

S	M	T	W	T	F	S
5	1	2	3	4	5	
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

MAR - 2016

FEBRUARY • MONDAY

29

WK10 • 060-306

Hypothesis

A premise or claim that we want to test (or) experiment

Null hypothesis: (In statistics) Null hypothesis means default hypothesis

$H_0 \rightarrow$ currently accepted value for a parameter

Alternative hypothesis: (H_a) \rightarrow also called Research hypothesis
 Involves the claim to be tested

eg:

It is believed that a candy machine makes chocolate bars that are on average 5g. A worker claims that the machine after maintenance no longer makes 5g bars.

Write H_0 and H_a ?

(default assumption) $H_0: \mu = 5g$

H_0 & H_a are mathematical opposites

$H_a: \mu \neq 5g$

possible outcomes of this test:

- Reject Null hypothesis H_0

- Fail to Reject Null hypothesis H_0

2016

01

TUESDAY • MARCH

WK10 • 061-305

M	T	W	T	F	S
7	8	9	10	11	12
14	15	16	17	18	19
21	22	23	24	25	26
28	29	30	31		

MAR - 2016

Test statistic - calculate from sample data
used to decide

- Ex) Take samples of 50 bars
- find Average Val
 - calculate test statistic

statistically significant - where do we draw the line to make a decision.

from 3 Random samples we are getting

Averages as below

- 1) Avg = 5.21 → more or less equal to H_0
 2) Avg = 5.79 so we can't reject H_0 .
 3) Avg = 7.21 → little far away from H_0

much far away from 5. So

we can definitely Reject H_0

so we could give thought of Rejecting H_0

Level of Confidence: $C = 95\%, 99\%$

How confident are we in our decision.

Level of Significance $\alpha = 1 - C$

so Level of Confidence $C = 95\%$

$$\Rightarrow \alpha = 1 - C$$

$$= 1 - 0.95 = 0.05$$

$$\boxed{\alpha = 0.05}$$

2016

S	M	T	W	T	F	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30					

Jan
APR - 2014

MARCH • WEDNESDAY

02

WK10 • 062-304

Basic Statistics

1) Descriptive statistics

2) Inferential statistics

1) Descriptive statistics

↳ methods of summarizing the information we have collected for an analysis.

We can summarize information by means of graphs, such as a pie chart, bar chart or numbers such as a mean, percentage, correlation coefficient.

2) Inferential statistics:

↳ It is about drawing conclusions about a population on the basis of only a limited number of cases.

Data and visualization

E.g.: Football

⇒ goals, winners, stopped penalties

Variables ⇒ characteristics of something or someone
cases ⇒ something or someone

2016

03

THURSDAY • MARCH

WK10 • 063-303

M	T	W	T
7	8	9	10
14	15	16	17
21	22	23	24
28	29	30	31

MAR-2016

- player themselves are
with both sides
- 9 [cases] player 1 player 2 player 3 ...
variables individual & variables environment
- 10 character body weight age ...
+ stus color ...
- 11 [variables] goals
- 12 - players matches behavior game over
- * A case can be individual player or Team
- * e variable means features (w) characteristics of those cases (player or team)
- * Variable need to vary (variables should not be same, e.g.) country of the team, here the case is particular Team. For the particular team has many players and all these players are from same country only, so we can treat the "country" feature/characteristics as variable, that feature is called "constant".

of Levels of measurement:

diff kinds of variable representing strongly diverging characteristics, for this reason it's of essential importance to distinguish diff kinds of measurement

2016

S	M	T	W	T	F	S
1	2					
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30

APR - 2016

MARCH • FRIDAY

04

WK10 • 064-302

nominal level : made up of various categories that differ from each other. There is no order, however. (means no ranking order)

⇒ It's not possible to argue out that one category is better or worse, or more, or less than another.

E.g 1 : nationality of the football players.

[Spanish, French, Italy, German, India, USA]

There is no order ranking order but differ from each other

E.g 2 : Gender of the football players, city the football meeting come from.

ordinal level : There is not only the difference b/w the categories of the variable, there is also an order.

E.g 1 : winners of the competition. which team is no 1, no 2, no 3. [1, 2, 3].

It's differ from each other and has order, but we don't know anything about the different b/w categories means how much the no 1 was better than no 2.

categorical variables :

both nominal and ordinal levels are called categorical variables.

2016

05

SATURDAY • MARCH

WK 10 • 065-301

M	T	W	T	F	S
1	2	3	4	5	
7	8	9	10	11	12
14	15	16	17	18	19
21	22	23	24	25	26
28	29	30	31		

MAR-2016

Interval level: different categories and has ranking order + also has similar intervals b/w the categories.

e.g. 1: Age of the football players

[player 1 \geq 16, player 2 \geq 18,
player 3 \geq 12, player 4 \geq 10]

{ 16, 18 } differs and we can say player 2 is older than

{ 10, 12 } same players 1

ratio level: it is similar to interval levels but in addition has a meaningful zero point.

e.g. 1: players body height measured in cm's

A height of 0 cm means there is no height at all.

