

## Contingency Tables and Introduction to Stats

---


24/09/2022

---

---

---

---



avg. age =  $\frac{a_1 + a_2 + a_3 \dots + a_n}{n}$   
 Index

Population: all possible subjects

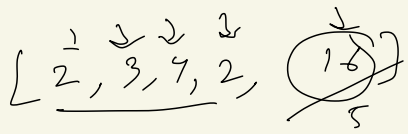
Sample: collected subset for our experiments.

- 1 - 12L
- 2 - 13L
- 3 - 12L
- 4 - 15L
- ...
- 20 - 17L

Summary statistics

Summary stats

Mean (avg.) avg. CTC. was 10L



CTC =  $\frac{27}{5}$   
 = 5.4L

Median [2, 2, 3, 4, 16]

Median CTC -> 3L

Mode: [20, 0, 0, 0, 30, 10, 5]

Mode = 0

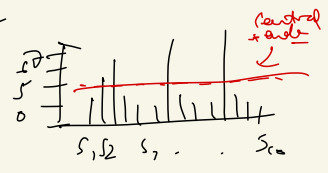
[2, 2, 3, 4, 16]

range of play that is offered in the class: (2 to 16L)

S.d. / variance: wide high range -> high variance

1, 3  
 -4, 4  
 60%

mean (median) ← (central tendency)  
 S.d. / variance ← variability



measure:

Population ( $\bar{x}$ )

Parameter ( $\mu$ )

(any prop. of underlying population) ( $\sigma$ )

Sample ( $\hat{x}$ )

Statistic ( $\hat{x}$ )

(estimating the property from sample  $\rightarrow$  statistic)

avg. wt of all Indians  $\rightarrow$  parameter (humanly Impractical)

sample of 1M Indians  $\rightarrow$  statistic (Practical)

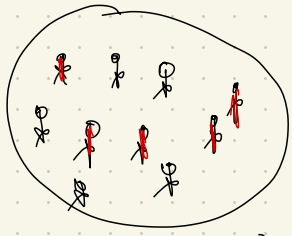
method to find out statistic from sample is called **Estimation**

answer will be estimate

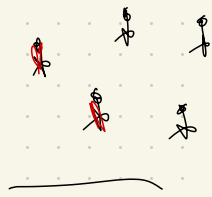
Sample of 1M Indian to find avg wt

$\rightarrow$  statistic

Val. of avg  $\rightarrow$  estimate



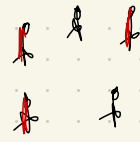
Population



prop. of ~~all~~

$$PPL \text{ w/ diabetes} = \left(\frac{2}{5}\right)$$

$$(\hat{p}) = .40$$

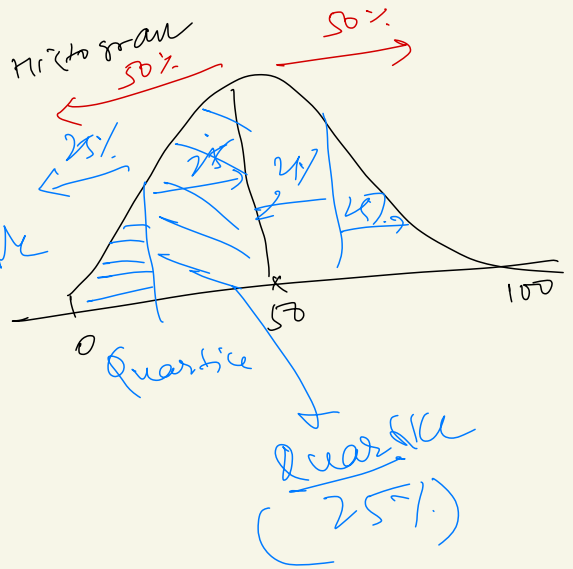
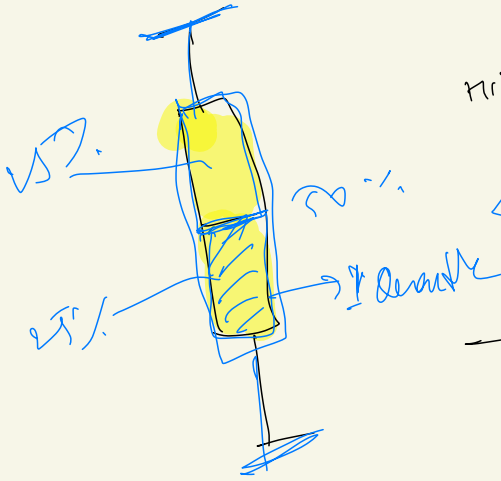


$$\hat{p} = \left(\frac{3}{5}\right) = 60\%$$

$\downarrow$  diabetic

	M	F
Person	.	.
Person	.	.

prop. of ppl who are diabetic are (35% - 40%) (95% Confid.)



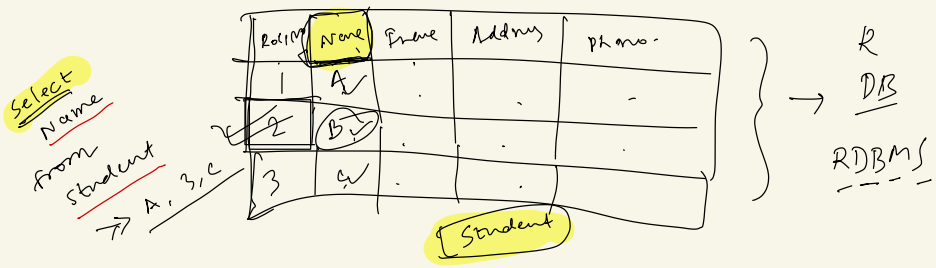
# Structured Query Language

Relational DB

Filtering rows & columns

Querying using SQL

	Attendance		Admission		Programs	
1						
2						
3						





11 Normalization

Primary key

<u>name</u>	class (year)	<u>Dept</u>	HOD	% marks
ABC	B.A.S	DS	/	70%
CDE	B.A.S	DS	/	95%

HOD

DS	RU
MG	BR



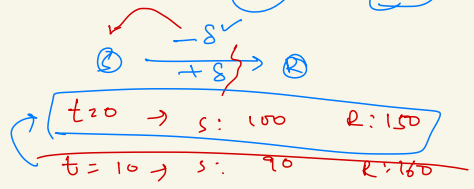
$\leftarrow \frac{amt + id,}{a/c} \leftarrow \frac{ac/bat}{\$}$

↓ transaction in its entirety  
 A C I D  
 ↓ ↓ ↓ ↓

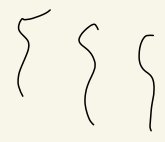
CBS

a/c no	balance
1	100 200
2	200
3	400

+30 }  
 -30 }  
 (800) → (800)



✓ A    ✓ B    ✓ C  
 +10   +10   +10  
 ✓ -10

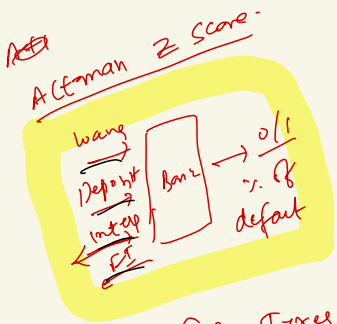


# Elementary Matrix Operations

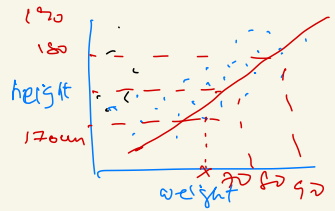
Module 2

(numpy) → python library for numerical processing

Goal: build some "system" which will predict height of a person given its weight.



ht (cm)	wt (kg)
170	70
150	65
180	80
⋮	⋮
200	?

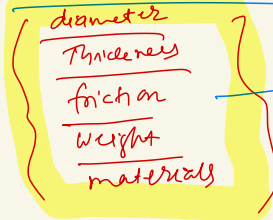


height & weight

Car Types



Durability of Car Types



probability of how car type spring failures

Date

{ current  
eng.  
mist }

Model

Results  
(Predictions)

	D	T	F	W	M	F
①	.	.	.	.	.	⊗
②	.	.	.	.	.	.
③	.	.	.	.	.	.
④	.	.	.	.	.	.



Data will have features -

muscular body, high wt, body hair, force (deep), etc.

5 cols /  
3 rows

Mus.	Wt	B.H	V	ply
1	58	1	1	M
0	-	-	-	F
0	-	-	-	F

1 - True (Present)  
0 - False (absent)

String of data: (Matrix / Matrices)  
singular / plural

1	2	3
4	5	6
7	8	9

(3x3)

2	3	6
2	1	6

C  
O  
L  
U  
M  
N

row

"Matrix"  
(3 cols, 3 rows)

rows: 2 (2x3)  
cols: 3

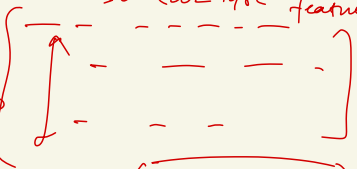
[ 1 ] → 1 row, 1 col.

[ 1 2 3 ] → (1 row, 3 col) (1x3)

Size of matrix:  
(row x cols)

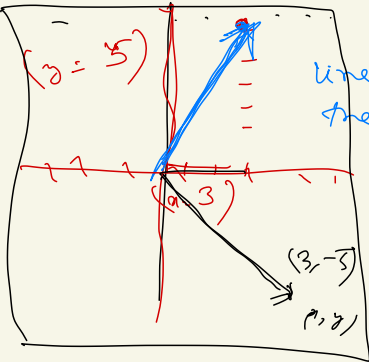
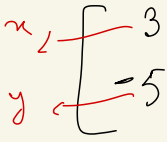
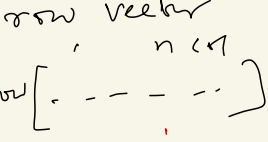
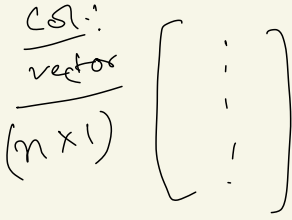
Vector: (direction + magnitude)  
speed, velocity, acceleration

30,000 cars



(30,000 x 50) cap. (12) LR.

n x 1 → col vector  
1 x n → row vector



2x + 3y = 1  
line connecting the pt (3,5) with origin (0,0) is vector.

Matrix (operation)

Addition

$$\underline{A} + \underline{B} \quad \begin{bmatrix} 2 & 2 \\ 1 & 0 \end{bmatrix} + \begin{bmatrix} 3 & 0 \\ 1 & 2 \end{bmatrix}$$

(2x2)      (2x2)

$$\begin{bmatrix} 3 & 2 & 1 \\ 0 & 1 & 2 \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0 \\ 2 & 2 \end{bmatrix}$$

(2x3)      (2x2)

Subtraction

$$\underline{A} - \underline{B}$$

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

Multiplication:

$$\begin{bmatrix} 2 & 3 \\ 4 & 5 \end{bmatrix} \times \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 2 \cdot 1 + 3 \cdot 0 & 2 \cdot 2 + 3 \cdot 1 \\ 4 \cdot 1 + 5 \cdot 0 & 4 \cdot 2 + 5 \cdot 1 \end{bmatrix} = \begin{bmatrix} 2 & 7 \\ 4 & 13 \end{bmatrix}$$

Matrix multiplication rule:

$$\begin{matrix} \underline{A} & & \underline{B} \\ \underline{m \times n} & & \underline{n \times p} \end{matrix}$$

no. of cols of 1<sup>st</sup> matrix =  
no. of rows of 2<sup>nd</sup> matrix

$$\begin{bmatrix} 2 & 3 & 2 \\ 4 & 1 & 6 \end{bmatrix} \begin{bmatrix} 2 & 3 & 2 \\ 4 & 1 & 6 \end{bmatrix}$$

2x3      ✗      2x3

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

$$\underline{3A} = \begin{bmatrix} 3 & 6 \\ 9 & 12 \end{bmatrix}$$

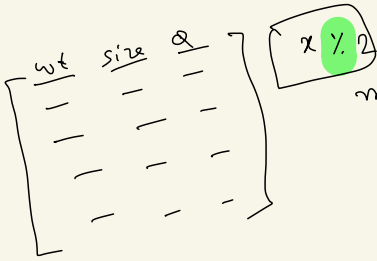
~~Transp~~: Transpose:

