1$$

## Linear Algebra

What is linear algebra?

& use

How to do?

## Matrix Addition

$$\begin{bmatrix} 2 & 2 \\ 4 & 1 \end{bmatrix} + \begin{bmatrix} 2 & 4 \\ 2 & 3 \end{bmatrix}$$

$$3x + 2y = 8$$

$$5x + 4y = 18$$

$$\begin{bmatrix} 3 & 2 \\ 5 & 4 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 8 \\ 18 \end{bmatrix}$$

# Matrix Multiplication  $\rightarrow$  not index multiplication  
 $A = 3 \times 4$        $B$  if  $m \times n$        $C_{\text{cols}} \times n$        $(n) \times p$        $\rightarrow$  Result is  $m \times p$

$$\begin{bmatrix} \vec{2} & \vec{2} \\ \vec{4} & \vec{1} \end{bmatrix} \times \begin{bmatrix} 2 & 3 \\ 1 & 4 \end{bmatrix} = \begin{bmatrix} ? & ? & ? \\ ? & ? & ? \end{bmatrix}$$

use in P&B  
 $\bullet$  Transpose  $A = \begin{bmatrix} 2 & 2 \\ 4 & 1 \end{bmatrix} \Rightarrow A^T = \begin{bmatrix} 2 & 4 \\ 2 & 1 \end{bmatrix}$

use in PCA  
 $\bullet$  Determinant of Matrix

$$\hookrightarrow A = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$$

$$\text{Det}(A) = a \begin{vmatrix} e & f \\ g & h \end{vmatrix} - b \begin{vmatrix} d & f \\ g & i \end{vmatrix} + c \begin{vmatrix} d & e \\ g & h \end{vmatrix}$$

$$\Rightarrow aei + bfg + cdh - afh - bdi -cef$$

$$/ a \times (ei - fh)$$

in Python (`det(A)`)

(11)

### Inverse Matrix

$$2 \rightarrow \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad \begin{pmatrix} A \cdot A^{-1} = 0 \end{pmatrix} \text{ (whole)}$$

$$3 \rightarrow \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad \begin{pmatrix} \text{for } 2 \times 2 \text{ (square matrix only)} \\ 2 \times 2, 3 \times 3, 4 \times 4 \end{pmatrix}$$

$$\rightarrow \text{Inv}(A) = A^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

~~Transpose of A~~  
Adjoint of A

10M up 8  
again come to same

$\rightarrow$  for Order:  $3 \times 3$

10M  $\rightarrow$  what done is 0  
 $\rightarrow$  but det(A) = 0  
20M

$$\therefore A = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$$

order 3 or above  
• the  $\text{Inv}(A) = \frac{1}{\det(A)}$

$$\text{Inv}(A) = \frac{1}{\det(A)} \begin{vmatrix} |ef| & |df| & |de| \\ |hi| & |gi| & |gh| \\ |bc| & |ac| & |ab| \end{vmatrix}$$

$$\begin{vmatrix} |bc| & |ac| & |ab| \\ |hi| & |gi| & |gh| \\ |ef| & |df| & |de| \end{vmatrix}$$

$$\textcircled{Q} \quad A = \begin{bmatrix} 4 & 7 \\ 2 & 6 \end{bmatrix} = A^{-1}$$

$$= \frac{1}{4 \times 6 - 7 \times 2} \begin{bmatrix} 6 & -7 \\ -2 & 4 \end{bmatrix}$$

$$= \frac{1}{10} \begin{bmatrix} 6 & -7 \\ -2 & 4 \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} 0.6 & -0.7 \\ -0.2 & 0.4 \end{bmatrix}$$

$$\Rightarrow A \cdot A^{-1} = I \quad \xrightarrow{\text{order matrix}} \begin{bmatrix} 4 & 7 \\ 2 & 6 \end{bmatrix} \cdot \begin{bmatrix} 0.6 & -0.7 \\ -0.2 & 0.4 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

### Relationship

Vector  $\xrightarrow{\text{Convert into matrix}}$  Matrix

$$\bullet v_1 = 3i + 4j \rightarrow [3 \ 4]$$

$$v_2 = 2i + 3j \rightarrow [1 \ 3]$$

Z Score  $\rightarrow$  Standardized (using matrix)

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} s_y & 0 \\ 0 & s_x \end{bmatrix}$$

scale up  
scaledown

### Rotation:

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

$\rightarrow$  decompose matrix

\* Vector to matrix calculation  
np. linear model

$$2x + y - 2 = 2$$

$$x + 3y + 2z = 1$$

$$x + y + z = 2$$

$$\begin{bmatrix} 2 & 1 & -1 & 2 \\ 1 & 3 & 2 & 1 \\ 1 & 1 & 1 & 2 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix}$$

- approach  
 (1) Row echelon methods  
 → Column echelon methods  
 (2) Inverse method

(12)

Version 2  
Bray ✓

# CRISP DM framework

Cross industry standard process for Data Mining

- ⑧ ML → Can you predict this item demand for next 6 months

- 1st do POC (Proof of concept)

1st do EDA → univariate, bivariate, MultiVariate  
 Obj measure

## ① Data Preparation

- ① Missing Value treat
- ② Duplicates value
- ③ Outlier → Class imbalance  
feature engineering
  - ① Select relevant feature
  - ② Create new feature → speed, ~~DOB~~, new attr
  - ③ Modify existing feature

Age	years remain
15 year	
16 year	
15	

## ⑤ Data Encoding

String → number → ~~onehot~~

## ① Missing Values.

- ① Missing completely random
- ② MPR (Missing not at random)  
salary

## → Deals

- ① Connect bus
- ② missing val < 1% → drop

- ③ small data → Do Imputation

Numerical Categorical  
 outliers → treated → ~~is done~~ → ~~not do~~ → Mode  
 mean median

- ④ PR (constant mean)

↳ (Replace with 0) (Mean → Mean → So 0) → 0 or -1

- ⑤ Create new columns flags

| Null | my | N db  
| 1 | 2 | 0  
| 0 | 3 | 1  
| 1 | 0 | 0  
| 0 | 1 | 1  
| 1 | 0 | 0

OR No, that is not possible in that column

+ Explain the ML Pipeline

+ Explain Lifecycle of ML

④

- ⑥ if Missing Value in Column > 30%

→ drop the Column, → add Noise  
 → impute previous / next / Avg / Med → unwanted data which is not naturally 0

Interpolation

(13)

## ⇒ Data Interpolation

- It is a process of estimating data in a range or getting unknown values from known values.

⇒ 3 Way

→ filling the data using same trend with help of interpolation

(1) Linear Interpolation

(2) Cubic

(3) Polynomial

 $x = [1, 2, 3, 4, 5]$  $y = [1, 3, 5, 7, 9]$ 

(known &amp; scattered)

 $x\_new = np.linspace(1, 5, 10)$  $y\_interp = np.interp(x\_new, x, y)$ (⇒ Scatter: ( $x\_new$ ,  $y\_interp$ ))

Cubic

OR

use from `scipy.interpolate import interp1d` $f = interp1d(x, y, kind='cubic')$  $x\_new = np.linspace(1, 5, 10)$  $y\_interp = f(x\_new)$ 

- Use Case → Line Spline
- missing data → ① scattered → ② based on trend of data
  - when not enough data
- ↳ your airplane create

Q:

$$\text{line } ax + b = 0$$

$$\text{quad } ax^2 + bx + c = 0$$

$$\text{cub } ax^3 + bx^2 + cx + d = 0$$

Polynomial

$$P = np.polyfit(x, y, 2)$$

→  $x\_new = np.linspace(1, 5, 10)$   
 $y\_imp = np.polyval(P, x\_new)$   
 P.H.Salt: ( $x\_new$ ,  $y\_interp$ )

## Outliers Treatment

① Dropping the outliers

distplot, boxplot

② Capping → replacing outliers with the nearest values that is not

③ replace with mean &amp; median

④ Scaling and transformation

or

 $\text{Q1} := df['Salary'].quantile(0.25)$  $\text{IQR} = Q3 - Q1$  $\text{lower\_fence} = Q1 - 1.5 \times IQR$  $\text{upper\_fence} = Q3 + 1.5 \times IQR$ 

Every thin box is not outlier

↳  $\text{df}[\text{df['Salary'} >= \text{lower\_fence}) \& (\text{df['Salary'} <= \text{upper\_fence})]$ -  $\text{df['Salary\_mean\_impute']} = np.where$ [( $\text{df['Salary'} > \text{upper\_fence})$  |( $\text{df['Salary'} < \text{lower\_fence})$ ),  $\text{df['Salary'].mean()}$ ,

df['Salary'])

CAP  
↳ deleted cap. $\text{lower\_cap} = df['Salary'].quantile(0.05)$  $\text{upper\_cap} = df[1, 0.95]$

(11)

- Eigen Value & Eigen Vector → useful in PCA → for reduce dimensionality reduction techniques.
- Eigen vectors do not change direction if some transformation is applied to them