Note: age cannot be 0 as there is no zero point. So therefore age is an interval variable.

06 SUNDAY

Quantitative variables:

both interval and ratio levels are called quantitative variables, because the categories are represented by numerical variables. (or) values

2016

S	M	T	W	T	F	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30					

APR - 2016

MARCH • MONDAY

07

WK11 • 067-299

Quantitative variables

Discrete

Continuous

A set of separate numbers & infinite region of values.

e.g.: no of goals scored by a variable is continuous

by Football players. it the possible values of

can be 1, 2, 3, ... but for the variables born an

but we can't say 1.23, 2.4... interval.

e.g.: height of the players
can be 170 cm, 170.25 cm,
170.0245 cm.& we don't have separate
numbers but an infinite
region of values.

Statistical methods

Difference Rank order Similar Intervals meaningful zero point

Categorical nominal + - - -

ordinal + + - -

Quantitative interval + + + -

ratio + + + +

→ Discrete - set of separate numbers

→ continuous - infinite region of values

2016

08

TUESDAY • MARCH

WK11 • 068-298

M	T	W	T	F	S
-	-	1	2	3	4
7	8	9	10	11	12
14	15	16	17	18	19
21	22	23	24	25	26
28	29	30	31		

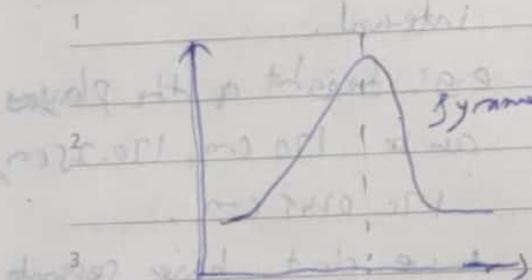
MAR - 2016

Frequency Table!

Shows the how the values are distributed over the cases.

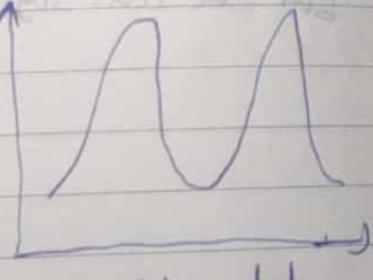
Graphs:

Qualitative variables categorical \Rightarrow pie, bar chart
 Quantitative variables \Rightarrow histograms
 (Continuous variable)

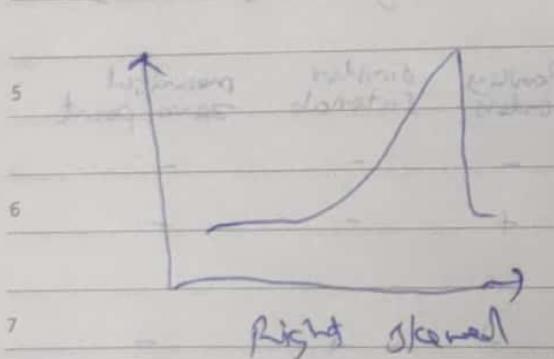


Similar to Bell shape

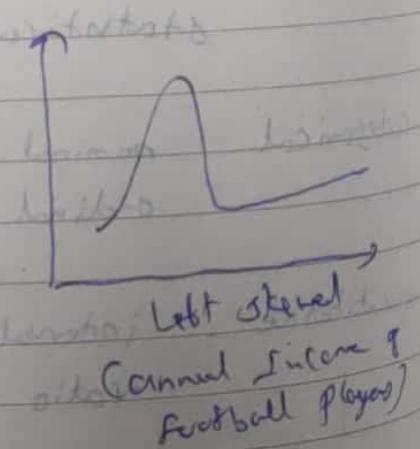
Unimodal



Bimodal



Right skewed



(Example: Income of football players)

measures of central tendency:

Median, mode, mean

\rightarrow describe the center of a distribution

Median, mode, mean

2016

S	M	T	W	T	F	S
1	2					
4	5	6	7	8	9	
3	12	13	14	15	16	
11	18	19	20	21	22	23
17	25	26	27	28	29	30
24						

APR - 2016

MARCH • WEDNESDAY

09

WK11 • 069-297

→ graph that nicely presents the variability of a distribution is the "box plot".

summary

Summarizing a distribution graphs

center

mode

→ categorical

median → influential
outlier/skewed

quantitative

median

$\bar{x} = \text{Ex}$

Information about the variability of the data!