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

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

$$\begin{array}{l} \swarrow \\ \underline{(2 \times 3)} \\ A = \begin{bmatrix} 2 & 3 & 2 \\ 1 & 0 & 1 \end{bmatrix} \\ \downarrow \\ \underline{(3 \times 2)} \\ A^T = \begin{bmatrix} 2 & 1 \\ 3 & 0 \\ 2 & 1 \end{bmatrix} \end{array}$$

$$\begin{array}{l} A \\ \begin{bmatrix} 2 & 3 \\ 1 & 2 \end{bmatrix} \times \begin{array}{l} C \\ \begin{bmatrix} 3 \\ 0 \end{bmatrix} \end{array} \\ \hline = \begin{bmatrix} \phantom{0} \\ \phantom{0} \end{bmatrix} \end{array}$$



modulo operator

$$2) \begin{array}{r} 21 \\ 20 \\ \hline 1 \end{array} \quad (10)$$

remainder  
(answer)

$$21 \% 2 = 1$$

$$7 \% 3 = 1$$

$$19 \% 4 = \underline{\underline{3}}$$

python  
numpy  
array

1	2
3	4

matrices  
(2x2)

matrix

1	2	3	1
4	5	6	2
9	8	0	0

(3, 4)

(3x4)

no. of rows = 3

no. of cols = 4

→ [2]

1	2	3	4	5
6	7	8	9	1
2	1	1	1	1
1	1	1	1	1

(4x5)

dim = 2

(MxN)

[1 2 3]

row vector

1
2
3

col. vector

(n x 1)

$$A = \begin{bmatrix} 2 & 3 & 4 \\ 5 & 6 & 7 \\ 2 & 1 & 2 \end{bmatrix}$$

Transpose: interchanges row & cols

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

Transpose of A:

$$(A^T) = \begin{bmatrix} 2 & 5 & 2 \\ 3 & 6 & 1 \\ 4 & 7 & 2 \end{bmatrix}$$

Trace of Matrix:

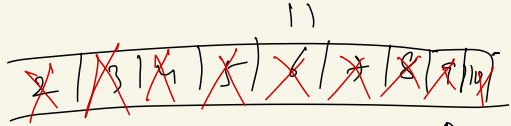
$$\text{tr}(A) = 2 + 6 + 2 = 10$$

$$\text{tr}(A) = \text{tr}(A^T)$$

Prime Number  
divisible by 1 or itself

21  
identify whether a number is prime or not?

- ~~21 ÷ 2~~ ✗
- 21 ÷ 3 ✓
- 19 ÷ 2
- 19 ÷ 3
- 19 ÷ 4
- ...
- 19 -



take the number - divisible by all number starting from 2 through (no. - 1)

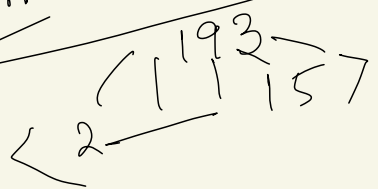
$$2 \dots 10$$

if not divisible Prime

(n)  
Testing primality  
Test from '2' to  $\sqrt{n} + 1$

$$113$$

$$\sqrt{113} = 10.6 \dots$$



2, 3, 4, 5, ..., 9, 10, 11



# Matrix Multiplication:

$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$\begin{matrix} (2 \times 3) \\ \hline A \end{matrix}$ 
 $\begin{matrix} (3 \times 3) \\ \hline B \end{matrix}$

$$= \begin{bmatrix} 1 \cdot 1 + 2 \cdot 0 + 3 \cdot 0 & 2 & 3 \\ 2 & 1 & 0 \end{bmatrix}$$

$2 \times 3$

$$A \Rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 0 \end{bmatrix}$$

Trace of matrix:

usually defined

$$\begin{bmatrix} 2 & 3 & 2 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

Transpose:

$$A^T = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 3 & 0 \end{bmatrix}$$

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

$3 \times 3$  (Identity matrix)

It has all 1's in its diagonal and 0's elsewhere (square matrix)

① matrix is square

$$\frac{A}{m \times m} \times \frac{B}{m \times m}$$

$(a \times b) \quad (b \times \_)$

$$\frac{a \times b}{a \times k} \quad \frac{b \times k}{a \times k}$$

$$\frac{2 \times 3}{2 \times 3} \times \frac{3 \times 3}{3 \times 3}$$

$\downarrow$

$$2 \times 3 \quad 3 \times 1$$

$(2 \times 1)$

$$= 2 + 0 + 1 = \boxed{3}$$

Matrix Vector multiplication

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

$3 \times 3 \quad \underline{\underline{3 \times 1}}$

$$= \begin{bmatrix} 2 + 0 + 1 \\ 1 + 0 + 0 \\ 0 + 0 + 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix}$$

$= 3 \times 1$

2x2 identity matrix :

5x5

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

$$\begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$$

$$\begin{bmatrix} -1 & 9 \\ 11 & 6 \end{bmatrix}$$

2x2

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

$I_2$

$$= \begin{bmatrix} -1 & 9 \\ 11 & 6 \end{bmatrix}$$

$$A \times I = A$$

$$A \times B \neq B \times A$$

(2x2) (2x2)

(commutative)

$$5 \times 3 = 15$$

$$3 \times 5 = 15$$

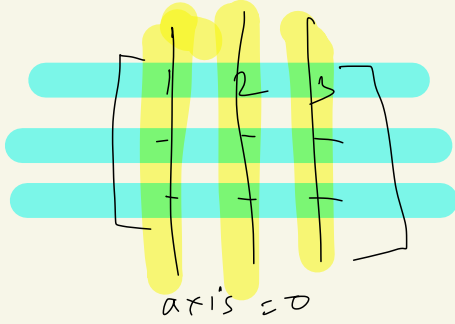
$$A = 3 \times 5$$

$$B = 5 \times 1$$

$$A \times B = 3 \times 1$$

$$B \times A = 5 \times 3$$

no. of axis =  
no. of dimension  
{0, 1}  
{0, 1, 2}



row axis = 1

col = 0



## Inner Product

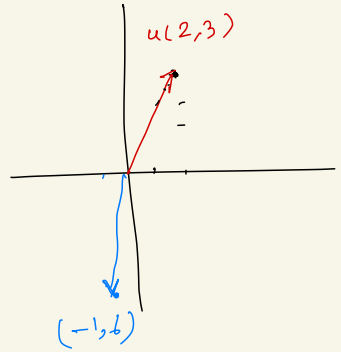
u:

$$u = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad v = \begin{bmatrix} -1 \\ 6 \end{bmatrix}$$

$$\begin{aligned} u^T v &= \begin{bmatrix} 2 & 3 \end{bmatrix} \begin{bmatrix} -1 \\ 6 \end{bmatrix} \\ &= 2(-1) + 3(6) \\ &= -2 + 18 \end{aligned}$$

$$u^T v = \boxed{16}$$

(-∞ ∞)



## Outer Product: u v^T

Inverse of matrix :

$$\frac{1}{\det(A)} \cdot \overline{(\text{adj}) A}$$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$\det(A) = \underline{\underline{ab - cd}}$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 6 \end{bmatrix}$$

$$\frac{1}{\det(A)} \begin{bmatrix} 6 & -3 \\ 2 & 1 \end{bmatrix}$$

$$\frac{\det}{\begin{bmatrix} 6 & 2 \\ 3 & 6 \end{bmatrix}}$$

$$\begin{bmatrix} (6)(6) - (3)(2) \\ (1)(1) \end{bmatrix} = \begin{bmatrix} 30 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 6 & 2 \\ 3 & 6 \end{bmatrix} = \begin{bmatrix} 6-6 \\ 0 \end{bmatrix}$$

$$u = [2 \ 3 \ 6]$$

$$\begin{aligned} \|u\| &= \sqrt{2^2 + 3^2 + 6^2} \\ \text{norm} &= \sqrt{4 + 9 + 36} \\ &= \sqrt{49} = 7 \end{aligned}$$

$$\|u\|_2 = \sqrt{x_1^2 + x_2^2 + x_3^2}$$

$$\|u\|_3 = \sqrt[3]{x_1^3 + x_2^3 + x_3^3}$$

$$\|u\|_1 = x_1 + x_2 + x_3$$

$$\begin{array}{c} \downarrow 1 \\ \downarrow 2 \\ \downarrow 3 \end{array} \begin{bmatrix} 0 & 3 & 6 \\ 1 & 4 & 7 \\ 2 & 5 & 8 \end{bmatrix}$$

Column major form

Linearizing a matrix:

0, 1, 2, 3, 4, 5, 6, 7, 8

no. of rows = 3

element = row

$$\begin{aligned} \text{element} &= \text{row} + (\text{column} \times \text{no. of rows}) \\ &= 1 + (1 \times 3) \\ &= 4 \\ &= 2 + 1 \times 3 = 5 \end{aligned}$$

row #	column #	element
0	0	0
1	0	1
2	0	2
0	1	3
1	1	4
2	1	5
	⋮	
		6

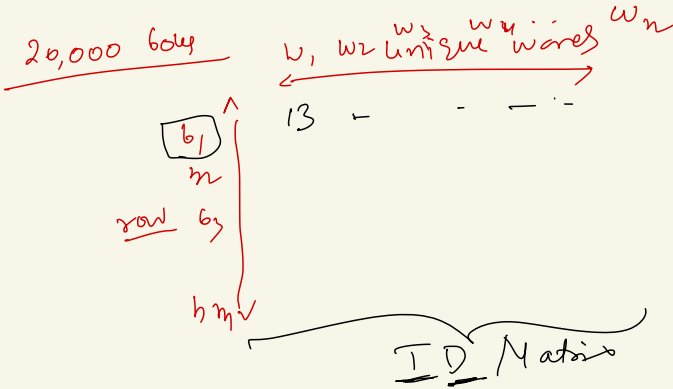
	$\tau_1$	$\tau_2$	$\tau_3$	$\tau_4$	...	$\tau_m$
$v_1$	✓	x				
$v_2$			x			
$v_3$				✓		
$\vdots$						
$v_n$	x					✓

$10 \rightarrow 100 \text{ sec}$   
 $11 \rightarrow 121 \text{ sec}$   
 $12 \rightarrow 144 \text{ sec}$

No. of entries in  
 matrix  $< 5\%$   

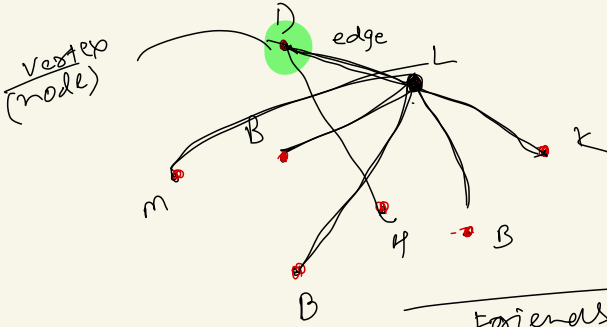

---

 $95\% \rightarrow \text{zero}$   
 " sparse "



# Graphs

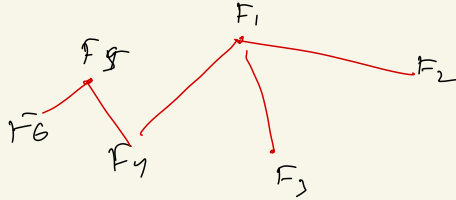
collection of nodes and edges  
(vertex)



connection  $\rightarrow$  edges  
nodes/vertices

## Friend's network

$u_1, p_1$   
 $u_2, p_2$   
 $u_3, p_3$



df:

col A	col B
1	A
2	B
3	C
D	d
d	D
A	1
B	2
C	3
<del>d</del>	<del>D</del>
<del>D</del>	<del>d</del>

df

A	B
B	A

df dupl

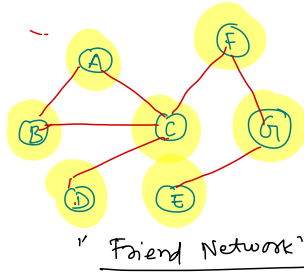
Note Book 10 (NumPg) : part 3 (Sparse matrices).

DS/ML → Data matrix

We want to



Friendship network in computer → Graphs are used.