→ for any square matrix  $A$ , if there is a matrix  $V$  such that if  $V$  multiply matrix  $V$  with a scalar  $\lambda$  that

$$AV = \lambda V \text{ holds true}$$

→  $\lambda$  is Eigen Value  
 $V$  Eigen Value

$V = v$   
 so direction  
 not change  
 only change  
 magnitude

$$g. \begin{bmatrix} 3 & 6 \\ 5 & 4 \end{bmatrix} \cdot V = \lambda V$$

$$V(A - \lambda I) = 0$$

$$V \left[ \begin{bmatrix} 3 & 6 \\ 5 & 4 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right] = 0$$

$$V \left( \begin{bmatrix} 3 & 6 \\ 5 & 4 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \right) = 0$$

$$V \begin{pmatrix} 3-\lambda & 6 \\ 5 & 4-\lambda \end{pmatrix} = 0$$

$$\text{det}(A - \lambda I) = 0$$

$$\sqrt{(3-\lambda)(4-\lambda) - 30} = 0$$

$$\sqrt{(3(4-\lambda) - \lambda(3-\lambda) - 30)} = 0$$

$$\sqrt{12 - 3\lambda - 4\lambda + \lambda^2 - 30} = 0$$

$\lambda = 9, -2$  so there is 2 Eigen Value

$$\begin{cases} \lambda_1 = 9 \\ \lambda_2 = -2 \end{cases}$$

PCA

$$\text{by } \begin{bmatrix} \text{cov}(a,a) & \text{cov}(a,b) \\ \text{cov}(b,a) & \text{cov}(b,b) \end{bmatrix}$$

$$\begin{bmatrix} \lambda_1 & \lambda_2 \\ \lambda_2 & \lambda_1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\begin{bmatrix} \lambda_1 & \lambda_2 \\ \lambda_2 & \lambda_1 \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$$

$$\begin{bmatrix} \lambda_1 & \lambda_2 \\ \lambda_2 & \lambda_1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \lambda_1 v_1 + \lambda_2 v_2$$

$$\begin{bmatrix} \lambda_1 & \lambda_2 \\ \lambda_2 & \lambda_1 \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} = \lambda_1 f_1 + \lambda_2 f_2$$

ex

$\begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix}$	$\begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$	$\begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \lambda_1 v_1 + \lambda_2 v_2$	$\lambda_1 = 9$	$\lambda_2 = -2$	$\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$
$\begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix}$	$\begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$	$\begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} = \lambda_1 f_1 + \lambda_2 f_2$	$\lambda_1 = 9$	$\lambda_2 = -2$	$\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$PC_1 = \lambda_1 f_1 + \lambda_2 f_2$  new data point with PCA

(15)

- Imbalance class
- feature extraction
- Engineering,
- Scaling & Data
- scaling
- Practical
- fdf in FCA
- simple linear
- Regression

→ Imbalance

for classification data  
 90% 80% 1000

10 distinct { what  
 1%, 2%, 3%, 10% } Minority = 90%. Among

→ Major, Minor

→ types

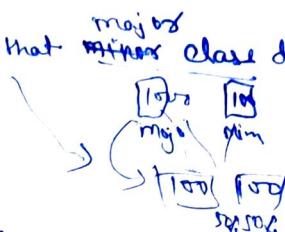
(1) Under Sampling

? Resample by that ~~major~~ <sup>major</sup> class data.

80

① DisAdv  
② Data loss

(2) Over sampling



(3) SMOTE

→ from sklearn.utils import resample

? resample(df\_minority, replace=True, n\_samples=len(df\_majority), random\_state=1)

• for sklearn.datasets import make\_classification

 $x, y = \text{make\_classification}(\text{n\_samples}=1000, \text{n\_redundant}=0,$   
 $\text{n\_features}=2, \text{n\_clusters\_per\_class}=1,$   
 $\text{weights}=[0.9, 0.1], \text{random\_state}=1)$ 

→ Pip install imblearn

• from imblearn.over\_sampling import SMOTE  
 $\rightarrow \text{oversample} = \text{SMOTE}()$  $x, y = \text{oversample.fit_resample}(\text{df\_final}[[\text{'F1'}, \text{'F2'}]], \text{df\_final}[\text{'target}])$ 

## feature extraction

Not All feature to use in Model creation

① Create new feature

② Modify the existing feature

 $\begin{matrix} 15+ \\ 16+ \\ 18+ \end{matrix} \rightarrow \begin{matrix} 15 \\ 18 \\ 16 \end{matrix}$  } Date | Age / Month / Year
 

distance	time	speed
50	2	25
100	3.5	28.57
=	=	=
=	=	=

↳ with increase in no of features -

- ① Model train cost ↑
- ②

(optional) feature Scaling (PCA → Interpretation easier) → less time in computation

why  
 same  
 scale

Age	No P
20	2
25	3
30	1
250	0

$$\begin{aligned} & \text{Interpretation easier} \\ & \text{less time in computation} \\ & \text{③ gradient descent optimization} \\ & \text{④ PCA} \\ & \text{⑤ matrix moving Avg} \\ & 2 \times 3 = 6 \\ & 5 \times 2 = 10 \\ & 3 \times 5 = 15 \\ & 250 \times 3 = ? \quad \rightarrow \text{more Compute} \end{aligned}$$

Types of scaling → in ML

① Standardisation (ML alg)

$$Z = \frac{x - \bar{x}}{\sigma} \rightarrow \text{SD} \rightarrow \text{all} = 0 \quad \sigma = 1$$

② Normalization (min-max scale) → DC Algo

$$X_{\text{scale}} = \frac{x_i - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}} \quad [0, 1]$$

③ Unit Vector

$$(2, 3) \text{ on } \sqrt{2^2 + 3^2} = \sqrt{13}$$

$$= \sqrt{13} = \sqrt{13}$$

$$\text{unit} \rightarrow \frac{1}{\sqrt{13}} \cdot \begin{pmatrix} 2 \\ 3 \end{pmatrix}$$

→ Select the right feature

(1) Filter method

(2) Embed method → Lasso regression (for feature selection)

(3) Wrapper method → Recursive feature elimination  
forward selection      backward selection

$f_1, f_2, f_3, f_4$  | f5, f6

(1) Correlation

$$(f_i, y)$$

(2) Multicollinearity (VIF)

$$\lambda_i = \frac{1}{\sum_{j=1}^n R_{ij}^2}$$

correlation among features  
correlated

in turn

### Data Encoding

(1) Nominal / Ordinal / Dummy variable → Categorical to Numerical

① No order in the data

disadv  
• Column has only  
goodcols

(2) Labelled and ordered - 1, 2, 3, 4, 5 → Discrete → ML can assume order

(3) Target Guided ordinal encoding

useful for level of teaching  
nominal higher categories

↳ select top 5 as top 10 unique categories in feature, if no y

Sol: select top 5

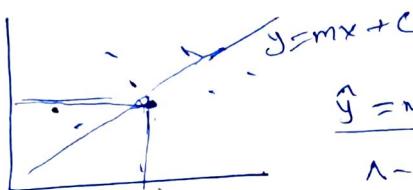
# Linear Regression

17

- Intuition
- Mathematical
- Implementation

→ Simple LR

ML  
 SL / OSL  
 Reg class  $y_{\text{not}}$   
 Conf  $y_{\text{du}}$



$$\hat{y} = mx + c \quad \text{or} \quad \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$\hat{\beta}_0$  - estimates

$$f(y) = mx + c$$

$m$  Predict coefficient  
 $c$  Intercept

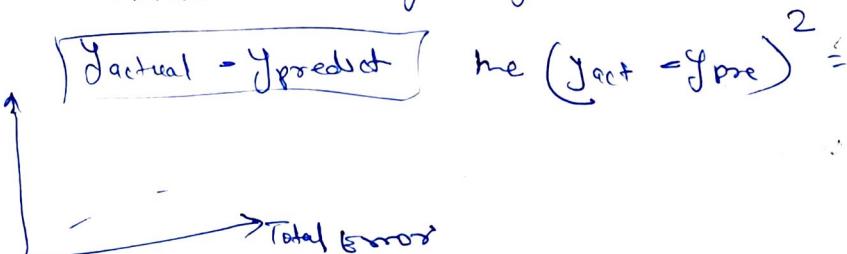
$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (\hat{\beta}_1 = m, \hat{\beta}_0 = c)$$

$$h_0(x) = \theta_0 + \theta_1 x$$

$$\text{or } h_0(x) = b + w x$$



Q Which line will give you least error??



Ans

So we want to have that combination of  $(m, c)$  /  $(\theta_0, \theta_1, b, w)$  /  $b_0, b_1$  where the overall error is least and the line will be best line.