↳ range

⇒ $Q_3 - Q_1$ (last value - first value)

↳ interquartile range ⇒ IQR: $Q_3 - Q_1$

↳ leaves out the extreme values

maximum value that's not an outlier

IQR

Q_3

Q_2 (median)

→ whisker

minimum value that's not an outlier

2016

10

THURSDAY • MARCH

WK11 • 070-296

M	T	W	T	F	S	S
1	2	3	4	5	6	7
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

MAR - 2016

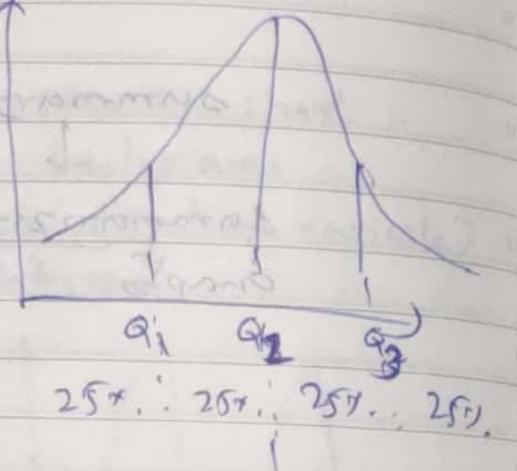
center of distribution + variability of distribution

of a distribution = more complete picture

values ordered from low to high

0, 5.6, 8.7, 14.1, 14.1,

(15), 17.2, 19.2, (19.3), 24.1, 27.7

Q₂ (median)Q₃

25%, 25%, 25%, 25%

outliers: < median

= value lower than $Q_1 - 1.5 \text{ (IQR)}$ (or)= value higher than $Q_3 + 1.5 \text{ (IQR)}$]

$$\text{variance } (\sigma^2) = \frac{\sum (x - \bar{x})^2}{n-1}$$

Longer variability means \rightarrow sample size

the more the values are spread out the around the mean

Standard deviation $s = \sqrt{\sigma^2} = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$

(most often measure of dispersion)
 & the average distance of an observation
 from the mean

2016

S	M	T	W	T	F	S
1	2					
4	5	6	7	8	9	
11	12	13	14	15	16	
18	19	20	21	22	23	
25	26	27	28	29	30	

APR - 2016

MARCH • FRIDAY

11

WK11 • 071-295

Z-scores

sometimes we want to know if a specific observation is common or exceptional. To answer that question, they express a score in terms of the number of standard deviations it is removed from the mean. This is called Z-score.

if we recode original scores into Z-scores we say that we standardize the variable

recode original scores into Z-scores

↓

standardization

replace the original scores by standard deviation (s) from the mean (\bar{x})

+ easy to see whether a specific score is relatively common or exceptional

mean \bar{x} = balance point

Z-scores

(-) scores mean below the mean value
(+) scores mean above the mean value.

$$z = \frac{x - \bar{x}}{s}$$

2016

12

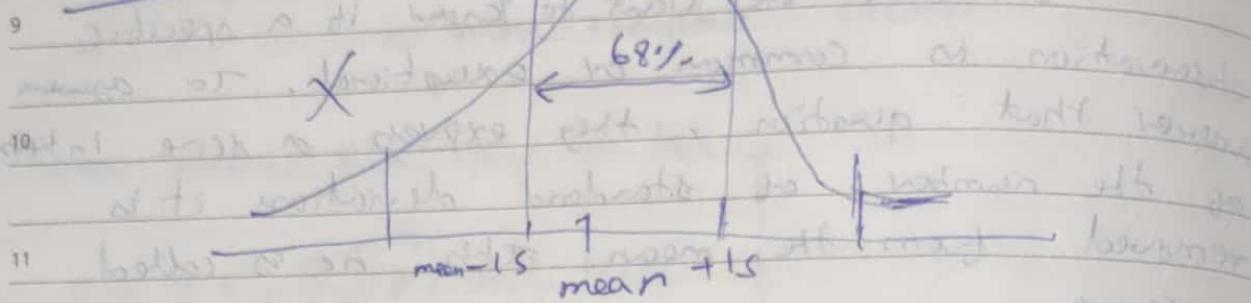
SATURDAY • MARCH

WK11 • 072-294

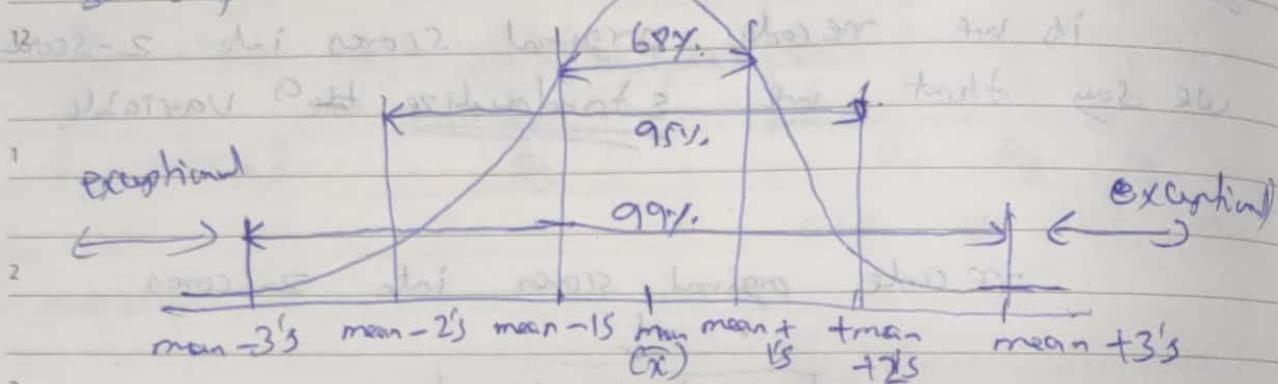
M	T	W	T	F	S
1	2	3	4	5	6
7	8	9	10	11	12
14	15	16	17	18	19
21	22	23	24	25	26
28	29	30	31		