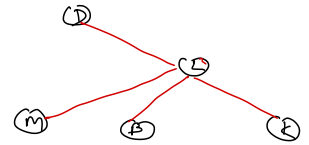


	M	B	J	E
M			1	
B	1			
J				1
E			1	

Collection of **nodes** and **edges**

Storing Graphs in Computers:

(i) Adjacency matrix



	A	B	C	D	E	F	G
A		1	1				
B	1		1				
C	1	1			1	1	
D				1			
E			1				
F							1
G						1	1

COO <sup>Point</sup> Format:

Row:	A	A	B	B		C	G
col:	B	C	A	C	-	-	F
val:	1	1	1	1		1	1

src      Target

1      a

2      b

4      d

d      4

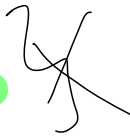
a      1

b      2

d      4

4      d

a



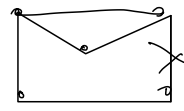
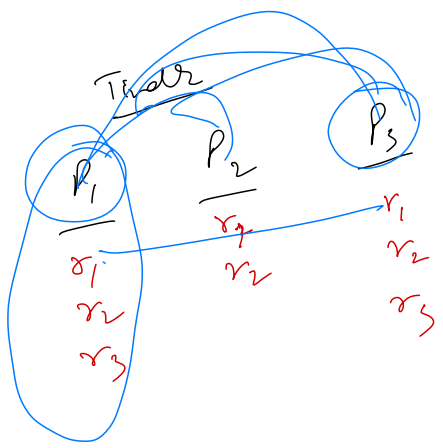
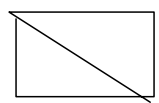
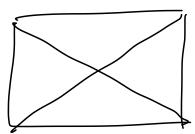
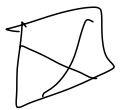
$$\begin{matrix}
 & \textcircled{0} & \textcircled{1} & \textcircled{2} \\
 0 & | & 0 & -2.5 & | & 1.2 \\
 1 & | & 0.1 & 1.0 & | & 0 \\
 2 & | & 6.0 & -1.0 & | & 0
 \end{matrix}
 \left[ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} \right] = \left[ \begin{matrix} 0+(-5) \cdot 1.2 \\ \dots \\ \dots \end{matrix} \right]$$

$3 \times 3 \qquad \qquad \qquad 3 \times 1$

Matrix  $\times$  Vector = vector  
 $(m \times n) (n \times 1) = (m \times 1)$

default dictionary

# (Don)	keys	val
	{ 0	→ 1: -2.5, 2: 1.2 }
	{ 1	→ 0: 0.1, 1: 1.0 }
	{ 2	→ 0: 6.0, 1: -1.0 }

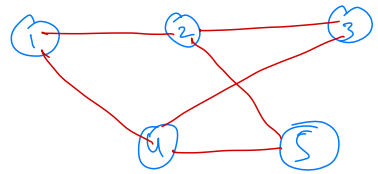


JEE

Roll no	Rank	
1.	5	<u>ACS</u>
2	1	<u>B-~</u>
3	2	

$\left\{ \frac{100!}{5!} \right\}$

100!

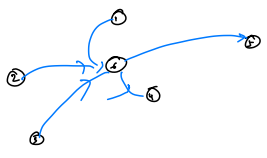


Page Rank:

Ranking algorithm of web pages.  
(Slow) → Execution of lookup/Query Retrieval (Fast)

Story:

- 1 - 30%
- 2 - ~~20%~~
- 3 - 40%
- 4 - 20%
- 5 - 5%
- 6 - 5%



website  
- Text/Image/MM  
- Links

$$\underline{x} \leftarrow \underline{x}A$$

$1 \times n$        $n \times n$

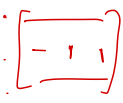
matrix  $\times$  vector  
= vector

$$2A = x$$

$A \times$  eigen vector =  $\lambda \cdot v$   
 eigen val of  $\lambda$  = eigen value



Row	2	2
col	2	3
val	1	1

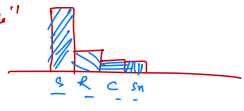


Markov chains

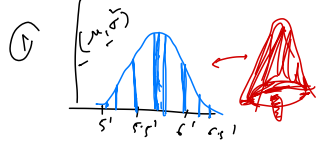
60%  $\frac{7 \times 7}{8}$  80%  
 $\frac{6}{7}$        $\frac{7 \times 7}{8}$        $\frac{80}{9}$  "85%"  
 (8) 90% "1st step base"  
 (7, 18) 2 step

memoryless

60% Sunny    20% Rainy    10% Cold    10% Snow



height of student in class:

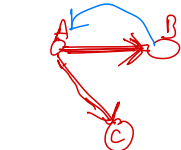


Symmetric;  
Normal distribution  
(Bell curve)  
Gaussian Dist.

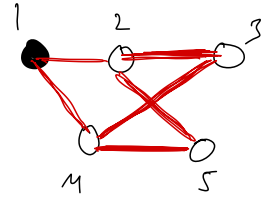
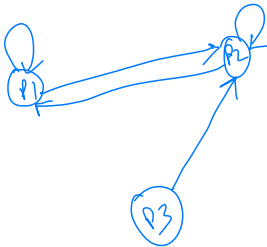
② Poisson Distribution:  
# white color vehicles  
{ 1000 → # " }



	A	B	C	D	E	F
A						
B						
C						
D						
E						
F						



	A	B	C
A			
B			
C			



Markovian:

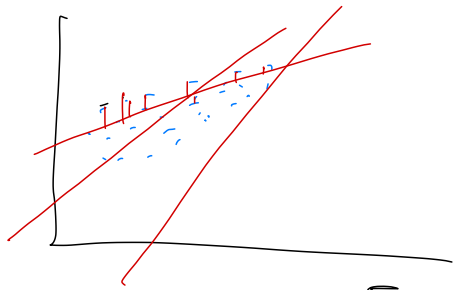
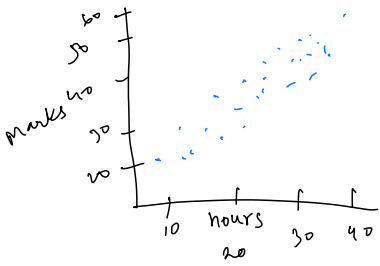
Future depends on past

1 steps

2 steps.

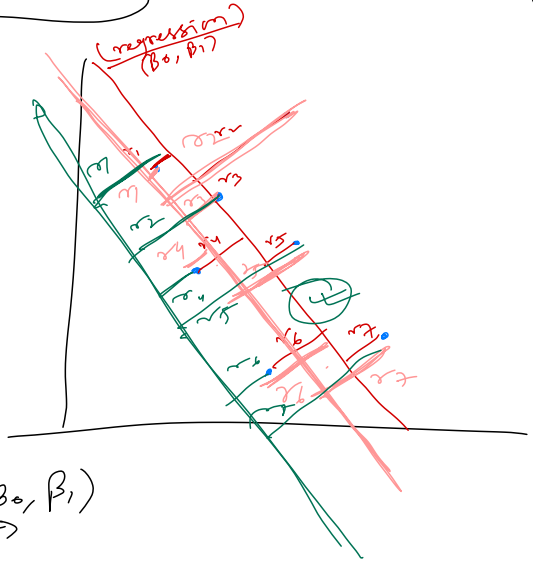






$$\sum_{i=1}^n x = \underbrace{x + x + x + \dots + x}_n = \sum x = \boxed{nx}$$

minimize  
the residuals  
 $r_1 + r_2 + \dots + r_n$   
should be as  
small as  
possible



Data  
(hours) (x)  
(marks) (y)



$(B_0, B_1)$

$$\min: 2x^2 - 3x + 5 \rightarrow 2 \cdot \frac{9}{16} - 3 \cdot \frac{3}{4} + 5 \quad (5-2)^2 = 9$$

$$\frac{d}{dx} (2x^2 - 3x + 5) = 0 \rightarrow \frac{4x}{8} - \frac{3}{4} + 5$$

$$\frac{4x - 3}{8} = 0 \rightarrow \frac{4x - 3}{8} + 5$$

$$x = \frac{3}{4}$$

$$= 5 - \frac{9}{8}$$

$$= \underline{\underline{3}}$$

$$\hat{y}_i = \beta_0 + \beta_1 x_i \quad (\text{estimated from data})$$

$(y_i - \hat{y}_i)$  : residual.

SSR (sum of square residuals)  
(to cancel the effect of opp. signs)

$$= \sum_i (y_i - \hat{y}_i)^2 \quad [(a-b)^2 = a^2 - 2ab + b^2]$$
$$= \sum_i (y_i - (\beta_0 + \beta_1 x_i))^2$$

$$= \sum_i \left[ y_i^2 - 2y_i(\beta_0 + \beta_1 x_i) + (\beta_0 + \beta_1 x_i)^2 \right]$$

$$= \sum_i \left[ y_i^2 - 2y_i\beta_0 - 2y_i\beta_1 x_i + \beta_0^2 + \beta_1^2 x_i^2 + 2\beta_0\beta_1 x_i \right]$$

$$\left[ \frac{\partial \text{SSR}}{\partial \beta_0} = \sum_i \left( 0 - 2y_i + 0 + 2\beta_0 + 2\beta_1 x_i \right) = 0 \right]$$

$$= 2 \left( \sum \beta_0 + \beta_1 x_i - y_i \right) = 0$$

$$\Rightarrow \sum (\beta_0 + \beta_1 x_i - y_i) = 0$$

$$\Rightarrow \sum_{i=1}^n \beta_0 = \sum_{i=1}^n (y_i - \beta_1 x_i)$$

$$\Rightarrow n\beta_0 = (y_1 + y_2 + y_3 + \dots + y_n) - \beta_1 (x_1 + x_2 + \dots + x_n)$$

$$\beta_0 = \frac{(y_1 + y_2 + \dots + y_n)}{n} - \beta_1 \cdot \frac{(x_1 + x_2 + \dots + x_n)}{n}$$

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

$$SSR = \sum_i \left[ y_i^2 - 2y_i \beta_0 - 2y_i \beta_1 x_i + \beta_0^2 + \beta_1^2 x_i^2 + 2\beta_0 \beta_1 x_i \right] \quad (1)$$

$$\frac{\partial SSR}{\partial \beta_1} = \sum_i \left[ 2\beta_1 x_i^2 + 2\beta_0 x_i - 2x_i y_i \right] = 0$$

$$= \sum_i \left[ -2x_i \left[ y_i - (\beta_0 + \beta_1 x_i) \right] \right] = 0 \quad (2)$$

substituting  $\beta_0$  in (2)

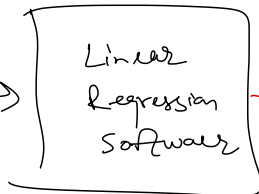
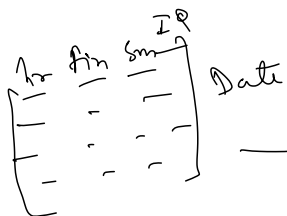
$$= \sum x_i \left[ y_i - (\bar{y} - \beta_1 \bar{x}) + \beta_1 x_i \right]$$

$$= \sum x_i \left[ y_i - \bar{y} - \beta_1 (x_i - \bar{x}) \right] = 0$$

$$= \sum x_i (y_i - \bar{y}) - \sum \beta_1 x_i (x_i - \bar{x}) = 0$$

$$\Rightarrow \sum x_i (y_i - \bar{y}) = \beta_1 \sum x_i (x_i - \bar{x})$$

$$\beta_1 = \frac{\sum x_i (y_i - \bar{y})}{\sum x_i (x_i - \bar{x})} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$



Estimate for (5)

Estimate for Betas.  
(slope & intercept)

# Multiple Linear Regression:

Predictor  $\rightarrow$  variables  
 hours studied,  
 IQ, financial status,  
 Smoking - since,

Marks scored =  $\beta_0 + \beta_1 \times \text{hours studied} +$

$\beta_2 \times \text{IQ} + \beta_3 \text{ fin} + \beta_4 \text{ Smoking}$

$\beta_5 \text{ Since (Yrs)}$

hrs studied	IQ	financ e	Smoking	Marks scored
-	-	-	-	-
-	-	-	-	-

response/prediction

MLR: has more than 1 predictor

$Y = \beta X + \epsilon$

Annotations:  
 -  $Y$ : response/prediction  
 -  $\beta$ : coefficient  
 -  $X$ : predictor (independent vars)  
 -  $\epsilon$ : "noise"

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

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

$y_0 = \beta_0 + \beta_1 x_0$   
 $y_1 = \beta_0 + \beta_1 x_1$   
 $y_2 = \beta_0 + \beta_1 x_2$

$$\begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \\ y_4 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_0 \\ 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_0 \\ \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\begin{bmatrix} 1 & x_0 & y_0 \\ \vdots & \vdots & \vdots \\ 1 & x_n & y_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

$(n \times 1) = (n) \times 2 \times (2 \times 1) + (n \times 1)$

This is not equivalent to  $y = mx + c$  if  $c \neq \epsilon$

$Y = X\beta + \epsilon$

Annotations:  
 -  $X$ : design matrix  
 -  $\beta$ : weights matrix  
 -  $\epsilon$ : noise/disturbance matrix

This is matrix noise/disturbance matrix.