So we say

$$\min \sum_{i=1}^n (y_i - y_{\text{pred}})^2$$

and this is

called OLS (Ordinary Least Squared Error)

↳ This is the simplest / most formulate initial formula to find best fit line.

$$\min \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}})^2 \rightarrow \text{min the sum of square diff for all } y_{\text{act}} \text{ & } y_{\text{pred}} \text{ for all data points.}$$

$$\min \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

Optimal  $(\hat{\beta}_0, \hat{\beta}_1)$  → best fit line

Q Which technique is to find best fit line in LR  
OLS

(18)

Now if i divide OLS

total no of data Point  
n

⇒

then it called

Mean Square Error  
MSE

$$\text{Min } \frac{1}{n} \sum (y_i - (B_0 + B_1 x_i))^2$$

so High MSE  
High cost↑

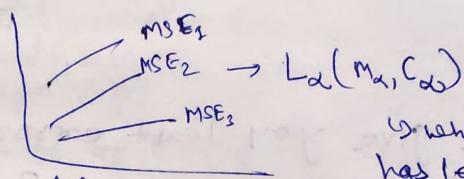
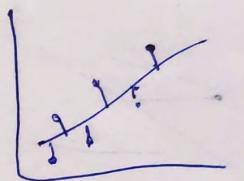
~~Cost~~  
Error

[Cost function]

Cost function :- A function/metric/Expression that ensures to penalize you whenever you do something wrong (more error)

① MSE

$$\rightarrow \frac{\text{SSE}}{n} = \frac{\text{OLS}}{n} = \frac{1}{n} \sum_{i=1}^n (y_{\text{act}} - (B_0 + B_1 x_i))^2$$



whichever has least MSE

so we optimize the process

↓ How

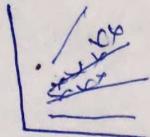
Gradient Descent

MSE  
outliers  
if outliers not treated  
use RMSE

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

$$\text{OLS} \rightarrow \text{SSE} \rightarrow \frac{\text{SSE}}{n} \rightarrow \text{Avg(MSE)} \approx \sqrt{\text{MSE}} \rightarrow \text{RMSE}$$

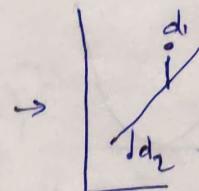
$R^2$   
Cost function



best fit line

(how)

minimize the error

(best representative of  
all the data points)

$E_{\text{total}} = d_1 + d_2$

(should be minimum)

$\text{Total Error} = \sum_{i=1}^n E_i$

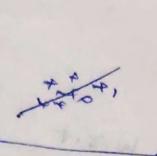
how calculate

$\sum_{i=1}^n (Y_{\text{act}} - Y_{\text{pred}})^2$

for Least errors

$= E(y_i - Mx_i - c)^2$

$\text{or } \beta_0, \beta_1$   
 $w, b \rightarrow D_L$   
 $\theta_0, \theta_1$



for Least Errors:  $E = \sum_{i=1}^n (y_i - Mx_i - c)^2$

$$\text{also write like } J(M, c) = \sum_{i=1}^n (y_i - Mx_i - c)^2$$

↓  $(\beta_0, \beta_1)$   
↓  $(\theta_0, \theta_1)$   
↓  $(e_0, e_1)$   
↓  $(w, b)$

So we solve by 2 approaches

Close form  
solnIterative form  
solution
 $\downarrow$   
Gradient  
Descent

→ formulate a problem

stated into math  
equation

like  $ay^2 + bx + c = 0$

$$d, f = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

# Total Error =  $\sum (y_i - Mx_i - c)^2$

$$\begin{cases} M \\ c \end{cases}$$

Scenario 1

$c=0$

$y=Mx$

If you change  
slope the  
error increase  
& Decrs

Scenario 2

$M=1$

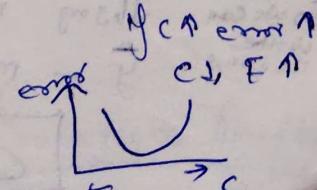
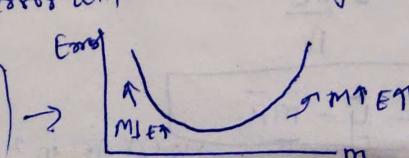
$y=x+c$

parallel lines

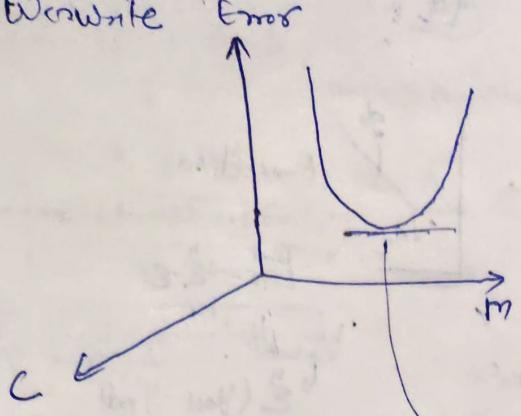
So observe.

→ if you increase M or  
decrease M error will  
change

We say  
Stop ↑ T. Error ↑  
Stop ↓ T. Error ↓



Wronwrite



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

$$J = n$$

$$y = x^2 = \text{parabola}$$

because we want least error

$$\begin{aligned} \text{Slope} &= 0 \rightarrow \text{so we do } \frac{\partial E}{\partial m} = 0 \text{ w.r.t } \\ \text{error} & \text{w.r.t } m \text{ & } c \\ \text{so slope} &= 0 \end{aligned}$$

Differentiation

$$\frac{\partial E}{\partial c} = 0$$

why equal to 0?

Error should be  
least  $\rightarrow$  so Slope  
of error = 0

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

w.r.t c

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

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

$$\Rightarrow \frac{\partial E}{\partial c} = 0 \Rightarrow$$

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

[divide by -2]

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

divide by n (no. of data points)

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

So we can  
write like this

$$\bar{y} - m\bar{x} - \frac{n}{2}c = 0$$

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

What is m?

$$nx^{n-1} \leftarrow x^n$$

$$\frac{dy}{dc} = \frac{\partial y}{\partial c} \frac{dx}{dc}$$

Constant factor  
is zero

(21)

So make equations for m differentiable

$$\frac{\partial E}{\partial m} = 0$$

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

$\bar{y} - m\bar{x} \rightarrow \text{we find}$

$$\hookrightarrow E = \sum (y_i - mx_i - \bar{y} + m\bar{x})^2$$

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

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

Solve this in any way  
divide by 2 both side

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

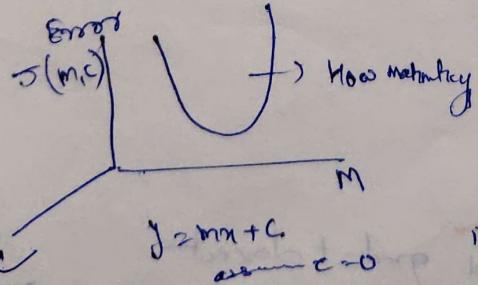
## If multiple variable

$x_1, x_2, x_3, x_4$

$$\frac{\partial E}{\partial x_1}, \frac{\partial E}{\partial x_2}, \frac{\partial E}{\partial x_3}, \dots, \frac{\partial E}{\partial x_n}$$

then more complex

Iterative solution



(Gradient descent)

$$\rightarrow J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{pred}})^2$$

$$\begin{array}{c|c|c} & \theta_0 = mn \\ \theta_1 & y & y \\ \hline 1 & | & | \\ 2 & | & | \\ 3 & | & | \end{array}$$

→ Initialize random m,c.

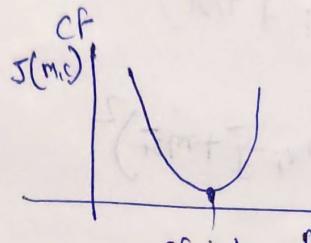
→ move in the direction.

→ least error will be optimal m,c.

}, How do by Convergence Algorithm

## \* Convergence Algorithm

- keep making new best fit line in the direction of dp until Error (CF) is reduced as compared to previous line.



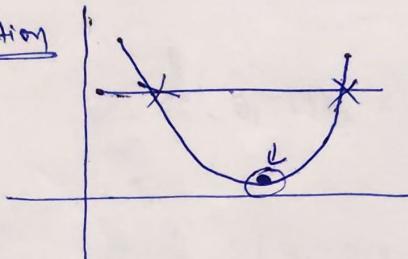
Point of minima  $\rightarrow$  global minima

$\downarrow$   
Till when ??

$\downarrow$   
where  $m, c$  is optimal

Error is least  
CF is least  
Cost function

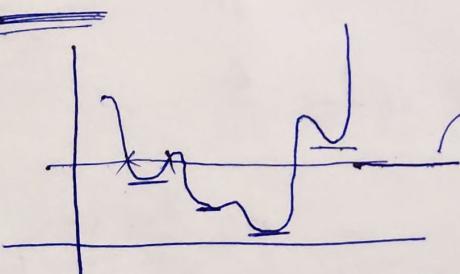
function



Convex function

$\rightarrow$  Cuts at only two points

Convex  
one minima



non-convex fn

Cuts more than two points

local, global minima

Pt - 3 → Global minima

Pt - 1, 2, 3 → local minima



$\Rightarrow$  Convergence Algorithms

Repetet until convergence

Cellol gradient descent



$$m_{\text{new}} = m_{\text{old}} - \eta \frac{\partial J(m,c)}{\partial m}$$

$$c_{\text{new}} = c_{\text{old}} - \eta \frac{\partial J(m,c)}{\partial c}$$

for this  
Time ↑  
Content ↑

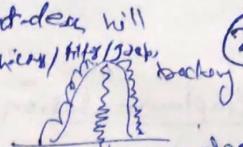
4th May 2:00

<http://blog.deeplearning4j.org/gradient-descent-will-developers-go-to/mkl-crashcourse/>

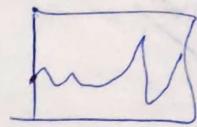
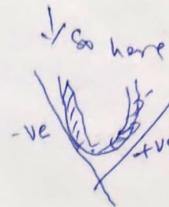
(23)

## Gradient Descent

Slope Coming down



→ sideset the T in Deep learning  
slope



Saddle Point

### Slope guides for Convergence

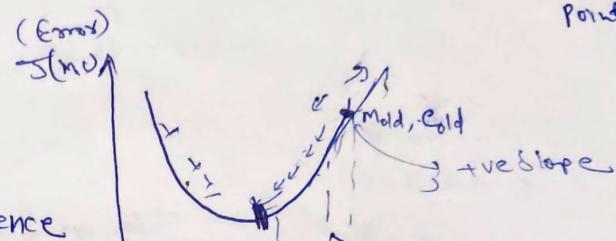
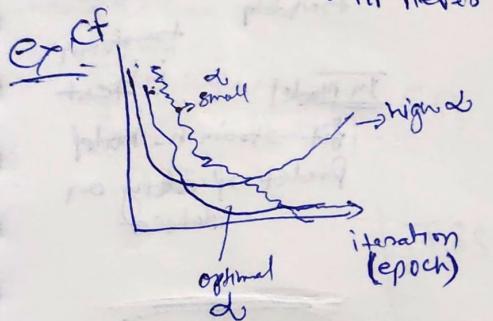
$$\frac{\partial f}{\partial \theta} = \eta \frac{\partial^2 f}{\partial \theta^2}$$

etc  
it is  $(0.1 - 0.001)$

learning rate  $\rightarrow$  It decides the convergence speed.