MAR - 2016

Rule:



Bell shaped

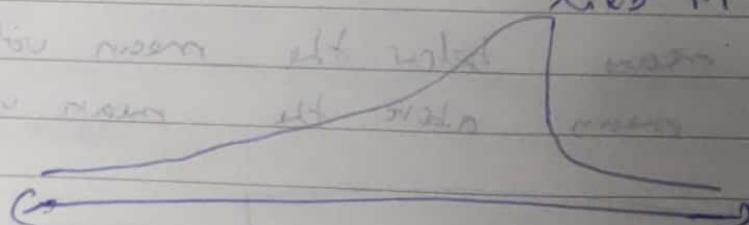


exceptional

exceptional

exceptional

13 SUNDAY



large positive z-scores

are more common b/c

more extreme values lie in Right side

2016

M	T	W	T	F	S
1	2				
5	6	7	8	9	
11	12	13	14	15	16
18	19	20	21	22	23
25	26	27	28	29	30

APR - 2016

MARCH • MONDAY

14

WK12 • 074-292

* Z-scores are more useful to compare different distributions.

Example:

City with 8 high schools, average grade for chemistry has b/w 6 to 10

School name Average grade chemistry

School 1

7.4

" 2

7.9

" 3

9.1

" 4

8.1

" 5

6.2

" 6

7.1

" 7

7.4

" 8

6.7

1) What does the distribution of the variable "average grade for chemistry" looks like?

2) What is the center of the distribution?

3) The variability of the distribution?

4) Construct a box plot

5) What is the z-score of School #3?

2016

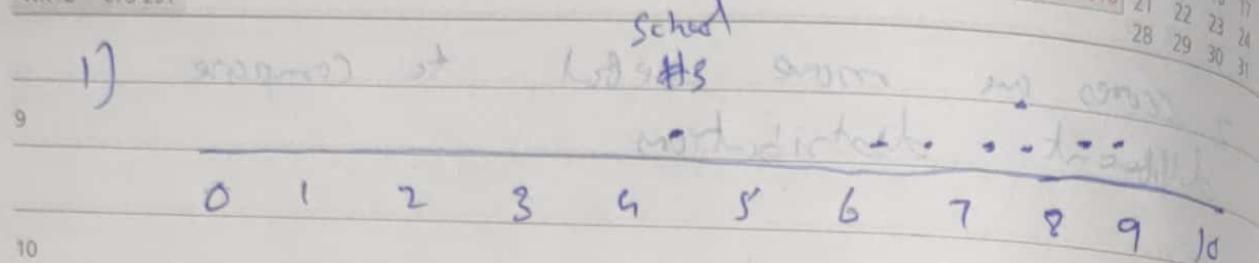
15

TUESDAY • MARCH

WK12 • 075-291

M	T	W	T
7	8	9	3
14	15	16	10
21	22	23	17
28	29	30	24

MAR - 2016



10) 2) mode $\Rightarrow 4.1, 6.2, 6.7, \underline{7.1}, \underline{7.4}, 7.9, 8.1$
 $\Rightarrow 7.4$ occurs twice

11) median $\rightarrow \frac{7.1 + 7.4}{2} = 7.25$

12) mean $\rightarrow \frac{\sum x}{n} = \frac{54.9}{8} = 6.86$

13) 3) Range \rightarrow highest value - lowest value
 $8.1 - 4.1 = 4$

14) IQR $\Rightarrow [4.1, 6.2, 6.7, 7.1] | [7.4, 7.4, 7.9, 8.1]$

15) Median $\rightarrow \frac{6.2 + 6.7}{2} = 6.45$

16) Q₁ $\rightarrow \frac{4.1 + 6.2}{2} = 5.15$

17) Q₃ $\rightarrow \frac{7.4 + 7.9}{2} = 7.65$

18) $IQR = Q_3 - Q_1 = 7.65 - 6.45 = 1.2$

19) Standard deviation (s) $= \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$