$y = \beta_0 + \beta_1 \times \text{hr} + \beta_2 \times \text{IQ} + \beta_3 \text{ fin} + \beta_4 \text{ Smoking}$

Annotations:  
 -  $\beta_0$ :  $i_1$   
 -  $\beta_1$ :  $i_2$   
 -  $\beta_3$ :  $i_3$   
 -  $\beta_4$ :  $i_4$

$\int x^2 + \int e^x$   
 $\frac{x^3}{3} + e^x + c$   
 $\frac{x^3}{3} + e^x + c$

$$y = \beta_0 + \beta_1 \times hr + \beta_2 \times IQ + \beta_3 \times fin + \beta_4 \times smoking$$

$y = \text{max}(C)$

$$y = (\beta_0) + \beta_1(hr) + (\beta_2 \times IQ) + (\beta_3 \times fin)$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & \text{hours} & IQ & fin & smoking \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

noise

does not mean intercept is "1".

$$\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_0 \\ 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_0 \\ \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$(n \times 2)$        $(2 \times 1)$

$$Y = X\beta + \epsilon$$

$\epsilon \neq \text{intercept}$

✓  $y_0 = \beta_0 + \beta_1 x_0 + \epsilon_0$  ← 1st row of dataset

✓  $y_1 = \beta_0 + \beta_1 x_1 + \epsilon_1$  ← 2nd row of dataset

✓  $\vdots$

$$\epsilon = (\hat{Y} - \hat{X}\hat{\beta})$$

$$d(\sin^2 x + \cos^2 x)$$

inner product:  $\Rightarrow (x^T x)$

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

to minimize  $\epsilon$ , let's find out inner product

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} (2 \times 1)$$

$$x = \begin{bmatrix} 1 & x_1 & x_2 \end{bmatrix} (n \times 2)$$

$$\epsilon^T \epsilon = (Y - X\beta)^T (Y - X\beta)$$

$$= (Y^T - (X\beta)^T) (Y - X\beta)$$

$$= (Y^T - \beta^T X^T) (Y - X\beta)$$

$$= Y^T Y - \underline{Y^T X \beta} - \underline{\beta^T X^T Y} + \beta^T X^T X \beta$$

$$(Y^T X \beta)^T = \beta^T X^T (Y^T)^T$$

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

$$\epsilon^T \epsilon = \underbrace{Y^T Y}_0 - 2 \beta^T \underbrace{X^T Y} + \beta^T \underbrace{X^T X} \beta$$

$$\frac{\partial \epsilon^T \epsilon}{\partial \beta} = -2 X^T Y + 2 X^T X \beta$$

$$0 = -2 X^T Y + 2 X^T X \beta$$

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

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

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

$$\begin{matrix} A & B \\ (\tilde{A} \tilde{B})^T \\ = B^T A^T \end{matrix}$$

$$\begin{matrix} (A - B)^T \\ = A^T - B^T \end{matrix}$$

$$\frac{d}{dx} (y^T x) = y^T$$

$$\frac{\partial a^T b}{\partial b} = \frac{\partial b^T a}{\partial b} = a$$

$$\frac{\partial b^T A b}{\partial \beta} = 2 A b = 2 b^T A$$

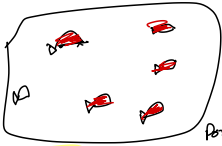
Ordinary  
Least Square

Normal form  
Equation

$$y = \beta x$$

$$\beta = \frac{y}{x} = \boxed{\beta_0 x^{-1} y}$$

(50 x 5)

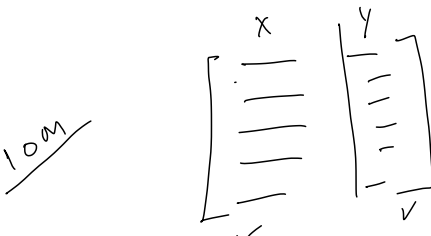


$$= \frac{y}{(X^T X)} X^T$$

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

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

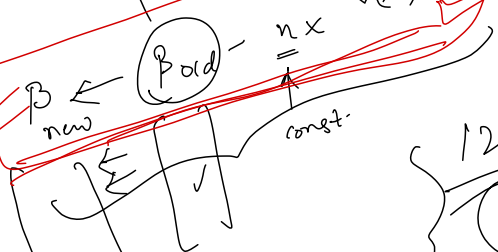
~~$A^T A$~~   
 $(A^T y)$   
 Inner product  
 "matrix"  
 50 examples  
 $y$   
 $x =$   
 $(n \times 2)^T$   
 $= (2 \times n) \times n$   
 $\rightarrow 2 \times 2$   
 $A^T A$   
 $= \boxed{B}$



Age, Smoking, FB, F.I.

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

100%  
 Iterative algorithm



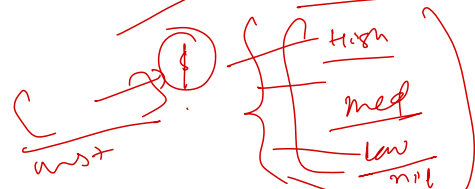
inner product  
 $A = m \times n$   
 $A^T = n \times m$   
 $A^T A$   
 $= n \times [m \times m] \times n$   
 $D = n \times n$

$$x^2 - (a+b)x + ab = 0$$

$$x^2 - b \pm \sqrt{b^2 - 4ac}$$

$$x^2 - 1 + 2x = 0$$

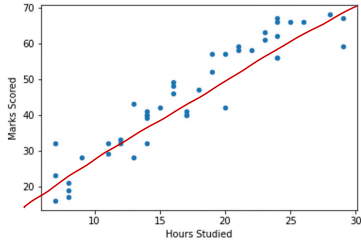
$$2x^2 - 2x + 3x^2 = 0$$



1.29  
 2.01

# Logistic Regression

Hours Studied	Marks Scored
12	32
16	46
16	49
13	43
29	59
8	17
24	56
14	40
7	23
15	42
::	::
19	57
24	62
14	40
24	67
19	52



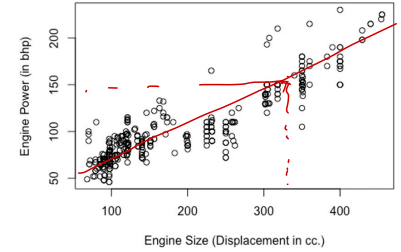
*marks = bot  $\beta$  x hrs.*

*"(Regression)"  
continuous  
(0,  $\infty$ )*

*"(Classification)"  
categorical*

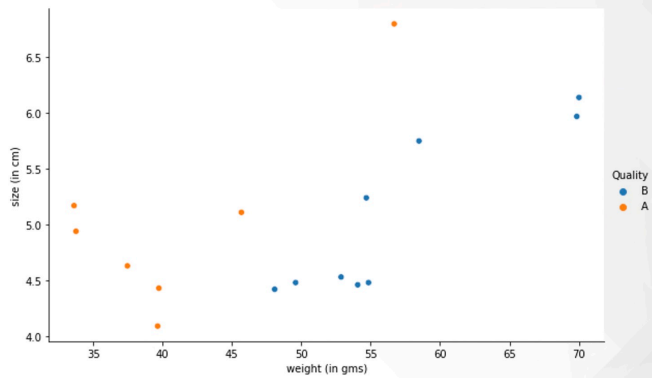
*$x_1, x_2, x_3, x_4, \dots, x_n$  (y)  
data*

Engine Size (cc)	Engine Power (bhp)
307	130
350	165
318	150
304	150
302	140
454	220
440	215
199	97
200	85
97	88



weight (in gms)	size (in cm)	Quality
54.84	4.48	B
48.07	4.42	B
54.68	5.24	B
54.06	4.46	B
70.02	6.14	B
45.67	5.11	A
69.86	5.97	B
37.45	4.63	A
33.60	5.17	A
49.58	4.48	B
39.63	4.09	A
39.72	4.43	A
58.48	5.75	B
33.74	4.94	A
52.86	4.53	B
56.69	6.80	A

Quality A is a better quality than Quality B



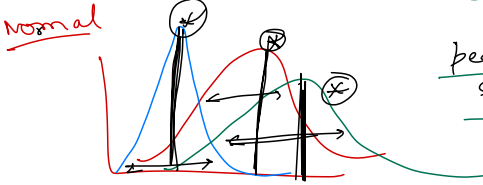
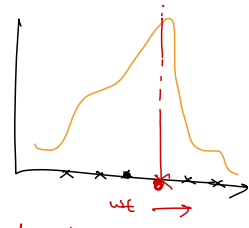
2 types of predictions — Continuous & Categorical



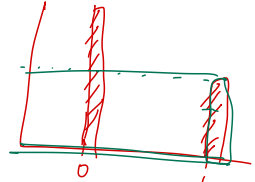
Given the dataset find most likely parameter of the distribution

Distribution: a graph showing prob.

$\mu, \sigma^2$  ✓ Normal distribution → distr. of height  
 Binomial " →  $(k, n, p)$  ✓  
 Poisson " →  $(\lambda)$  ✓



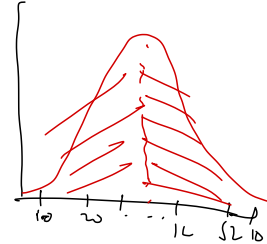
peak position → mean ( $\mu$ )  
 spread → variance ( $\sigma^2$ )



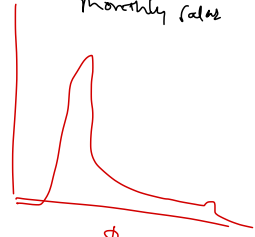
count of english speaking ppl



distr. of height of avg 30yr males in IN



monthly sales



$\$ \rightarrow$   
 "log normal"

$$\frac{d}{dx}(x^2 - 2x) = 2x - 2 \qquad \frac{d}{dx} \log x = \frac{1}{x}$$

$$\frac{d}{dx} \log(x - x^2) = \frac{1}{(x - x^2)} \times (1 - 2x)$$

Chain rule

$$\frac{d(uv)}{dx} = uv' + vu'$$

~~$$\frac{d(x \log p)}{dp} = x \log p$$~~

$$x \frac{d}{dp} \log(1-p) = \frac{1}{(1-p)} (0-1)$$

$$= \left( \frac{-1}{1-p} \right)$$

Coin

HTHTTH

define: prob. of head =  $p(H) = p$   
prob. of tail =  $p(T) = (1-p)$  ✓

MLE:  $L = p \times (1-p) \times p \times (1-p)^3 \times p$

$$L = p^3 (1-p)^4$$

$$\log L = \log(p^3 \cdot (1-p)^4)$$
$$= \log p^3 + \log (1-p)^4$$

$$\log L = 3 \log p + 4 \log (1-p)$$

$$\frac{\partial \log L}{\partial p} = \frac{3}{p} + \frac{4}{(1-p)} (-1) \quad (\text{from chain rule})$$

$$\Rightarrow \frac{3}{p} - \frac{4}{(1-p)} = 0$$

$$\Rightarrow \frac{3}{p} = \frac{4}{1-p}$$

$$3 - 3p = 4p$$

$$7p = 3$$

$$p = \frac{3}{7}$$

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

logarithm:

$$2^3 = 8$$

$$\log_2 8 = 3$$

$$\log_2 2^3 = 3$$

log number =  $a^b$   
(base)

$$\log_2 3^6 = 6$$

$$\frac{3^6}{3}$$

$$\log_2 3^3$$

$$\log_a a^b = b \log_a a$$

$$\log(asb) = \log a + \log b$$

$$\frac{d(\log p)}{dp} = \frac{1}{p}$$

$$\log \left( \frac{p}{1-p} \right) = \beta_0 + \beta_1 x$$

ln

$$\log_a b = c \Rightarrow b = a^c$$

log odds

$$\frac{p}{1-p} = e^{(\beta_0 + \beta_1 x)}$$

Betting

$(-\infty, \infty)$

$$\frac{1}{0} = \infty$$

# R/F vs R/V  
201 vs 180

$$\text{odds} = \frac{\pi_{\text{win}}}{\pi_{\text{loss}}} =$$

$$\text{odds of R/F winning} = \left( \frac{201}{180} \right) \frac{\text{win}}{\text{loss}} = \left( \frac{1}{0} \right)$$

$$p = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$= \frac{t}{1+t} = \frac{t/t}{(1+t)/t} = \frac{1}{\frac{1}{t} + 1}$$

$$= \frac{1}{1+t^{-1}}$$

$$= \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

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

(Sigmoid logistic)

**Example (Normal data).** Maximum likelihood estimation can be applied to a vector valued parameter. For a simple random sample of  $n$  normal random variables,

$$L(\mu, \sigma^2 | \mathbf{x}) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-(x_1 - \mu)^2}{2\sigma^2} \right) \cdots \left( \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-(x_n - \mu)^2}{2\sigma^2} \right) = \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2.$$

$$\ln L(\mu, \sigma^2 | \mathbf{x}) = -\frac{n}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2.$$

$$\frac{\partial}{\partial \mu} \ln L(\mu, \sigma^2 | \mathbf{x}) = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = \frac{1}{\sigma^2} n(\bar{x} - \mu)$$

Because the second partial derivative with respect to  $\mu$  is negative,

$$\hat{\mu}(\mathbf{x}) = \bar{x}$$

is the maximum likelihood estimator.

$$\frac{\partial}{\partial \sigma^2} \ln L(\mu, \sigma^2 | \mathbf{x}) = -\frac{n}{\sigma^2} + \frac{1}{(\sigma^2)^2} \sum_{i=1}^n (x_i - \mu)^2 = \frac{n}{(\sigma^2)^2} \left( \sigma^2 - \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \right).$$

Recalling that  $\hat{\mu}(\mathbf{x}) = \bar{x}$ , we obtain

$$\hat{\sigma}^2(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{x})^2.$$

**Example** Suppose that  $X$  is a discrete random variable with the following probability mass function: where  $0 \leq \theta \leq 1$  is a parameter. The following 10 independent observations

$X$	0	1	2	3
$P(X)$	$2\theta/3$	$\theta/3$	$2(1-\theta)/3$	$(1-\theta)/3$

were taken from such a distribution: (3,0,2,1,3,2,1,0,2,1). What is the maximum likelihood estimate of  $\theta$ .