$\eta$  - too small  $\rightarrow$  Very slow for Convergence

$\eta$  - too high  $\rightarrow$  Exploding gradient, problem  
 $\downarrow$  or  
will never converge



$$\begin{aligned} \text{Mold} &= 5 \\ \text{Mold} &= \text{Mold} - (\text{some value}) \\ &= 5 - 1 \end{aligned}$$

so you  $\downarrow \eta$

what?

$$\frac{\partial C_F}{\partial M_{old}}$$

$$\begin{aligned} \text{if } -\text{ve slope} \\ M_{new} &= M_{old} - \frac{\partial C_F}{\partial M_{old}} \\ &= M_{old} - (-\text{ve slope}) \\ M_{new} &< M_{old} \end{aligned}$$

Cantors plot

$$\begin{aligned} M_{new} &= M_{old} + \text{ve value} \end{aligned}$$

(24)

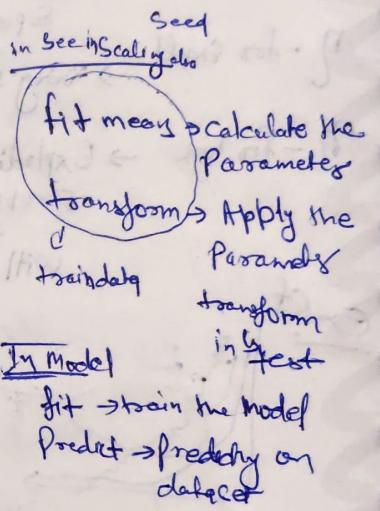
## → Linear Regression Implementation

Mark

- ML pipeline
- Read dataset
- Preprocess data
- clf →  $X, y$
- Train-test
- Scaling
- Model training
- Model evaluation
- Model prediction

(random\_state = 111)

reproduce the same result



evaluation

①  $R^2$

② RMSE & Root Mean Square Error

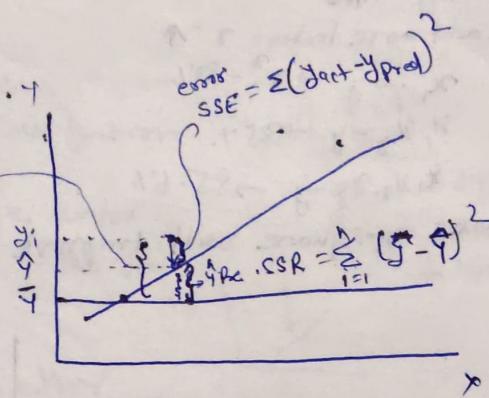
- ③ MSE
- ④ RMSE
- ⑤ MAE

## Evaluation Metric (Regression)

- ① R square
- ② Adjusted R-square
- ③ MSE
- ④ RMSE
- ⑤ MAE

### R-SQUARE

Total Sum of  
square =  $\sum (y_i - \bar{y})^2$   
TSS



$$\begin{aligned} \bar{y} &= \text{avg of } y \\ \hat{y} &= \text{Predicting } y \\ y_i &= \text{Actual } y \end{aligned}$$

Find  
R-square My Model is  
better.

→ can R square -ve?  
→ Yes but why

if y very far

$$\frac{1 - \frac{SSE}{TSS}}$$

Sum of square due to regression

SSR → is the  $\sum$  of square diff of ~~Actual~~  
~~y - Avg y - Predicted~~

SSE = is the  $\sum$  of square diff of  $y$  Actual  
 $y_i - y$  Predicted,

TSS = is the  $\sum$  of square diff of \*  
 $y_i$  actual -  $y$  Avg

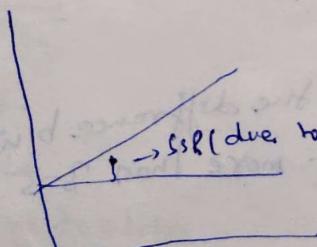
if  $R^2 = 85\%$   
so we say  
85% Variance  
will explain by  
this feature x.

$$\frac{SSE}{n} = \text{MSE}$$

Rsquare =  $1 - \frac{SSE}{TSS}$  or  $\frac{SSR}{TSS}$

also called  
Coefficient of  
determination →

out of total error,  
SSR is the variation explained  
by linear regression line.



Can I say this  
much of Variation is  
 $y$  is explained by  
Regression

R-Square =  $\frac{SSR}{TSS}$ , out of total variation, SSR is variation explained.

$$\% = \frac{SSR}{TSS} \times 100$$

→ The % variation in  $y$  explained by  $x$ .

Now Explain Variance  
Unexplain Variance

(26)

## ② Adjusted $r^2$

We know  $\rightarrow r^2$  $\Rightarrow \%$  age explained variance in  $y$  due to  $X$ .

Example: House price pred.

 $x_1, x_2, x_3, y$  $x_1, x_2, x_3, \dots, x_{100}$ add more features  $r^2 \uparrow$  $x_1 \rightarrow y \rightarrow r^2 - 80\%$  $x_1, x_2 \rightarrow y \rightarrow 85\%$  $x_1, x_2, x_3 \rightarrow y \rightarrow 85.6\%$ 

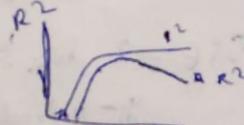
\* As we add more features  $r^2$ -square will improve / or remain constant

if  $x_1, x_2, \dots, x_{100}$ 

then how to know which feature is not contribute much.

(So use)

## Adjusted $r^2$ -square

 $\hookrightarrow J+$  penalizes  $r^2$  as we add new features.

$$A_r^2 = \frac{1 - (1 - R^2)(N-1)}{N-P-1}$$

 $N = \text{No of df}$  $P = \text{No of features}$ 

Example

$$\text{seen no } 1 \quad R^2 = 87, N = 11 \\ P = 2$$

$$R^2 = 87, N = 11, P = 8$$

$$A_r^2 = \frac{1 - (0.2)(10)}{11 - 8 - 1}$$

$$= 1 - \frac{2}{2} \Rightarrow 1 - 1 = 0$$

$$\text{Q } A_r^2 = \frac{1 - (1 - 0.8)(11-1)}{11-2-2} = \frac{1 - 0.2 \times 9}{8} = 0.75$$

$$\text{Q } A_r^2 < r^2$$

④ only add feature in the model if the difference b/w  $r^2$  square & adjusted  $r^2$  square is not more than 3-5%.

(27)

### (3) MSE (mean Square ERROR)

$$\frac{1}{n} \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}})^2$$

- ④ Can we use another  
Cost function  
yes

① Error is quantified

② Lower the MSE better the model will be

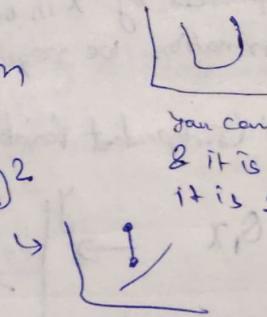
#### Advantage

① It is differentiable  
so it is used as Cost function

② Emphasis on large errors  
 $(y_{\text{act}} - y_{\text{pred}})^2$

#### Disadvantage

- ① Not robust to outliers  
② It is not in same unit as y



You can differentiate any point  
& it is convex so  
it is ↓ minimal.

$$\begin{aligned} y_{\text{act}} &= 3m \\ y_{\text{pred}} &= 6m \\ \text{err} &= (3-6)^2 \\ &= 9m^2 \end{aligned}$$

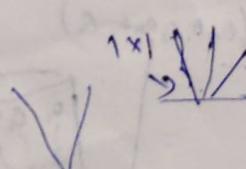
### (4) MAE

$$| = \frac{1}{n} \sum_{i=1}^n |y_{\text{act}} - y_{\text{pred}}|$$

#### Adv

- ↳ Less sensitive to outliers
- ↳ More interpretable
- ↳ It is of same unit

$$\left( \frac{1}{n} \sum_{i=1}^n |y_{\text{act}} - y_{\text{pred}}| \right)$$



#### DisAdv

- Not differentiable at x=0
- time consuming

→ Lower the MAE Better will be the model

### (5) RMSE (Root mean Square Errors)

$$\text{RMSE} = \sqrt{\text{MSE}}$$

#### Adv

- ① Same unit
- ② Differentiable
- ③ less sensitive to outliers

## Multi Linear Regression

(28)

agenda

→ Assumption of LR

→ MLR

→ Polynomial reg

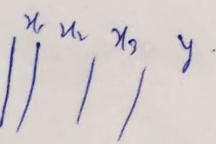
→ Lasso, Ridge, Elastic net

→ Cross Validation

→ Hyperparameters

→ Logistic regression

Cost function



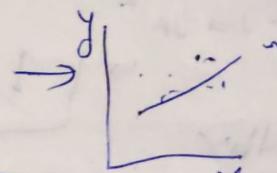
Why we go MLR

- All the information  $y$  is not captured by  $x$  in order to capture more & more information we require more  $x$ .