$s = 1.27$

2016

S	M	T	W	T	F	S
1			2			
4	5	6	7	8	9	
11	12	13	14	15	16	
18	19	20	21	22	23	
25	26	27	28	29	30	

APR - 2016

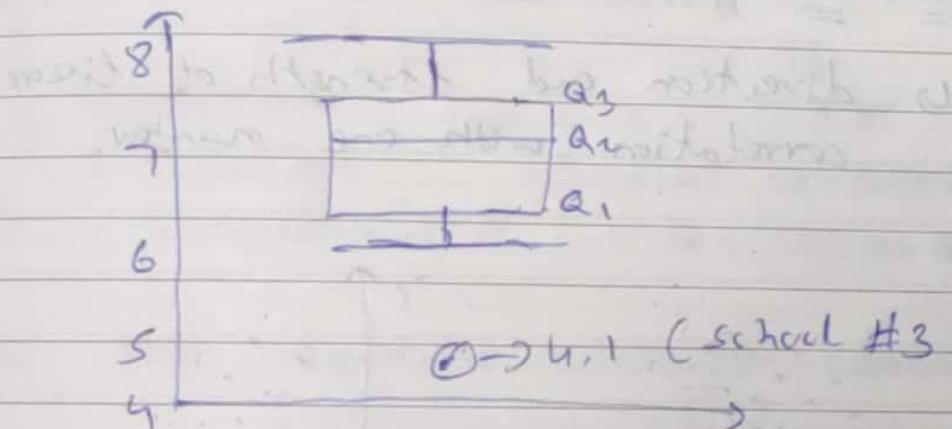
MARCH • WEDNESDAY

16

WK12 • 076-290

4) $Q_1 = 6.45$, $Q_2 = \text{median} = 7.25$, $Q_3 = 7.68$,
 $IQR = 1.23$

outlier = value lower than $Q_1 - 1.5(IQR)$
 or higher than $Q_3 + 1.5(IQR)$



5) $Z = \frac{x - \bar{x}}{s} = \frac{4.1 - 6.86}{1.23} = -2.17$

School #3



average of 4.1



lies 2.17 standard deviation below the mean and can therefore be conceived of as a rather exceptional value.

2016

17

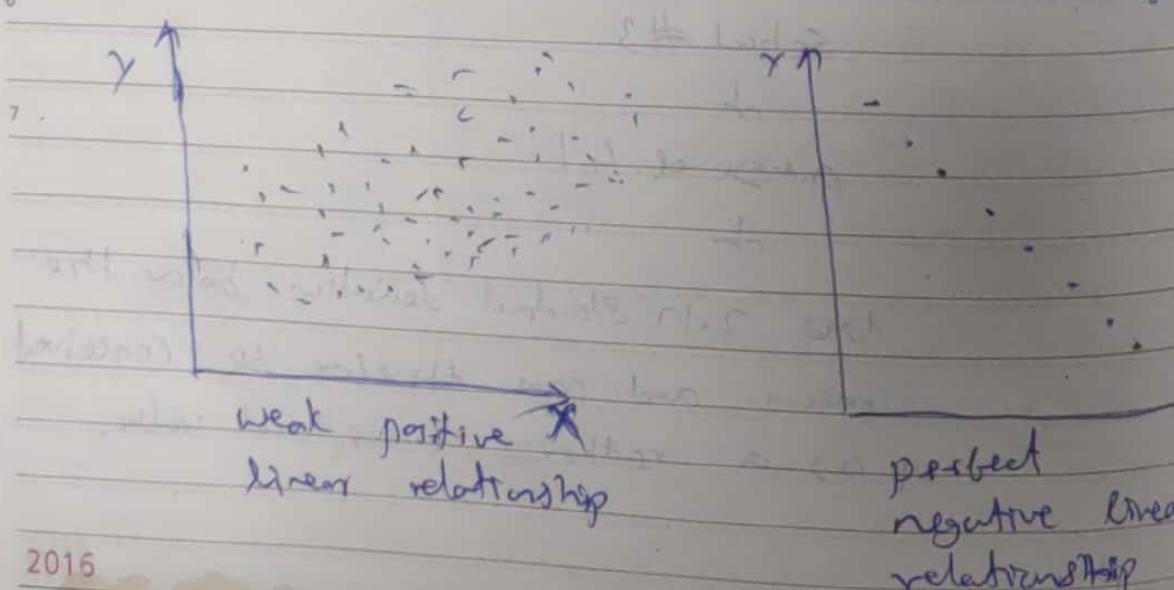
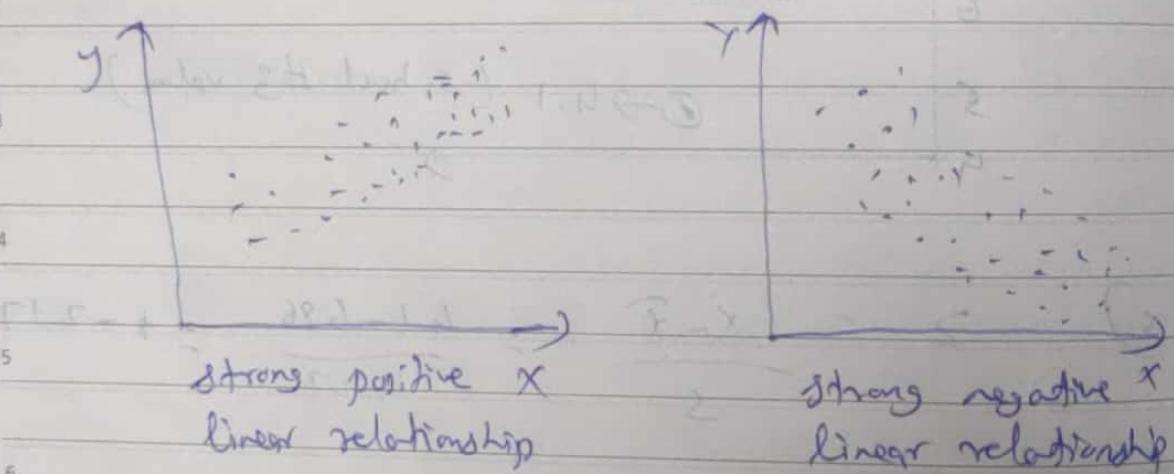
THURSDAY • MARCH