**Solution:** Since the sample is (3,0,2,1,3,2,1,0,2,1), the likelihood is

$$L(\theta) = P(X=3)P(X=0)P(X=2)P(X=1)P(X=3) \times P(X=2)P(X=1)P(X=0)P(X=2)P(X=1)$$

Substituting from the probability distribution given above, we have

$$L(\theta) = \prod_{i=1}^n P(X_i | \theta) = \left(\frac{2\theta}{3}\right)^2 \left(\frac{\theta}{3}\right)^3 \left(\frac{2(1-\theta)}{3}\right)^3 \left(\frac{1-\theta}{3}\right)^2$$

Let us look at the log likelihood function

$$\begin{aligned} l(\theta) &= \log L(\theta) = \sum_{i=1}^n \log P(X_i | \theta) \\ &= 2 \left( \log \frac{2}{3} + \log \theta \right) + 3 \left( \log \frac{1}{3} + \log \theta \right) + 3 \left( \log \frac{2}{3} + \log(1-\theta) \right) + 2 \left( \log \frac{1}{3} + \log(1-\theta) \right) \\ &= C + 5 \log \theta + 5 \log(1-\theta) \end{aligned}$$

where  $C$  is a constant which does not depend on  $\theta$ . It can be seen that the log likelihood function is easier to maximize compared to the likelihood function.

Let the derivative of  $l(\theta)$  with respect to  $\theta$  be zero:

$$\frac{dl(\theta)}{d\theta} = \frac{5}{\theta} - \frac{5}{1-\theta} = 0$$

and the solution gives us the MLE, which is  $\hat{\theta} = 0.5$ .

randomly predicted

so

$(50 \times 2)$

$(2 \times 1)$

=

$(50 \times 1)$

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

$S = \frac{\sigma}{\sigma^2}$       $S = \sum$

$$\sigma(3) = \frac{1}{1+e^{-3}}$$

$$\sigma(x) = \frac{1}{1+e^{-x}} \quad (\text{Sigmoid function})$$

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

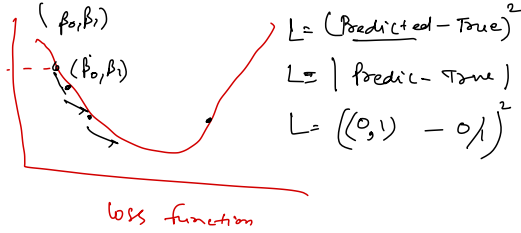
$$\frac{d}{dx} \sigma(x)$$

$$\frac{d}{dx} \left( \frac{1}{x} \right) = -\frac{1}{x^2} = -x^{-2}$$

chain rule of derivative

$$\frac{d}{dx} (u \cdot v) = u \cdot v' + u' \cdot v$$

loss:



$$\frac{d \sigma(x)}{dx} = \sigma(x) \cdot (1 - \sigma(x))$$

For each training data-point, we have a features  $x_i$  and an observed class,  $y_i$ .

⊛ The probability of the class is  $p$ , if  $y_i = 1$ , or  $1 - p$  if  $y_i = 0$

Cross Entropy  
loss L<sub>CE</sub>

$$\begin{aligned} P(Y|X) &= p(x_i)^{y_i} \cdot (1 - p(x_i))^{1-y_i} \\ P(Y|X) &= \hat{y}^{y_i} \cdot (1 - \hat{y})^{(1-y_i)} \end{aligned}$$

$$\begin{aligned} \text{head } \hat{y} & \\ y_i = 1 & \rightarrow p \\ y_i = 0 & \rightarrow 1 - p \end{aligned}$$

$$\hat{y} = \sigma(\beta_0 + \beta_1 \cdot x)$$

$\frac{1}{1+e^{-(\beta_0+\beta_1 \cdot x)}}$  estimate/prediction  $\rightarrow \hat{y}$   
true / actual  $\rightarrow y$

$$\log l = y \log \hat{y} + (1 - y) \log(1 - \hat{y}) \quad (\text{using log-identities})$$

$$= y \log \sigma(\beta_0 + \beta_1 \cdot x) + (1 - y) \log(1 - \sigma(\beta_0 + \beta_1 \cdot x))$$

$$\frac{\partial \log l}{\partial \beta_j} = \frac{y}{\sigma(\beta_0 + \beta_1 \cdot x)} \frac{\partial \sigma(\beta_0 + \beta_1 \cdot x)}{\partial \beta_j} + \frac{(1 - y)}{(1 - \sigma(\beta_0 + \beta_1 \cdot x))} \frac{\partial (1 - \sigma(\beta_0 + \beta_1 \cdot x))}{\partial \beta_j}$$

$$= \frac{y}{\sigma(\beta_0 + \beta_1 \cdot x)} \frac{\partial \sigma(\beta_0 + \beta_1 \cdot x)}{\partial \beta_j} - \frac{(1 - y)}{1 - \sigma(\beta_0 + \beta_1 \cdot x)} \frac{\partial \sigma(\beta_0 + \beta_1 \cdot x)}{\partial \beta_j}$$

$$= \left[ \frac{y}{\sigma(\beta_0 + \beta_1 \cdot x)} - \frac{1 - y}{1 - \sigma(\beta_0 + \beta_1 \cdot x)} \right] \frac{\partial \sigma(\beta_0 + \beta_1 \cdot x)}{\partial \beta_j}$$

$$= \frac{y - y\sigma(\beta_0 + \beta_1 \cdot x) - \sigma(\beta_0 + \beta_1 \cdot x) + y\sigma(\beta_0 + \beta_1 \cdot x)}{\sigma(\beta_0 + \beta_1 \cdot x)(1 - \sigma(\beta_0 + \beta_1 \cdot x))} \cdot \sigma(\beta_0 + \beta_1 \cdot x)(1 - \sigma(\beta_0 + \beta_1 \cdot x)) \frac{\partial(\beta_0 + \beta_1 \cdot x)}{\partial \beta_j}$$

$$= \frac{y - y\sigma(\beta_0 + \beta_1 \cdot x) - \sigma(\beta_0 + \beta_1 \cdot x) + y\sigma(\beta_0 + \beta_1 \cdot x)}{\sigma(\beta_0 + \beta_1 \cdot x)(1 - \sigma(\beta_0 + \beta_1 \cdot x))} \cdot \sigma(\beta_0 + \beta_1 \cdot x)(1 - \sigma(\beta_0 + \beta_1 \cdot x)) \cdot x$$

$$= (y - \sigma(\beta_0 + \beta_1 \cdot x)) \cdot x$$

$$\frac{\partial \log l}{\partial \beta_j} \Rightarrow (\hat{y} - y)x$$

Stochastic Gradient Descent

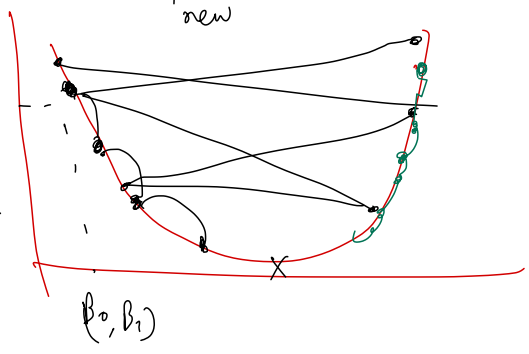
$$\beta = \beta - \eta \cdot \nabla l(\beta)$$

# Gradient Descent :

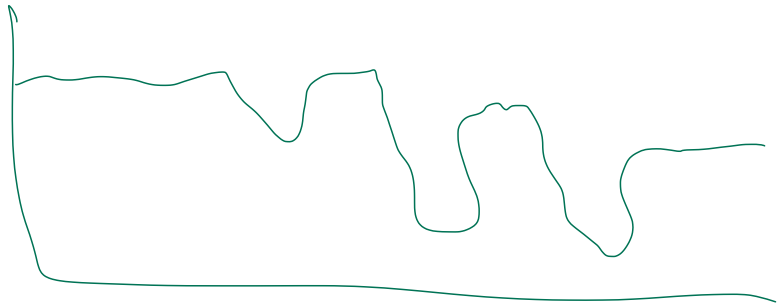
Loss

$$\beta_{\text{new}} = \beta_{\text{old}} - \eta \frac{\partial \text{Loss}}{\partial \beta}$$

learning rate -> high -> finding loss possibility of missing min.



low ->



BLUE -> estimator  
 Best Linear Unbiased Estimator

Dataset

Logistic Regression

$$\frac{[0, 1]}{0.6}$$

Pass (1)  
 vs  
 Fail (0)

2	1	3	4
-3	1	4	7
9	1	5	9

$$\sqrt{(3-2)^2 + (1-1)^2 + (4-3)^2 + (2-4)^2}$$

# Big Picture of ML modelling

Business Dis in → Question

"whom should be distribute our loans to"

Data

Relationship  
(# fr. w/ from)

Age, Tenure, Annual Income, Credit score, Gender, Yrs. of Educ,

Male - 1  
Female - 0

35	3	15	700	1	15	0
40	10	10	800	1	12	1

loan status: 0, 1 (No) (Yes)

X → K-1 →

loan\_details

"Dataset" (data.xlsx)

"Classification Model"

Logistic Regression  
↓  
(85%)

Decision Tree  
↓  
(90%)

SVM  
↓  
(88%)

KNN  
↓  
(70%)

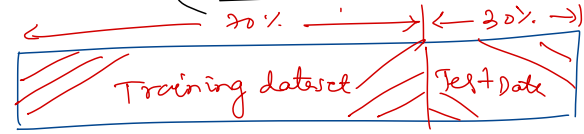
95%

mycode = 111  
df = sql & mycode

```

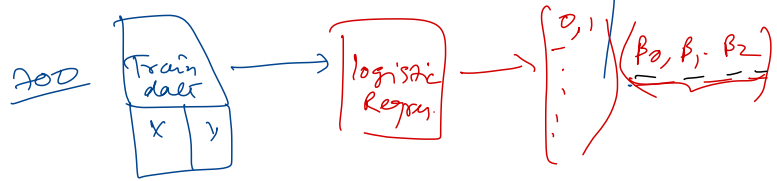
Select
Age, T, A, C, G, Y
from
loan_details
    
```

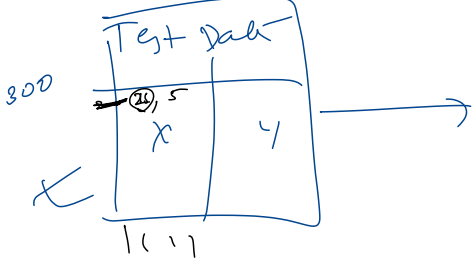
"Shuffle the dataset"



Missing values  
- Inputs missing value using mean

ratio: 70%





Learned Logistic Regression

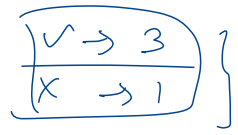
$$f(x) = 0.51 + (0.71 \times \text{Age}) + (2.0 \times \text{Temp})$$

2.5 + 5

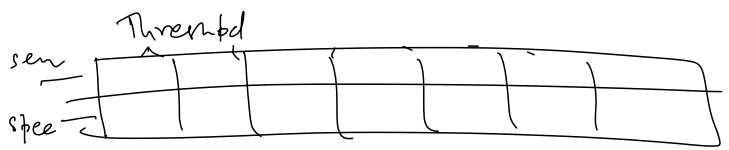
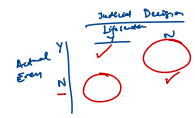
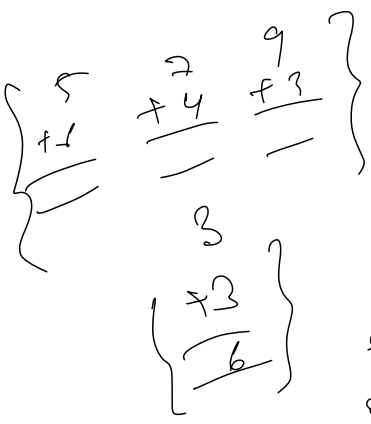
m/o	ALT
1	0
1	1
0	0
0	0

(0, 1)

only 0.66 → 1



$(\frac{3}{4}) = 0.75$



(Test dataset)  
"Ground Truth"

Real Label

0.5

"model output"

Predicted Label

		Positive	Negative
Predicted Label	Positive	True Positive (TP)	False Positive (FP)
	Negative	False Negative (FN)	True Negative (TN)

$$\text{Precision} = \frac{\sum TP}{\sum TP + \sum FP}$$

$$F1 = \frac{2 \times P \times R}{P + R}$$

$$\text{Recall} = \frac{\sum TP}{\sum TP + \sum FN}$$

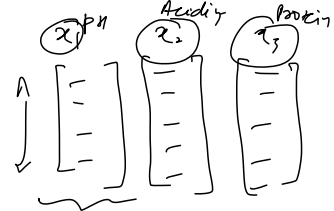
$$\text{Accuracy} = \frac{\sum TP + \sum TN}{\sum TP + \sum FP + \sum FN + \sum TN}$$



# Clustering

# (unsupervised)

(no labels)



$$S = \{x_1, x_2, \dots, x_n\} \text{ Dataset}$$

Given dataset and a number of clusters ' $k$ '

want's

Divide the dataset into ' $k$ ' groups

$$A_1, A_2, A_3, \dots, A_k$$

$k$  clusters / groups / partitions

Such that  
Want to have!