MLR → More than 1 IV (Independent Variable) is used.

1-IV

$$y_{\text{pre}} = \theta_0 + \theta_1 x$$



→ intercept will only change

MVR

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

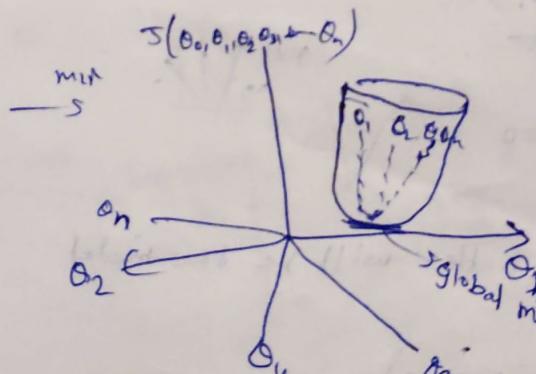
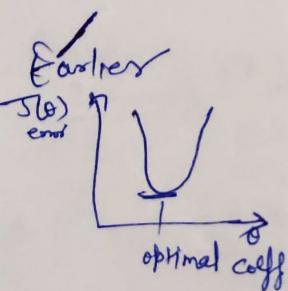
So what is  $\theta_0, \dots, \theta_n$  → are optimal coefficients

Cost function of this

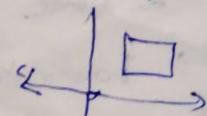
$$CF = \frac{1}{n} \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}})^2$$

$$(y_{\text{act}} - (\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n))^2$$

MLR Visuals



It is Plane Hyperplane



In MLR, the Intercept is the point where hyperplane intersects the vertical axis.

eg marks obtain by ast even if he didn't study

$$y = 50 + 2.5x$$

$\rightarrow \theta_0 \Rightarrow$  with limit ↑

in # of hours studied, the marks increase by 2.5 units

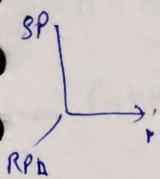
$$y_{\text{pred}} = 70 + 3.2(\text{Age})$$

Example

① # no of hours | Marks

② Selection of year | Age

Then how to interpret in MLR



$$Y_{pred} = \theta_0 + \theta_1 X_1 + \theta_2 X_2$$

Mathematically change but fundamentally same

$\theta_1 (0.52) \rightarrow$  with 1 unit increase in age of car

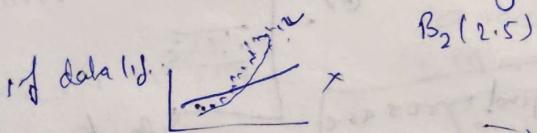
the SP of car decreased by 0.5 units

units on Avg 'keeping RPM Constant'

means 0

Whatever you taking others are constant

7x13/25



$$B_1(0.52)$$

## ■ Polynomial Regression :

① Simple polynomial regression (1 Dependent Var (DV))  
1 Independent Var (IV)

- H.W. & test
- System
- Explain/Unexplain
- Variance  $R^2$
- Stats model
- Feature Selection

Polynomial degree 0,  $h_0(x) = \theta_0 x^0$

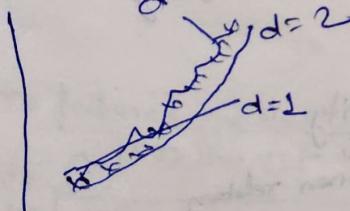
$$\therefore 1 \rightarrow h_0(x) = \theta_0 x^0 + \theta_1 x_1^1 \rightarrow \text{Simple LR}$$

$$\therefore 2 \rightarrow h_0(x) = \theta_0 x^0 + \theta_1 x_1 + \theta_2 x_2^2$$

$$\therefore 3 \rightarrow h_0(x) = \theta_0 x^0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3$$

$$\vdots \quad \vdots \quad \vdots \quad \vdots$$

$$d=3 - \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$



It captures the non-linear relationship

\* As you increase the degree, you might get an overfitting model.

# Non-linear relationship

→ also other model  
Decision Tree  
non-linear

PR → Disadvantage Overfitting the model

## For multiple X

Polynomial degree 2:  $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 +$

$$\theta_4 x_1^2 + \theta_5 x_2^2 + \theta_6 x_3^2 + \theta_7 x_1 x_2$$

$$+ \theta_8 x_1 x_3 + \theta_9 x_2 x_3$$

$x_1, x_2$

↓  
choose product  
keeping power min

(2)

## # Assumption of Linear Regression

① Linearity  $\rightarrow$  X and Y should have linear relationship.

② Independence  $\rightarrow$  Observations (rows) are independent of each other.

③ Homoscedasticity  $\rightarrow$  Also known as Constant Variances. The variances of errors are constant ex: Durbin-Watson test  $\rightarrow$  to find errors are homosced. or not.

Opposite  $\rightarrow$  Heteroscedasticity

Durbin-Watson test  $\rightarrow$  to find errors are homosced. or not.

Homosced. is not variance.

④ Normality of error: - errors should be normally distributed.

$f_1, f_2, f_3 \rightarrow y$   
not among points

⑤ Whether part  $u_1$  is contributing to y & what part  $u_2$  is contributing to y  
Can we say  $u_1 = u_2$   $\downarrow$  not answer

$$8x_1 + 2u_2 = y$$

$$10x_1 = y$$

$$2u_1 + 8u_2 = y$$

$$10u_2 = y$$

\* so we say multi-collinearity  
may together linear relation

$u_1 \rightarrow x_2$   
 $u_1 \rightarrow y$  } Correlation

$$x_1 \approx (u_1, u_2)$$

$$u_1 \approx (u_1, u_2, u_3) \quad \{ \text{multi-collinearity}$$

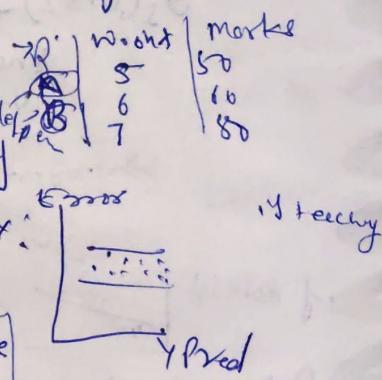
$\rightarrow$  where a feature exhibits a linear relationship with more than one variable

$$u_1 \sim u_2, u_3, \dots$$

$$u_5 \sim u_1, u_2, u_3, \dots, u_4$$

(What)  $\rightarrow$  One feature explained by other feature.

(Why)  $\rightarrow$  You can not interpret correctly what is the contribution of each individual feature w.r.t y



31

you only lose  
interpretationit is severe problem  $\rightarrow$  Ans NO  
severe $\hookrightarrow$  No effect on predictionThen Problem  $\rightarrow$  you lose interpretability

$$\textcircled{1} \quad \begin{matrix} x_1 & x_2 & x_3 & x_4 & y \\ x & \end{matrix}$$

 $\textcircled{2}$  Computationally expensive  $\rightarrow$  model training
Solution 1  $\textcircled{1}$  Detect Multi-Collinearity or Measure by

VIF (Variance inflation factor)

a measure of amount of multi-collinearity  
in regression

$$\boxed{VIF_i = \frac{1}{1-R_i^2}}$$

 $R^2 \rightarrow$  % Variation in y  
explained by XHow  $x_1, x_2, x_3, x_4$ y  $\sim x_1, x_2, x_3, x_4$  $x_2 \sim x_1, x_3, x_4$  $x_3 \sim x_1, x_2, x_4$  $x_4 \sim x_1, x_2, x_3$ The  $x_j$  is treated as y & calculate the R-square

$$y_j \sim (x_1, x_2, x_3, \dots, x_n)$$

VIF  $\sim 0$  to  $\infty$ VIF  $> 10$ , then  
drop feature one  
by one not  
together.So  $\rightarrow$  Since VIF doesn't impact prediction, you should always ask  
business team before dropping.

$$\boxed{2 \text{ Question}} \rightarrow \text{Why } VIF \geq 10 \rightarrow 10 = \frac{1}{1-R^2} \Rightarrow 1-R^2 = \frac{1}{10} \Rightarrow R^2 = 1 - \frac{1}{10} = \frac{9}{10} = 0.9$$

$$(x_1) \sim x_1, x_2, x_3$$

 $\hookrightarrow$  90% of variance in X is explained by  $x_2, x_3, x_4$ 

if 1000 features then how.