WK12 • 077-289

MAR - 2016

M	T	W	T	F	S
1	2	3	4	5	
7	8	9	10	11	12
14	15	16	17	18	19
21	22	23	24	25	26
28	29	30	31		

- 9 Contingency table → nominal variables
ordinal
- 10 scatterplot → quantitative variables
- 11 Pearson's R (correlation)
- 12 =
↳ direction and strength of linear correlation with one number.



2016

S	M	T	W	T	F	S
1	2					5
3	4	5	6	7	8	9
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	

APR - 2016

MARCH FRIDAY

18

WK 12 • 078-288

Y ↑
 X ↓
 Pearson's R shows a linear relationship between two variables.

Pearson's R \Rightarrow how strong or weak correlation.

↳ direction \rightarrow strength
 + = positive \rightarrow -1 = perfect negative
 - = negative \rightarrow +1 = perfect positive

e.g. body weight vs chocolate consumption.

(y) (x)
 (dependent variable) (independent variable)

Person	X	Y	ZX	ZY	$ZX \cdot ZY$
P1	50	50	-1.01	-1.15	1.17
P2	100	70	-0.56	-0.07	0.04
P3	200	70	0.34	-0.07	-0.02
P4	300	95	1.24	1.29	1.60

Pearson's R \Rightarrow
$$r = \frac{\sum ZX \cdot ZY}{n-1}$$

$r = 2.78/3 = 0.93$

2016

19

SATURDAY • MARCH

WK12 • 079-287

MAR-2016

M	T	W	T	F
1	2	3	4	
7	8	9	10	11
14	15	16	17	18
21	22	23	24	25
28	29	30	31	

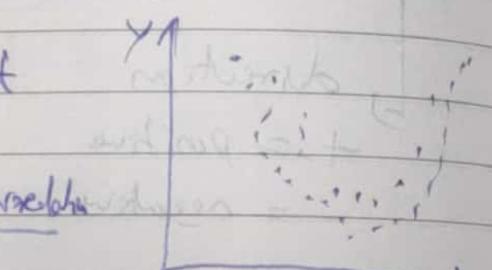
so for the example $r=0.93$ says that x and y are strongly correlated, i.e. strong positive linear relationship.

Note: check scatterplot before calculate pearson's

note: e.g.; when curvilinear case it doesn't work

if $r^2 = 0.15$ means not weak negative correlation.

but if it is weak linear correlation



$r^2 = 0.15 = \text{weak linear correlation}$

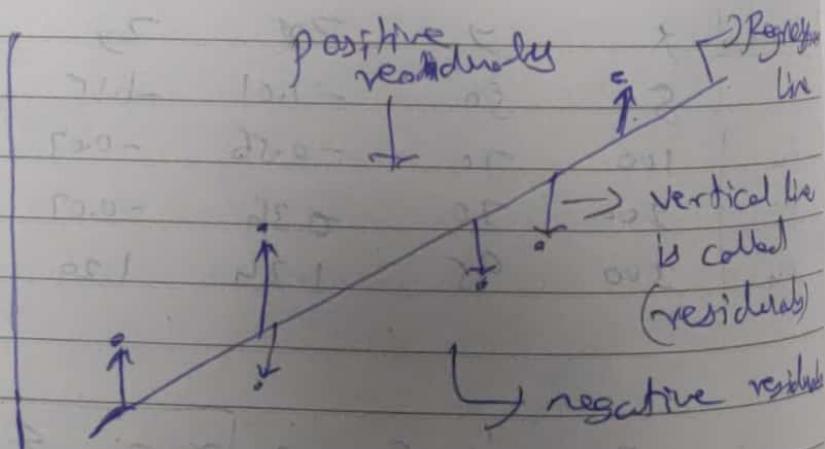
Regression:

Find Regression Line

\rightarrow $y = f(x)$

line with the smallest sum of squared residuals

20 SUNDAY



OLS \Rightarrow (ordinary Least square Analysis)

\Rightarrow used to find the regression line

2016

S	M	T	W	T	F	S
1	2					
3	4	5	6	7	8	9
11	12	13	14	15	16	
18	19	20	21	22	23	
25	26	27	28	29	30	

APR - 2016

MARCH • MONDAY

21

WK13 • 081-285

describes

Regression Line: linear relationship between two variables.