(I)  $A_i \neq \phi$  (empty)

[no clusters should be empty]

(II)  $A_i \cap A_j = \phi \quad \forall i, j$

[no two clusters should have common data]

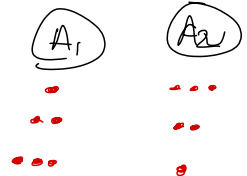
(III)  $\bigcup_{i=1}^k A_i = S$  (All points are associated with some or other cluster)

$$L(P(A_1, A_2, \dots, A_k)) = \sum_{i=1}^k \sum_{x \in A_i} d(x_i, A_i)$$

take point in  $A_i$  and find the distance b/w  $x$  and  $A_i$  mean ( $i$ th cluster) and do this for all such clusters

F No. of clusters = 2, no. of points 4

# of clusters  $\rightarrow \frac{2^{n-2}}{2} \Rightarrow \frac{2^{4-2}}{2} = 2$



$$S(n, k) = \frac{1}{k!} \sum_{i=0}^k (-1)^i \binom{k}{i} (k-i)^n$$

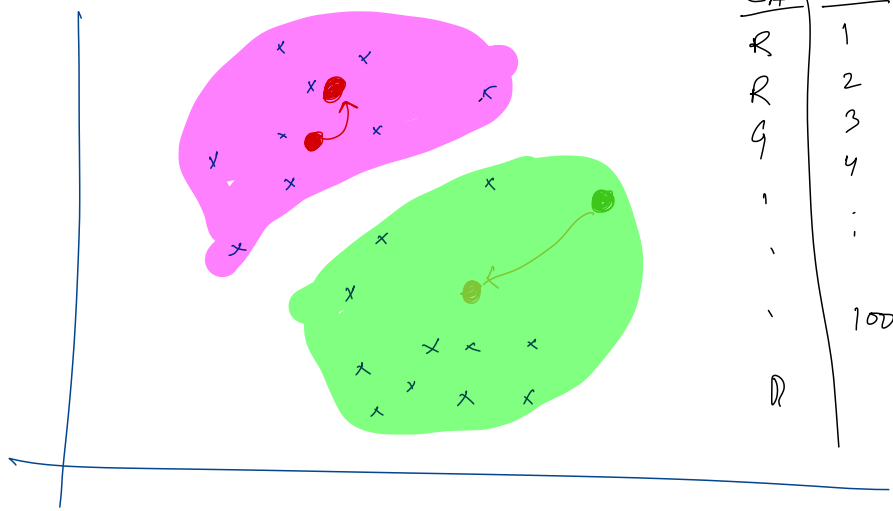
Cluster n data

points in

k clusters / partitions

(Stirling) number of  
2nd kind

$$= O(k^n)$$



CA	Pi	GA	RA
R	1	5	3
R	2	2	1
9	3	.	.
.	4	.	.
.	:	.	.
.	100	.	.
D			

Jaccard Similarity

$$A = \{3, 4, 7, 9\}$$

$$A \cap B = \{3, 4\}$$

$$B = \{1, 2, 3, 4\}$$

$$A \cup B =$$

$$A \cup B = \{1, 2, 3, 4, 7, 9\}$$

$$A = \{2, 3, 4\}$$

$$B = \{1, 3\}$$

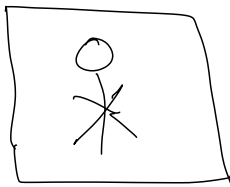
$$= \frac{0}{0}$$

$$= 1$$

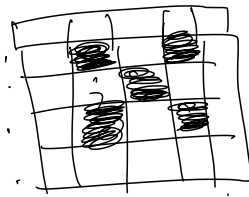
$$\frac{|A \cap B|}{|A \cup B|} = \frac{2}{6} = \left(\frac{1}{3}\right) = 33\%$$

$$= 1 - 33\%$$

$$= 67\%$$



Black/white



(5x5)

(25px)

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

(0 - 255)

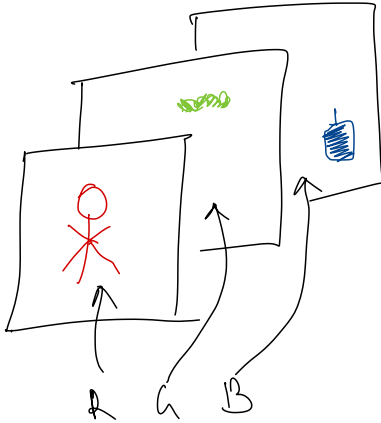
dark  
black

light  
white

gray scale  
images

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 70 & 50 & & \end{bmatrix}$$

[0, 255]

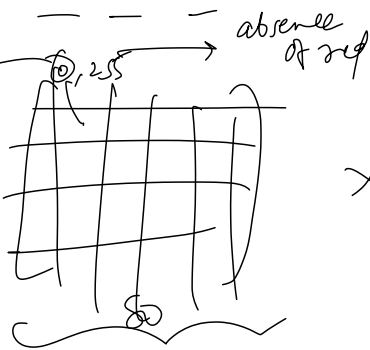


R G B



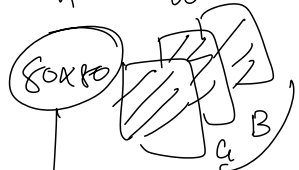
pure red

80

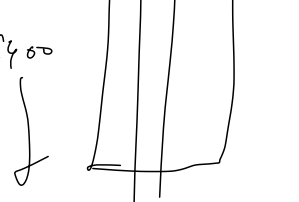


x 3

80 x 80 x 3  
h w r/g/b



6400



# Dimensionality Reduction

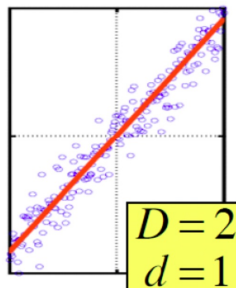
$$\rightarrow \begin{vmatrix} 2-\lambda & -1 \\ -1 & 2-\lambda \end{vmatrix}$$
$$\Rightarrow (\lambda - 2)^2 - (-1)(-1) = 0$$

$\Rightarrow$   $\lambda$ 's are e values

customer	day	We	Th	Fr	Sa	Su
		7/10/96	7/11/96	7/12/96	7/13/96	7/14/96
ABC Inc.		1	1	1	0	0
DEF Ltd.		2	2	2	0	0
GHI Inc.		1	1	1	0	0
KLM Co.		5	5	5	0	0
Smith		0	0	0	2	2
Johnson		0	0	0	3	3
Thompson		0	0	0	1	1

Goal of dimensionality reduction is to discover the axis of data!

Rather than representing every point with 2 coordinates we represent each point with 1 coordinate



Rank: No. of linearly independent ( $\frac{\text{columns}}{\text{rows}}$ ) of "A".

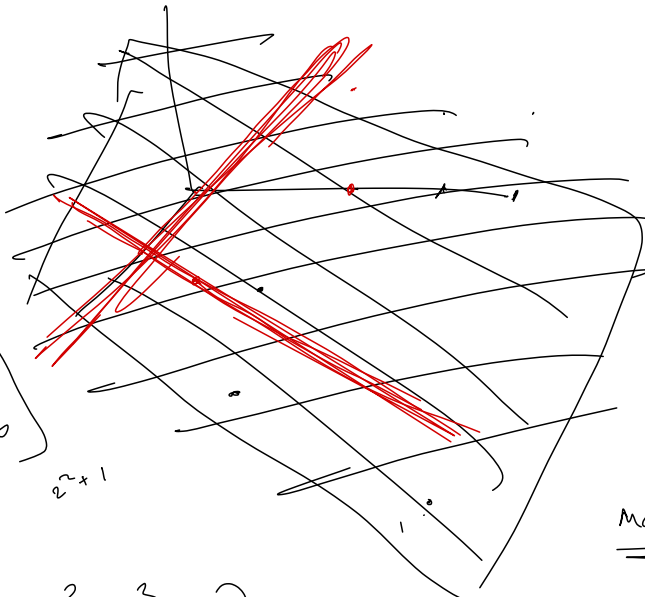
$$A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 1 & 3 & 5 \end{bmatrix}$$

$R_3 \leftarrow R_1 + R_2$   
 $\text{Rank}(A) = \textcircled{2}$

Basis:

$$\left[ \begin{array}{ccc} 1 & 2 & 3 \\ \hline \end{array} \right] \quad \left[ \begin{array}{ccc} 0 & 1 & 2 \\ \hline \end{array} \right]$$

new coordinates will be:  $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix}$



$$\begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad 2 \times 1$$

Proof:

$$a^2 + b^2 = (a+b)^2 - 2ab$$

$$\text{LHS} = a^2 + b^2 \times 0 = 0$$

$$\text{RHS} = \dots \times 0 = 0$$

Matrix Decomposition

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 0 & 1 & 2 \end{bmatrix} = \begin{bmatrix} A \end{bmatrix} \times \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$2 \times 3$        $3 \times 3$        $3 \times 3$

$$\begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

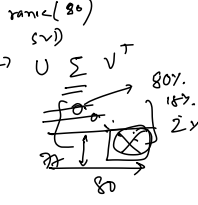
$$\begin{bmatrix} \frac{3}{10} & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0.5 \end{bmatrix}$$

$$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

$$\begin{aligned} &= 3^2 + 2^2 + \left(\frac{1}{6}\right)^2 \\ &= 9 + 4 + 0.25 \\ &= 13.25 \end{aligned}$$

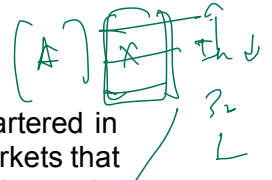
$$\begin{aligned} &\frac{(2)}{1.25} \quad ; \quad \frac{(2)}{12.25} \quad \frac{(0.5)}{1.5} \\ &\frac{9}{13.25} \quad \left| \quad \frac{1}{12.25} \quad \left| \quad \frac{1}{1 \times 12.25} \right. \right. \\ &\underline{67\%} \quad \underline{30\%} \quad \underline{3\%} \end{aligned}$$

$\ln \sigma$   
 SMP  
 (SMB)  
 95%



15L  
Turn  
369

$15 \times (1.30)^4 = 42L$  → 2L → 1000



## Data Driven Consumer Analytics at Superstore

**Superstore** is a multinational groceries and general merchandise retailer headquartered in Garden City, England. Shops of **Superstore** are larger, mainly out-of-town hypermarkets that stock nearly all product ranges, although some are in the heart of town centers and inner-city locations. Recently they have started a home shopping service through their website. **Superstore** observed a significant growth by opening the website operations.