No  $\textcircled{1}$  1000 time VIFlibrary  $\textcircled{2}$  RFE (Recursive Feature Elimination)

Implementation

② PCA

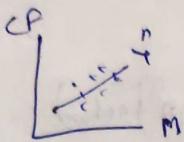
Kinds of Multi.  
• Data based  
• Structural

④

# Regularisation  $\Rightarrow$  To add something to regularise section

↓  
to regularize  
↓  
to penalize

or data  
↓  
train test } overfitting  
Acc1 Acc2  
low bias High Var



$M$  is very high  $\rightarrow$  model has memorized the dp  
+ it means overfitting

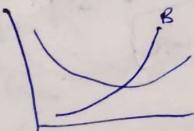
•  $R_{\text{model}}$   
  
 $y = mx + c \rightarrow M=0$  then  $y = \text{Constant}$   
 $\hookrightarrow$  No use of  $x$  in predicting  $y \rightarrow$  underfitting to  
High bias

or  
 $y = mx + c \rightarrow m = \infty$  ( $m$  is very, very large)

You are giving all importance to  $x_1 \rightarrow$  it is memorising  $x_1$   
Overfitting

So what we do

$$\frac{0}{x} \rightarrow \frac{m - \alpha}{x}$$



→ introduces some bias  
 $\downarrow$   
 error in fit

$$CF \rightarrow \underbrace{\frac{1}{n} \sum_{i=1}^n (y_{\text{act}} - mx - c)^2}_{\text{error 20}} + \underbrace{\lambda m^2}_{\text{so we add some}}$$

or introduce some bias

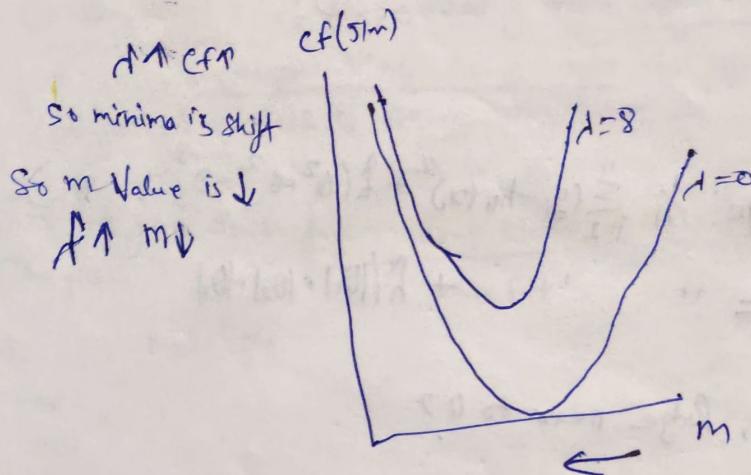
- ① Ridge  $- CF + \lambda m^2 \rightarrow (L_2 \text{ Regularization}, L_2 \text{ Norm})$  bias / Training error
- ② Lasso  $- CF + \lambda |m|$
- ③ Elastic-net  $- CF + \lambda m^2 + \lambda |m|$

~~$CF$~~  = Ridge =  $CF + \lambda m^2$

$\lambda$  decide by regularization in hyperparameters during

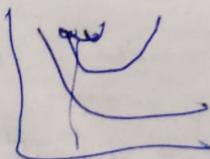
# Related

## Relationship b/w $\lambda$ , Cost function, m



## ② Lasso (L1 penalty, L1 norm)

$$\text{Lasso} = \text{cf} + R \cdot \sum_{i=1}^n |\text{slope}|$$



↳ when also decrease in  $m \downarrow$   
→ As you ↑ λ the less significant feature will be removed in Lasso →

$\lambda \uparrow \text{Cf} \uparrow$  minima shift

$R \uparrow$  slope ↓

Slope can be 0.

## Ridge (

→  $L_2$  norm / penalty

→  $\lambda \uparrow m \downarrow$

→ m will never become 0, it will be close to 0

→ Reduces the overfitting (by reducing the coefficient)

→ effective in handling multi-collinearity  
(reducing ~~coefficient~~ all the coeff.)

$x_1 \sim [x_1, x_2]$   $x_1$  is also passed  $x_2, x_3$  is also passed  
minimization happen → overfitting → Ridge reduces overfitting

Discadv → Doesn't make least important feature to 0. → reduce multicollinearity as well

Coefficient to 0. → overcome by Lasso → feature Selection

Lasso (L1/m)  
→ L1 norm / Penalty

→  $\lambda \uparrow m \downarrow$

→ m will become 0

Adv leads to feature Selection

- removing insignificant feature
- Dis Adv → leads to Automatic feature Selection / sparsity

③ Elastic net → combination → Reducing overfitting ( $L_2$ )  
→ Feature Selection ( $L_1$ )

$\lambda_{L_2}$   
Hyperparameter

→ General formulas → Ridge:  $\frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x))^2 + \lambda (\theta_1^2 + \theta_2^2 + \theta_3^2 + \dots + \theta_n^2)$   
Lasso:  $\frac{1}{n} \sum_{i=1}^n |y_i - h_{\theta}(x)| + \lambda (|\theta_1| + |\theta_2| + |\theta_3|)$

④ Why Lasso makes some 0, Ridge near to 0?

May 18 mon -

- Cross Validation & Hyper parameters study
- Logistic regression depth
- Evaluation metric
  - ↳ Practical
  - Multi class classification

# Cross Validation: Experimenting with different arrangement of same data to build different models of same algorithms.

r-square

\* When if you don't fix random state each time it will change.

↳ for

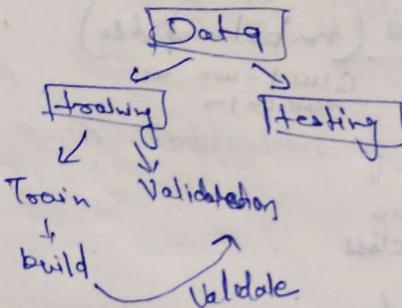
→ Exp1 - Model 1

Exp2

Exp3

...

\* Many experiments and average out the metrics increases the confidence in the model.



### Example

tooting of different need	tooting		Validation	Train	Val	r-squr	/decreas
	Exp1	Exp2					
	"	"					
	"	"					
	"	"					
	"	"					
	"	"					

↳ Exp 1 data  
not repeat in  
Exp 2 data

↳ Below arrange  
of same data  
in different  
way

### Types

#### ① Leave One Out Cross Validation (LOOCV)

Exp1	1 2 3 4	-	-	-	100	not use in training but use	→ test data is unseen data
Exp2	1 2	3 4	-	-	100	1-validate 99 - testing data	
Validation	1 2	3 4	-	-	100	2nd data for validate 99 - train data	

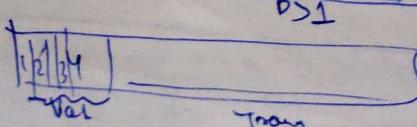
Avg - r-square [100] Mode

Adv

dis Adv -

- ① Time Complexity
- ② Model Overfitting

#### ② Leave P Out Cross Validation



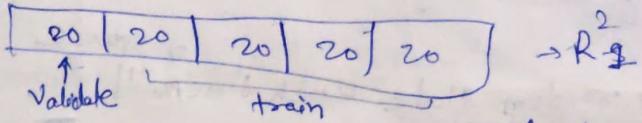
### ③ K-fold Cross Validation

$n=100 \Rightarrow k \text{ group } (k=5) \rightarrow$  No of groups that you want to divide the data

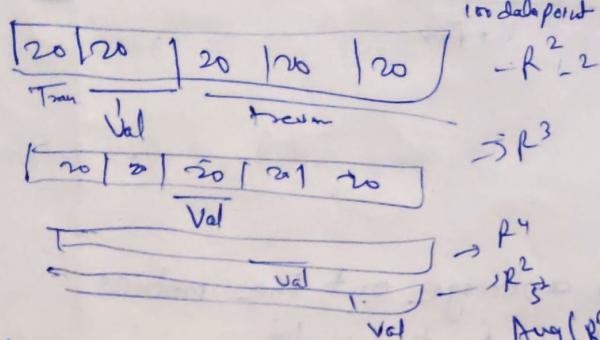
Example

$$\frac{100}{5} = 20 \text{ data point}$$

Exp-1



Exp-2



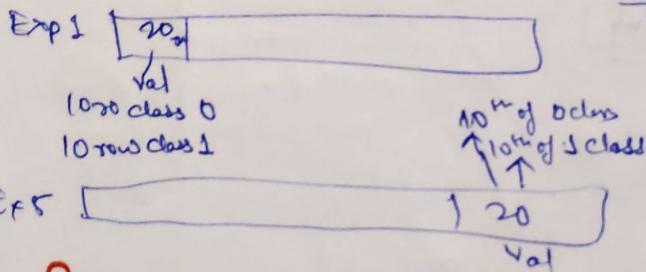
Adv

- So time complex
- Reduce

### ④ Stratified k-fold cross Validation (imbalanced data)

Ex:  $k=5, n=100$

$$\frac{\text{Class 1} = 40}{\text{Class 0} = 60}$$

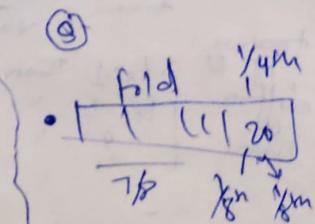


## \* Hyper Parameter

What is \*  
why necessary \*

$$\rightarrow \text{you minimize ridge}, \sum_{i=0}^n (y_i - y_{pi})^2 + h(\text{shape})^2$$

$C, M$  is learned  
from data  
(during the training  
phase)



def \*  
It's External Configuration  
of a model that are not  
learned from the data  
but are set prior to  
training process

• you don't know what value you put in

$$h = (1, 2, 4, 10, 20, 100, 1000)$$

Tune the  $\rightarrow$  best performance  $\rightarrow$  Hyperparameter tuned.

This process called Hyper parameter tuning

These parameters called Hyperparameters

like

Model

Random forest  
Impact tree  
Decision tree

Relation is hyperparameters & CV

(\*) Hyperparameters tuning with cross Validation

How → by

- ① Grid Search CV (Grid search + Cross Validation)
- ② Randomized search CV

### ① Grid Search CV

$$\lambda \rightarrow 1, 2, 3, 5, 10$$

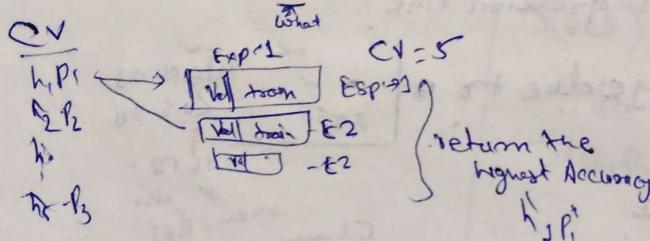
$$p = p_1, p_2, p_3$$

	$P_1$	$P_2$	$P_3$	$P_4$
$P_1$	1	2	3	4
$P_2$	5	6	7	8
$P_3$	9	10	11	12

which give high R<sup>2</sup> value

# All possible combination of hyperparameters is taken

→ So where does CV come to role → inf.