Compute Regression line: method of least squares

$$\hat{y} = a + bx$$

$$b = r \left(\frac{s_y}{s_x} \right) \quad \text{standard deviation of } y$$

$$a = \bar{y} - b(\bar{x})$$

R squared (R^2)

↳ tells you how much better a regression line predicts the value of a dependent variable than the mean of the variable

↳ tells you how much of the variance in your dependent variable is explained by your independent variable.

e.g., $r^2 = 0.69$ means \Rightarrow prediction error is 69%. smaller than when you use the mean
 \Rightarrow 69% of the variance in the values can be predicted by the previous values.

the higher the value means higher variance 2016

S	M	T	W	F	S
1	2	3	4	5	6
7	8	9	10	11	12
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30

Jan
MAR - 2016

MARCH • WEDNESDAY

23

WK13 • 083-283

Randomness:

Randomness it is not a property of a phenomenon.

humans → are not very good in assessing randomness.

- over interpretation of randomness (apophenia)
- bad in constructing random data.

probability:-

→ probability is way to quantify randomness.

→ $0 \leq \text{probability} \leq 1$

→ $\sum \text{probability} = 1$

* important terms:

'experiment', 'event', 'independent trials',
'relative frequency'.

→ law of Large numbers

→ keep calm and carry on many trials.

sample space, events & tree diagrams:

sample space → all possible outcomes for the experiment (or) collections of all possible elementary (or) combined outcomes

=) events outcomes for a random phenomenon

(or)

subset of sample space

2016

24

THURSDAY • MARCH

WK13 • 084-282

MAR-2016

M	T	W	T	F	S	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

★ note:

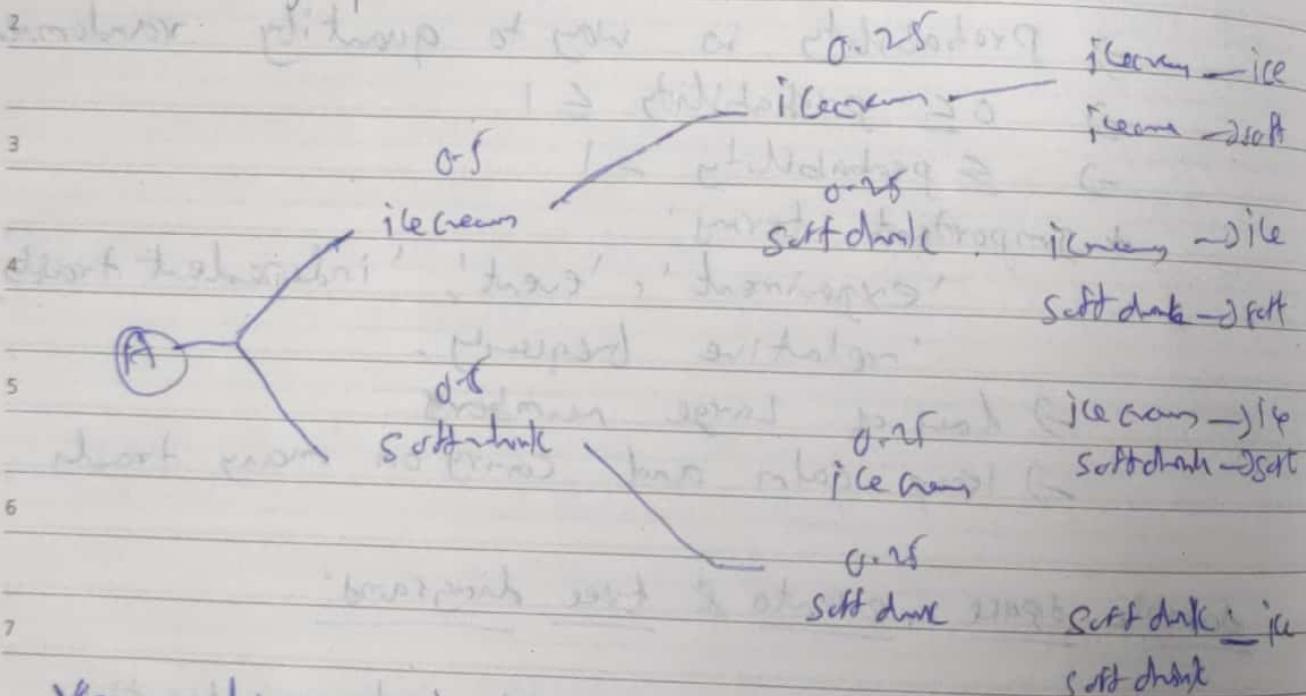
9 I never know whether your assessment of the probability was correct only three trials.