0 → 269

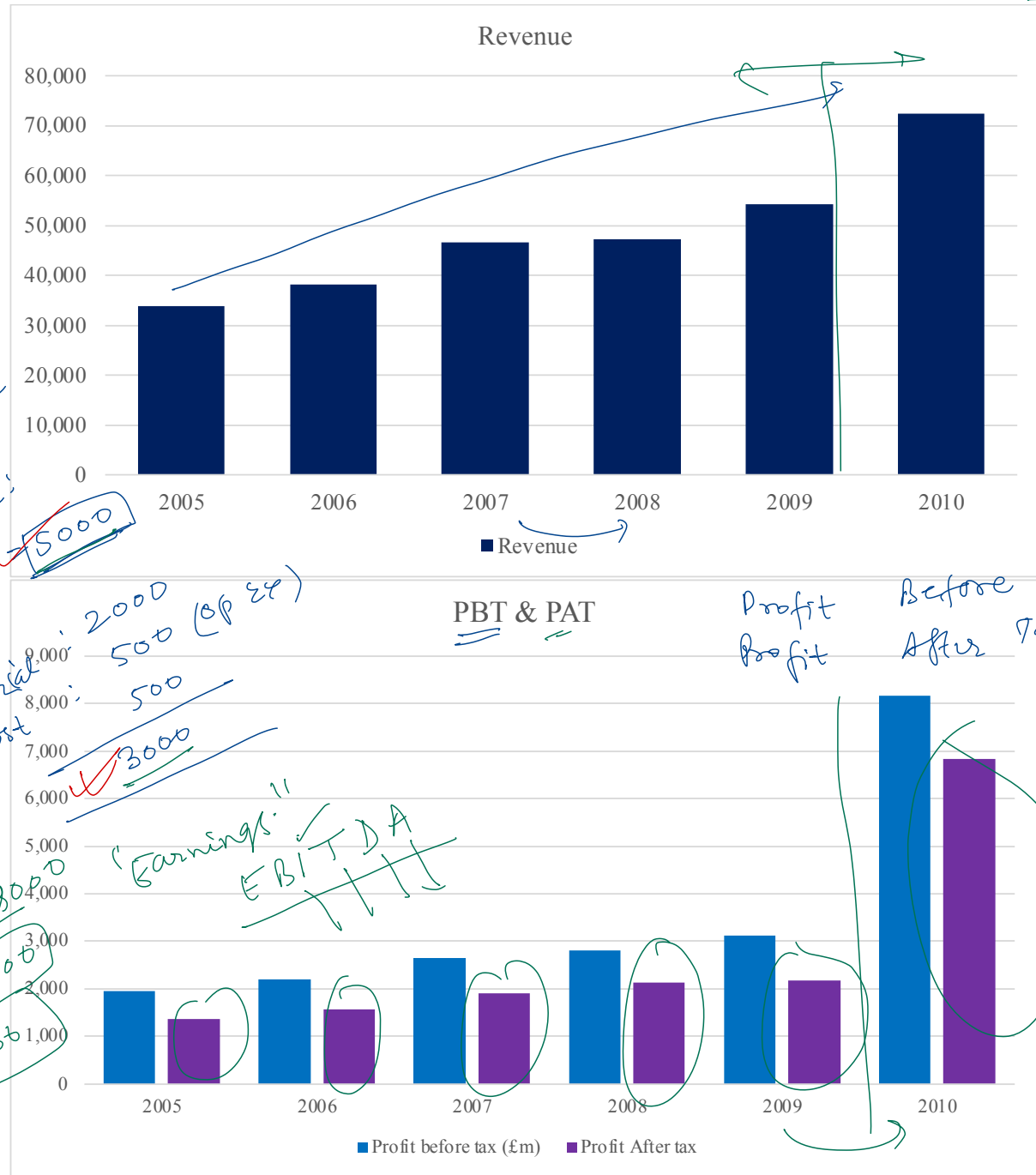


Figure 1 Revenue & Profit (in £ millions)

According to NY Bank retail analyst, "**Superstore** has pulled off a trick (online store) that I'm not aware of any other retailer had thought earlier". However due to internet penetration & democratization, many competitors supermarket chains have opened their website operations and this is driving lots of consumer leaving for other competing stores. This is known as

consumer **churn**. Formally, Churn is the phenomenon where customers of a business no longer purchase or interact with the business.

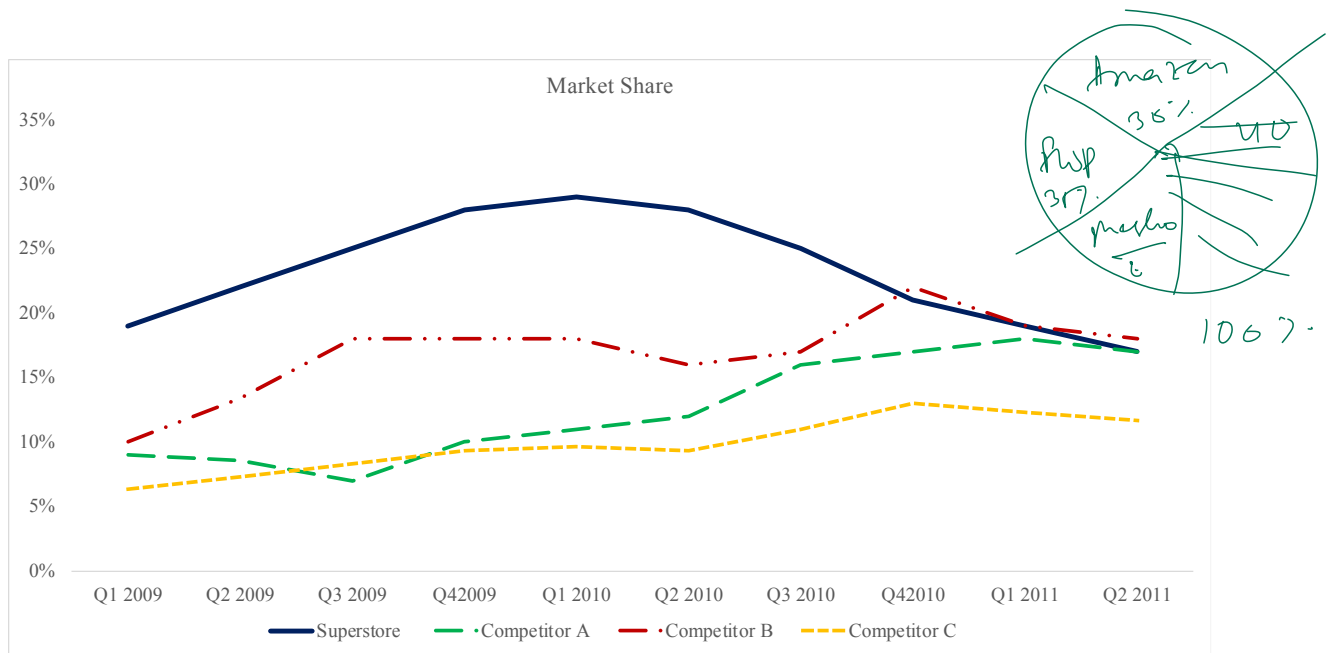


Figure 2 Market Share

**Superstore** consulted the MBB management consulting firm to help with their falling market share. MBB remarked “**Superstore** has plenty of consumer data. It is now imperative to use a data driven strategy to define a marketing campaign”.

You joined the company in the aftermath of the adaptation of a standardized solution that the consulting company suggested, and are tasked with **defining data driven marketing campaign to mediate consumer attrition**.

Your fundamental task is to ---

1. Understand the consumer base according to purchase & spending patterns and call out any specific trends/spikes.
2. Group your customers into high engaged, mid and low engaged. (Hint: There is a very intuitive RFM framework. See more details in Appendix 1.)
3. Build a data driven mechanism to identify which customers are likely to churn.
4. You are also given £25000 to spend for marketing campaigns. You would like to give some coupons as incentive for the customers to perform their next transaction (hence reducing the churn). Which all customers would you like to target.

Here is the purchase transaction data containing 540k rows <https://github.com/vntkumar8/musical-spoon/raw/main/Online%20Retail.xlsx>

**Snapshot of data**

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
536365	85123A	WHITE HANGING HEAI	6	01/12/10 8:26	2.55	17850	United Kingdom
536365	71053	WHITE METAL LANTER	6	01/12/10 8:26	3.39	17850	United Kingdom
536365	84406B	CREAM CUPID HEARTS	8	01/12/10 8:26	2.75	17850	United Kingdom
536365	84029G	KNITTED UNION FLAG	6	01/12/10 8:26	3.39	17850	United Kingdom



## Appendix 1

### **RFM Segmentation**

RFM segmentation is a method to identify groups of customers for special treatment. It is used as a tool to improve customer marketing.

### **What is RFM Segmentation?**

RFM analysis allows marketers to target specific clusters of customers with communications that are much more relevant for their particular behavior – and thus generate much higher rates of response, plus increased loyalty and customer lifetime value. Like other segmentation methods, an RFM model is a powerful way to identify groups of customers for special treatment. RFM stands for recency, frequency and monetary – more about each of these shortly. Marketers typically have extensive data on their existing customers – such as purchase history, browsing history, prior campaign response patterns and demographics – that can be used to identify specific groups of customers that can be addressed with offers very relevant to each.

While there are countless ways to perform segmentation, RFM analysis is popular for three reasons:

- It utilizes objective, numerical scales that yield a concise and informative high-level depiction of customers.
- It is simple – marketers can use it effectively without the need for data scientists or sophisticated software.
- It is intuitive – the output of this segmentation method is easy to understand and interpret.

### **What are Recency, Frequency and Monetary?**

Underlying the RFM segmentation technique is the idea that marketers can gain an extensive understanding of their customers by analyzing three quantifiable factors. These are:

- **Recency:** How much time has elapsed since a customer's last activity or transaction with the brand? Activity is usually a purchase, although variations are sometimes used, e.g., the last visit to a website or use of a mobile app. In most cases, the more recently a customer has interacted or transacted with a brand, the more likely that customer will be responsive to communications from the brand.
- **Frequency:** How often has a customer transacted or interacted with the brand during a particular period of time? Clearly, customers with frequent activities are more engaged, and probably more loyal, than customers who rarely do so. And one-time-only customers are in a class of their own.
- **Monetary:** Also referred to as "monetary value," this factor reflects how much a customer has spent with the brand during a particular period of time. Big spenders should usually be treated differently than customers who spend little. Looking at monetary divided by frequency indicates the average purchase amount – an important secondary factor to consider when segmenting customers.

YouTube tutorial: <https://youtu.be/i-HNJZeOOMY>

Churn: Netflix, Spotify, Apple Music, DOD, DFTT

Subscription based

cust-id → 1x

Amazon Ecommerce

PayTM

Assumptions

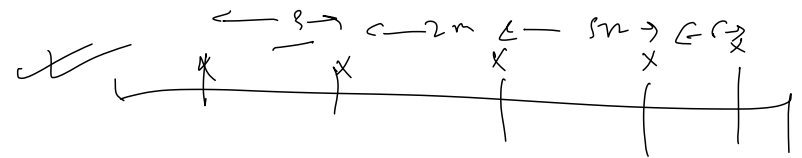
11 orders (5000)

1 Customer - 22 yrs (1000)

$$\frac{5000}{1000} = 5$$

(24m)

$$2021 \rightarrow \frac{(500)}{(100)} = 5$$

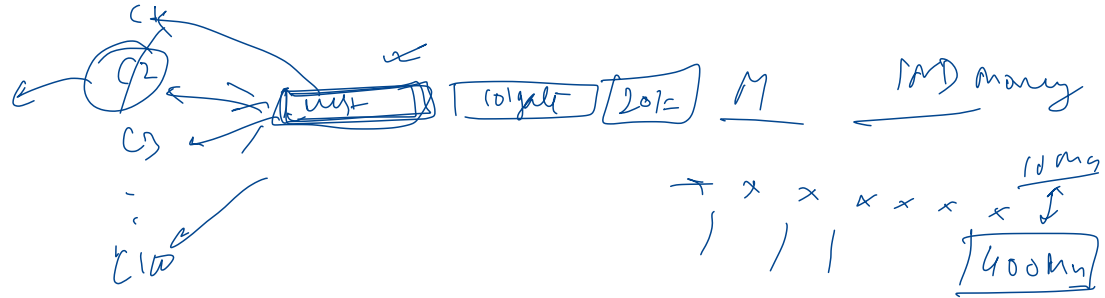
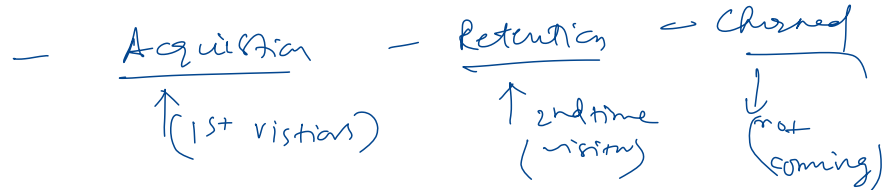


$$\frac{3, 2, 5, 6}{4} = 4$$

$$\frac{2, 3, 5, 6}{4} = 4$$

5%

Marketing Analytics



100% →

	1	2	3	4	5	6
May - 2010	✓					
Apr - 2010	✓					
Mar - 2010	✓					
Feb - 2010						X
Jan - 2010	100%	25%	20%	15%	4%	
Dec - 2009	100%	30%	20%	15%	5%	2%

Dec

How many active →

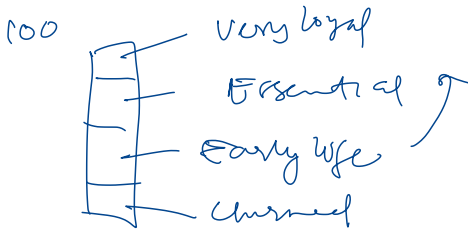
30% → Jan

20% → Feb

11% → Mar

RPM:

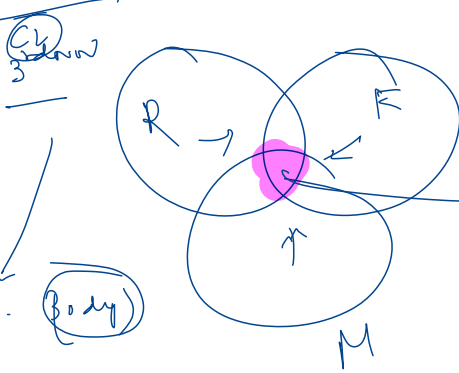
Method to group/segment your customers into bins



Recency: when did customer make its last purchase txn.