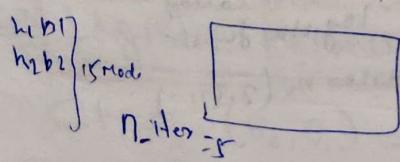


→ for each combination of  $\lambda$  &  $p$  → k-fold CV will happen

### ② Randomised Search CV

→ We will not see all possible combination.

→ Out of all the possible combination only Select some random combination



$$5 \times 5 = 25 \text{ Models}$$

Adv → Time complexity decreases

→ If data is small → instead of train test split → you can do with k-fold cross validation.

→ Lasso ridge & Elastic net → Cross Validation

18 May  $\rightarrow$  2:25

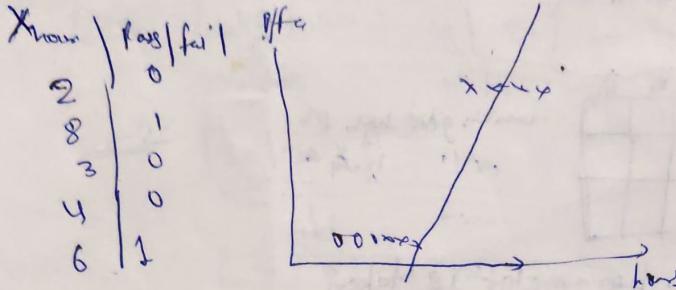
## Classification

$\rightarrow$  binary

$\rightarrow$  Multi Class

$\rightarrow$  SVM  
 $\rightarrow$  Boosted  
 $\rightarrow$  XGBoost

### ① Logistic Regression



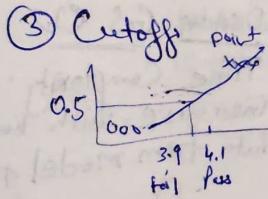
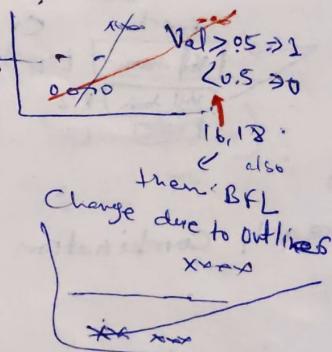
here we cannot do R  
regression line.

② Why we can not use regression line?

③ Best fit line will change due to presence of outliers.

④ Range of output ( $y \rightarrow 0 \text{ or } 1$ )

↳ But regression line goes  $-2$  to  $+2$



$\Rightarrow$  So squashing of line do to solve

2 problem

- $x \text{ too large} \rightarrow 0, 1$
- outliers

no problem.  
now

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

$\Rightarrow$  S shaped curve  
or  
Sigmoid curve

↳ this curve called  
Logistic function  
 $e^{-x}$  (2.371...)

explain:  
1. Logistic Regression

$$1 + e^x$$

$$-x \text{ and}$$

$$x = -\alpha \Rightarrow 1$$

$$1 + e^{-\alpha} = \frac{1}{1 + e^\alpha} = \frac{1}{1 + \frac{1}{\alpha}} = 0$$

$$x = \alpha \Rightarrow \frac{1}{1 + e^{-\alpha}} = \frac{1}{1 + \frac{1}{e^\alpha}} = \frac{1}{1 + \frac{1}{\alpha}} = \frac{1}{1 + 0} = 1$$

$$h_\theta(x) = \theta_0 + \theta_1 x \rightarrow \text{Best fit line}$$

$$h_\theta(x) = \sigma(\theta_0 + \theta_1 x)$$

$$h_\theta(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$

$\rightarrow$  Logistic Regression model

Cost function

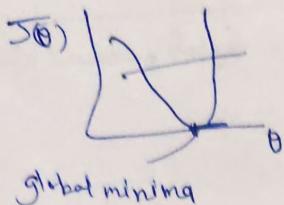
$$\frac{1}{n} \sum_{i=1}^n t(y_{act} - y_{pred})^2$$

To get optimal  $\theta_0$  &  $\theta_1$ , minimize the Cost function,  
linear Reg.

### Linear Regression

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i))^2$$

Convex function



### Logistic Regression

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i))^2$$

$$h_\theta(x_i) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_i)}}$$

but this is non convex  
so there is local minima



So here Cost function is use

### Log Loss function

$$J(\theta_0, \theta_1) = -y_i \log(h_\theta(x_i)) - (1-y_i) \log(1-h_\theta(x_i))$$

$$J(\theta_0, \theta_1) = \begin{cases} -\log h_\theta(x) & \text{if } y_i = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y_i = 0 \end{cases}$$

To minimize these  $\theta_0$  &  $\theta_1$ ,  
by convergence Algorithm

repeat until convergence

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j}$$

to get optimal  $\theta_0, \theta_1$

For Multiple Variable

$$h_\theta(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n)}}$$

### Logistic Regression with Regularisation

$$CF = J(\theta_0, \theta_1) = -y \log(h_\theta(x))$$

① Lasso:

② Ridge:

③ elastic net:

$$C = \frac{1}{\lambda}$$

Convex function  
when horizontal line  
cut 2 point



## Logistic Regression

$$\frac{1}{1+e^{-z}} \rightarrow \text{sigmoid}$$

No of hours studied	Passed
5	1
2	0
3	0
6	1
7	1
8	1
9	1
10	1

if N samples  $\begin{cases} 0 \\ 1 \end{cases}$

- For Samples labelled 1 : Estimate  $\theta_0$  such that  $\hat{p}(x)$  is as close as 1 possible.
- for Sample labelled 0 ; Estimate  $\theta_1$  such that  $1 - \hat{p}(x)$  is as close to 1 possible.  
 $\hat{p}(x)$  is close to 0 possible

$$\prod_{i=1}^n p(x_i) \cdot \prod_{i=1}^n 1 - p(x_i)$$

it means multiply all the labels

$$L(\theta) = \prod_{i=1}^n (p(x_i)^{y_i} \times (1-p(x_i))^{1-y_i})$$

Likelihood

do log both side

$$\log(L(\theta)) = \sum_{i=1}^n y_i \log p(x_i) + (1-y_i) \log(1-p(x_i))$$

log likelihood estimate

$\Rightarrow x \log p$

## Maximum Likelihood estimation

So we -ive So the maximum error reverse to minimum

→

$$\underset{\text{Maxima}}{\bigcup} (\text{Lve}) \prod$$

# 19 May - Classification Evaluation Metrics & Implementation

[is dataset]

res. orderly

→ newton-cg, sag, saga → these are optimization algorithms

→ Lasso, Ridge → alpha = h

$$C = \frac{1}{R} \quad (\text{because } \sum_i D_i \text{ is summable})$$

C smaller values specify stronger regularization

$$\frac{P}{1-P} \text{ odds} \Rightarrow \frac{P(\text{succ})}{P(\text{failure})}$$

$$\log \frac{P}{1-P}$$

$P \uparrow$  log odds  $\uparrow$

$T \uparrow$  odd  $\uparrow$  log  $\uparrow$

monotonic relationship  
→ one direction

→ misclassification

→ 2nd significant  
→ 2nd row  
→ 1st column

0.0	0.0
0.0	0.0

$$\text{FOP} = \frac{0+0.0}{0.0+1} = 0.0$$

→ same sign + needs break  $\rightarrow$   
add punctuation of phrasal?

F-biased towards two  
→ confusion matrix

## Evaluation metrics for classification

① Confusion matrix

② Accuracy / misclassification rate

③ Precision

④ Recall

⑤ F-Beta Score

⑥ True Positive rate (sensitivity)

⑦ False Positive rate

⑧ True Negative rate (specificity)

⑨ ROC-AUC

⑩ Precision-Recall / sensitivity-specificity trade off.

1:37

		Actual Value	
		True(1)	False(0)
Predicted Value	True(1)	TP	FP
	False(0)	FN	TN

	0	1	2
0	True <sub>0</sub>	false <sub>0</sub>	false <sub>0</sub>
1	false <sub>1</sub>	True <sub>1</sub>	false <sub>1</sub>
2	false <sub>2</sub>	false <sub>2</sub>	True <sub>2</sub>

$$\textcircled{D} \quad \underline{\text{Accuracy}} = \frac{TP + TN}{TP + FP + FN + TN}$$

\*Misclassification rate = opposite of Accuracy

$$\frac{FP + FN}{\text{total}} = 1 - \text{Accuracy} = 1 -$$

### ③ Precision :-

Impulse Class  
 $\downarrow$   
 1000dp → 900 c1  
 $\searrow$   
 100 class } ) →

Out of all predicted ones, how many are actual ones

$$\text{what accuracy without model} \\ \frac{4}{4} = \frac{900+0}{1000} = 90\%$$

Accuracy doesn't give more priority to minority class

Out of all actual value, how many are correctly predicted.

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

## Precision of the codes:

④ Recall: Out of all predicted values, how many are correctly predicted with actual value.