Frequency: how often does customer make their purchase

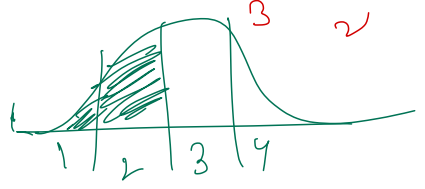
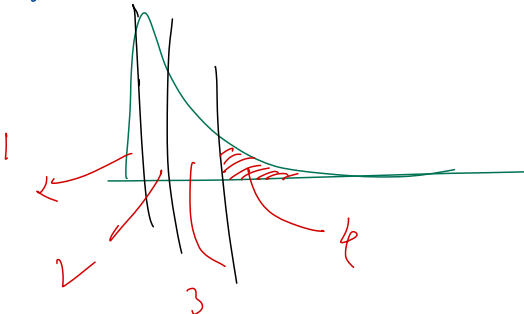
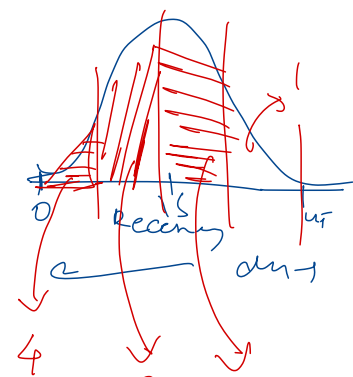
Monetary val.: how much money customer has spent



High val. Customer  
✓ higher spend,  
very frequently  
and last purchase  
was very recent!

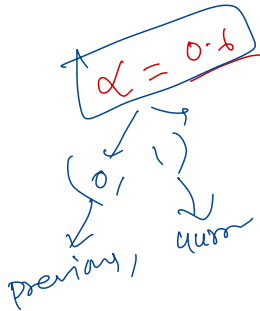
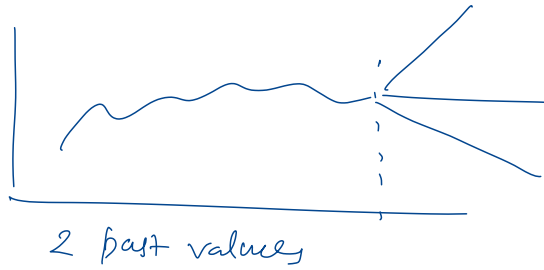
Recency ↓      frequency ↑      Monetary ↑

cust\_id, ~~id~~ item, date, price.  
0  
Recency 15 day  
to by (1) ← → Last order date (frequency)



# Exponential Smoothing

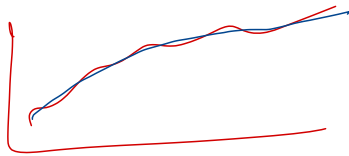
'Forecasting'



current past  
 [ 1, 3, 3, 6, 7, 8, 9, ... ]

$$\alpha \times \text{Current value} + (1-\alpha) \times \text{Past value}$$

$$0.6 \times 3 + 0.4 \times 3$$



$$0.6 \times 6 + 0.4 \times 3 = 3.6 + 1.2 = 4.8$$

$$0.6 \times 3 + 0.4 \times 3 = 1.8 + 1.2 = 3$$

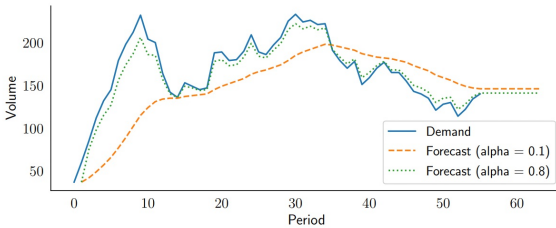


Figure 3.2: Simple smoothing

Simple/single exponential smoothing: This smoothing can be used for making forecasts based in a time series that has no trend and seasonality. Simple exponential smoothing does not do well when there is a trend in the data. Double exponential smoothing: This type of exponential smoothing comes with the support for trend components of time series.

## simple exponential smoothing.

In the math notation  $x_t$  is  $t$  is  $[t]$ , and  $\hat{x}_t$  is our prediction for  $x_t$ ;

Initial conditions

- $s_0 = x_0$ . This is our initial guess.
- $\hat{x}_0$  is undefined. We can't call the first guess a prediction since it's actually the first observation.

For  $t > 0$

$$\begin{aligned} s_t &= \alpha(x_t) + (1 - \alpha)s_{t-1} \\ \hat{x}_t &= s_{t-1} \end{aligned}$$

When  $\alpha$  is closer to 1 the model is more sensitive to recent observations. When  $\alpha$  is closer to 0 the model is more sensitive to past observations.

Simple/single exponential smoothing: This smoothing can be used for making forecasts based in a time series that has no trend and seasonality. Simple exponential smoothing does not do well when there is a trend in the data. Double exponential smoothing: This type of exponential smoothing comes with the support for trend components of time series.

Now we will implement double exponential smoothing. For our implementation the formula is as follows:

- $s_0 = x_0$
- $b_0 = 0$
- $\hat{x}_0$  is undefined

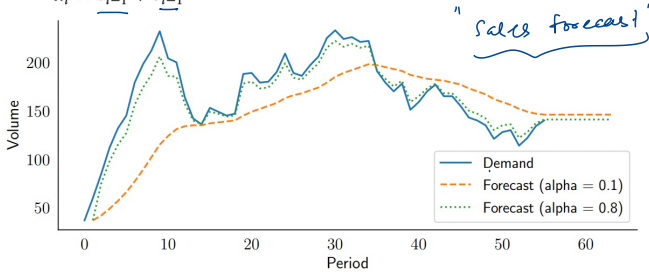
$$s_1 = \alpha x_1 + (1 - \alpha)(s_0 + b_0)$$

$$s_1 = \alpha x_1 + (1 - \alpha)(x_0 + 0)$$

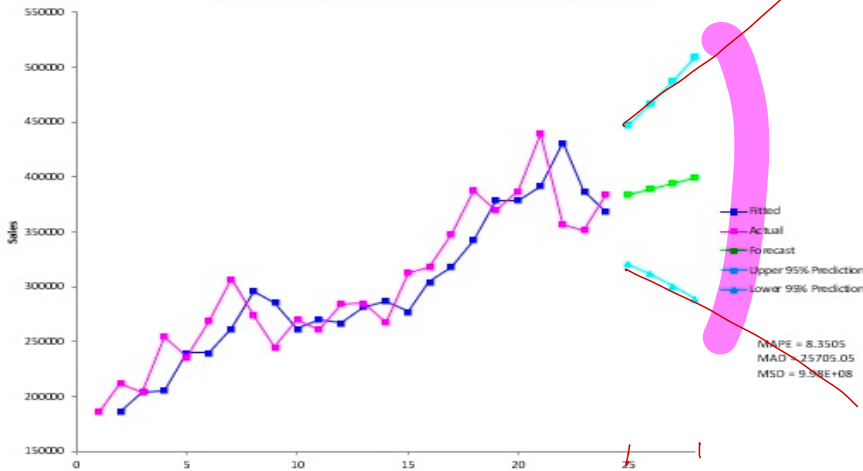
For  $t > 0$ :

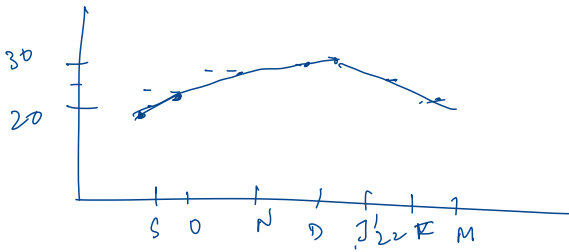
- $s_t = \alpha x_t + (1 - \alpha)(s_{t-1} + b_{t-1})$
- $b_t = \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1}$
- $\hat{x}_t = s_{t-1} + b_{t-1}$

$$s_{12} + b_{12} \Rightarrow (s_{12})_{b_{12}}$$



## Sales (Double Exponential Smoothing Model)





Consumer price index

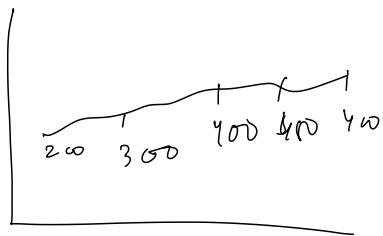
$(20, 24, 28, 29, 30, 29)$

$$\frac{24-20}{20}; \frac{28-24}{24}; \frac{29-28}{24}$$

$$\frac{28-20}{20}; \frac{29-24}{24}; \frac{30-28}{28}$$

$$\frac{1}{5}; \frac{1}{6}; \frac{1}{24}; \dots$$

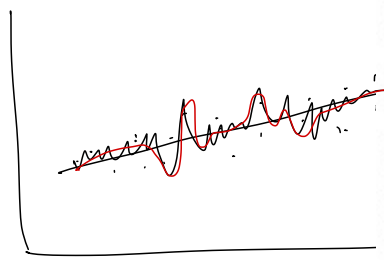
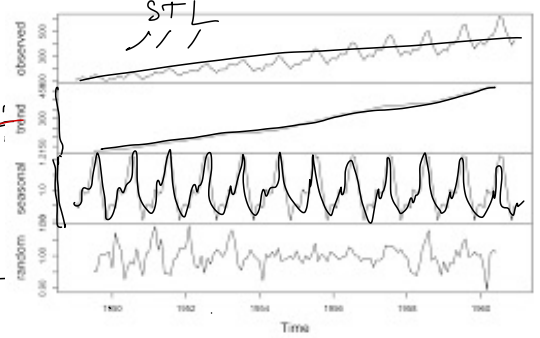
$$0.20\%; 12\%; \dots$$



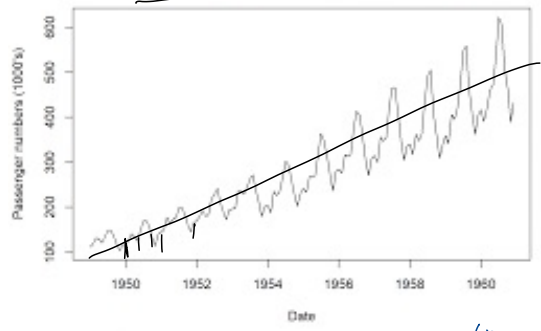
⊙  
- trend

Seasonality

Decomposition of multiplicative time series



Air Passenger numbers from 1949 to 1961



- Seasonality - MA

$R^2 \rightarrow (0.1) \rightarrow 9\%$   
 $MSE \rightarrow$  lower the better

Actual: 90, 150, 250, 400  
 Predicted: 100, 200, 280, 300

$$MSE = (100-90)^2 + (200-150)^2 + (280-250)^2 + (300-400)^2$$

Data  $\rightarrow$  Model  $\rightarrow$   $\sigma/P$

# Double Exp. Smoothing:

$$(\alpha, \beta) \rightarrow (0.2, 0.4, 0.6, 0.8)$$

$$(\underbrace{0.1, 0.2, 0.3, 0.4, 0.8, 1.1}_4)$$

$$\rightarrow \alpha \times 4 = \boxed{10}$$

$\alpha =$

$\alpha$	$\beta$	MSE
0.1	0.2	110
0.1	0.4	210
0.1	0.6	300
0.1	0.8	80
⋮	⋮	⋮

(Predicted - Actual)