TP  
TP + FM

## USECASE-1      Spam Classifier

So, we wanted a model with high precision  
 $\rightarrow$  Threshold / detection

Increase the cutoff → Precise User

Increase the cutoff -  $0.5 \rightarrow 0.9$  Precision  
Use all

$TP \Rightarrow$  Mail - Spam(1) }  $\rightarrow$  ~~Correct~~  $\rightarrow$  ~~Not~~ (T)  
 Model - Spam(2) }  $\rightarrow$  ~~Correct~~  $\rightarrow$  ~~Not~~ (F)

$TN \Rightarrow$  Mail - Not Spam(0) }  $\rightarrow$  ~~Correct~~  $\rightarrow$  ~~Not~~ (A-009) (P)  
 Model - Not Spam(0) }  $\rightarrow$  ~~Correct~~  $\rightarrow$  ~~Not~~ (P)

~~Import~~  
 $FP \Rightarrow$  Mail  $\rightarrow$  Not spam Blunder }  $\rightarrow$  Wrong  
 Model  $\rightarrow$  spam } Predict

$FN \Rightarrow$  mail - spam  $\rightarrow$  Model  $\rightarrow$  Not spam } Predict



19 May  
2:40

If change the criteria

50% → Spam

After Accuracy the 1.0.9

Minimise the case FP

Goes to 30% → Spam

## USE CASE 2 :

$$\begin{matrix} \text{Pred :} & 1 & 0 \\ & TP & FP \\ 0 & FN & TN \end{matrix}$$

$$TP = \text{Actual - diabetic} - \text{Not diabetic}$$

FN Model Not "

TN

FP - Act → No diabetic  
Model - Predicted diabetic

1st level check  
2nd level check  
for cross check

VFN → Act - diabetic  
model - No diabetic

Blunder ✓ So it  
is important

⇒ FN is more important  
Why??

Use Recall

How Conviction of a person in the Court trial.

Stock Mkt Crash or not

(5) F-beta Score  $\frac{(1+\beta^2) P \times R}{P+R}$    
 P → Precision  
 R → Recall

① if FP and FN both are important  $\beta=1$

$$F_1 \text{ Score} = 2 \frac{P \times R}{P+R}$$

② if FP is more important than FN

$$\beta=0.5 \quad f_{0.5} \text{ Score} = \frac{(1+0.25) P \times R}{P+R}$$

③ if FN is more important than FP

$$\beta=2 \quad f_2 \text{ Score} = \frac{(1+4) P \times R}{P+R}$$

⑥ True Positive Rate  $\rightarrow$  Out of all actual 1, it is actually predicted one.

↳ Sensitivity, Recall

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

Some part  
Total All Part

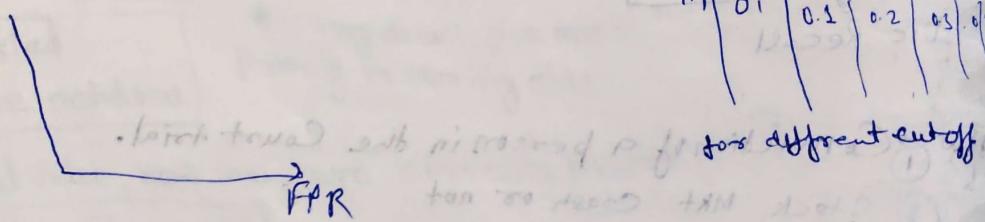
⑦ False Positive Rate:  $\frac{FP}{FP + TN}$

⑧ True Negative Rate  $TNR = \frac{TN}{FP + TN}$   
 ↳ Specificity  $\Rightarrow TNR = 1 - FPR$

⑨ ROC-AUC  $\rightarrow$  Receiver Operating Characteristic  
 Area under Curve.

TPR vs FPR

$\frac{TPR}{FPR}$

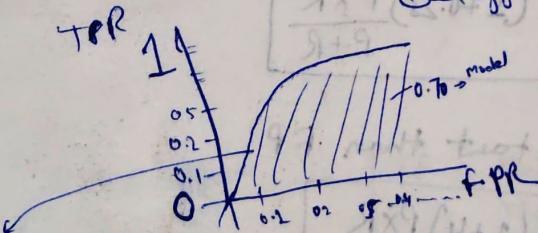


X	Y	Prob	cutoff
5	1	0.9	0.5 (yrc)
6	0	0.8	
7	1	0.6	
8	0	1	
9	1	1	
10	0	1	

for different cutoff calculate

0.1	0.2	0.3	0.4	.....	1
0	0	1	1		
1	1	0	0		
2	0	1	0		
3	1	0	1		
4	0	0	0		
5	1	1	1		
6	0	0	0		
7	1	0	0		
8	0	1	0		
9	1	0	0		
10	0	0	1		

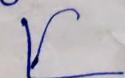
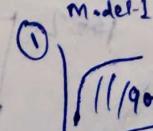
calculate all TPR & FPR at different cutoff & plot a chart



Not plot cutoff only  
 Plot those value of  
 TPR & FPR

0.70  $\rightarrow$  Model will correctly predict 70% of time.

• Higher the AUC better the model will be

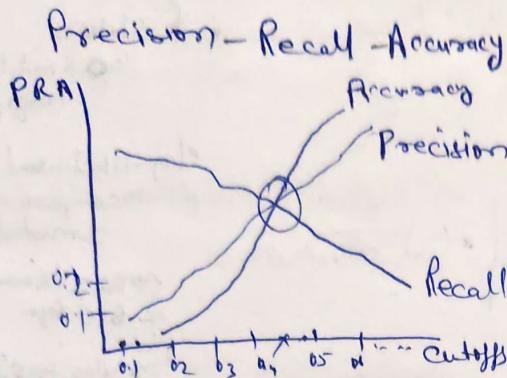


Cutoff  $\rightarrow$  by business

(2) Cutoff  $\rightarrow$  business team (b.)

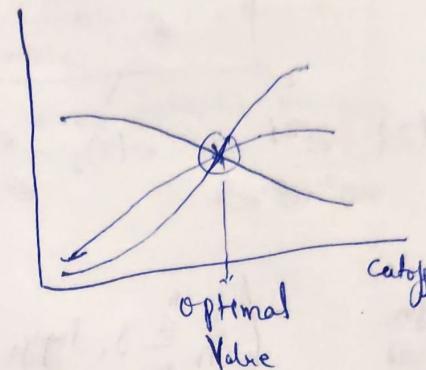
if No business team (default is 0.5)

$\Rightarrow$



OR

Sensitivity - Specificity - Accuracy



Where 'PRA' cut that is more accurately cut off and give the priority to all (PRA)

are  
both same  
concept

## class 25 May

### Multiclass Classification & Hyperparameters tuning

↳ OVR (One vs Rest)

Solution  $\rightarrow$  Multinomial

make data  $\rightarrow$  make\_classification (sample=1000, n\_features=10, n\_informative=5, n\_redundant=5, n\_classes=2)

- Linear Rf
- Logistic
- Hypothesis
- Multi-class
- Classification
- Decision tree
- Random Forest
- SGD
- ensemble
- PCA
- Annealing
- Unsupervised
- Time Series

ROC-AUC we for compare 2 models Performance

↳ One vs rest



Model 1  
10 folds  
20 fold

Model 2  
20 fold

20

R square

so we

Adjusted

R-Square

↓

Why because

Curse of dimension

↓ Time Compl.

↓ Overfitting

↓ No better

Interpretation

↳ Multinomial method / Softmax regression

Multinomial loss

specificity of Logistic

→ we don't decompose the problem into

binary classification.

different

$$\begin{cases} m_1 = 98\% \text{ (Acc, P, R)} \\ m_2 = 78\% \text{ (Acc, P, R)} \end{cases}$$

→ Modify the loss / cost function

→ Single Model

$$\rightarrow \text{Sigmoid } \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-z}} \Rightarrow \frac{1}{e^{z+1}} \frac{e^z}{1+e^z} !$$

Absolute	Relative
98% Sc	$\frac{TPR}{FPR}$
98% And	$\frac{FPR}{TPR}$
X no comp	↓
	Compare 2 model

- CHI-tot for (categorical) int Stats Mode)
  - OLS
  - f-test (Continued)
  - ↳ Comparing Variance

Softmax

$$\sigma(z)_j = \frac{e^{z_j}}{\sum_{j=1}^K e^{z_j}}, \quad j \text{ is no of class}$$

$$\sigma(z)_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}} \quad \sigma(z)_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}} \quad \sigma(z)_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

Cost fn (log loss) =  $-\left( \frac{1}{n} \sum_{i=1}^n y_i \log y_i + (1-y_i) \log(1-y_i) \right)$

Multiclass (multipl.)  $\Rightarrow -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K y_k^{(i)} \log \hat{y}_k^{(i)}$

$y_k^{(i)}$  → no of class  
 $i$  = no of datapoints

$x_1$	$x_2$	$y$	$y_{k=1}$	$y_{k=2}$	$y_{k=3}$
$x_{11}$	$x_{12}$	1	1	0	0
$x_{21}$	$x_{22}$	2	0	1	0
$x_{31}$	$x_{32}$	3	0	0	1

Class-1  $\hat{y}_1^{(1)} \log \hat{y}_1^{(1)} + \hat{y}_2^{(1)} \log \hat{y}_2^{(1)} + \hat{y}_3^{(1)} \log \hat{y}_3^{(1)} + \dots$

$P > 1+1$   
 when  $d \geq 30$

log-likelihood  
 compare 2 model

AIC (Akaike info)  
 BIC (Bayesian)  
 Similar to  $R^2$   
 but  $R^2$  better

- f-statistics
  - ↳ model is significant