10 +
11 no accurate probability estimate
remember

12 e.g.: i) Box have 2 softdrinks & ice cream,

13 ii) There are 4 persons waiting to take it from box

1 what is the probability of for 5th person's pick



Venn diagrams! following No 1 = soft drink

overlap = joint probability

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

2016

Note: joint means intersection

S	M	T	W	F	S
1	2				
3	4	5	6	7	8
9	10	11	12	13	14
15	16	17	18	19	20
21	22	23	24	25	26
27	28	29	30		

APR-2016

MARCH • FRIDAY

25

WK13 • 085-281

Joint and marginal probabilities

e.g.-	Gender	rest	active	swim	Total
	male	34	12	5	51
	female	45	8	9	62
	Total	79	20	14	113

into percentage

0.301	0.106	0.044	0.451
0.397	0.071	0.080	0.549
0.699	0.177	0.124	(Total)

Joint probabilities:

→ probabilities for the intersection or certain outcomes of the variables

Marginal probabilities:

→ probability for an outcome of each individual variable

Conditional probability : $P(A|B)$

the probability of an event, given that another event occurs.

$$P(\text{A given B}) = \frac{P(A \text{ and } B)}{P(B)}$$

2016

26

SATURDAY • MARCH

WK13 • 086-280

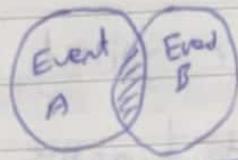
MAR - 2016

M	T	W	T	F	S	S
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

$$P(A|B) \leq P(B)$$

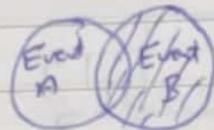
9

$$P(A|B) =$$



$$P(A \text{ and } B)$$

10



$$P(B)$$

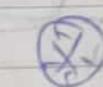
11

12 Independence \Rightarrow knowing the outcome of one event does not influence your knowledge about the outcome for the other.

1 disjoint \Rightarrow if one event occurs the others (dependent) cannot occur. So the joint probability is zero.

Variables whose possible values are numerical outcomes

Random Variable \rightarrow of a random phenomenon



probability distribution

continuous (infinite no. of possible values)

discrete (countable no. of distinct values)

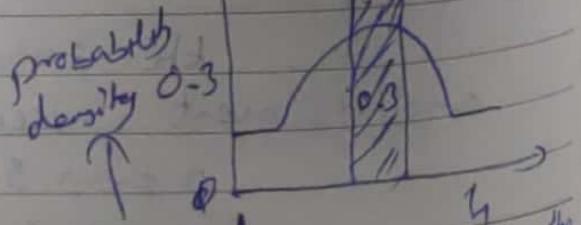
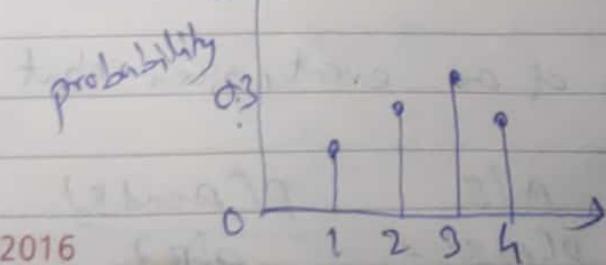
probability density function

probability mass function

(probability in % only)

p(2-3)

27 SUNDAY



probabilities are given by the surface area under the curve

S	M	T	W	T	F	S
1	2					
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30

APR - 2016

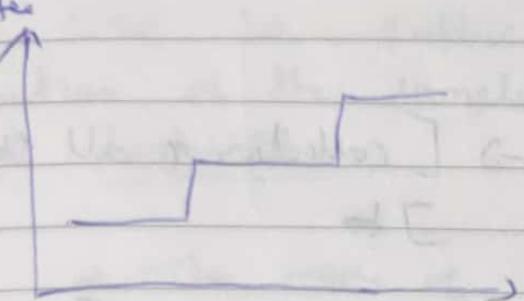
MARCH • MONDAY

28

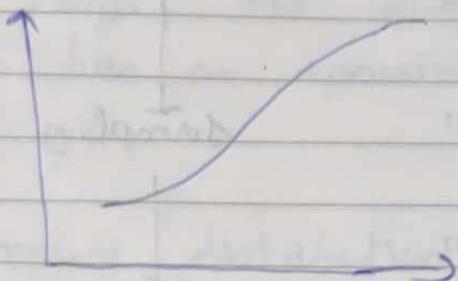
WK 14 • 088-278

Cumulative probability:

a) discrete



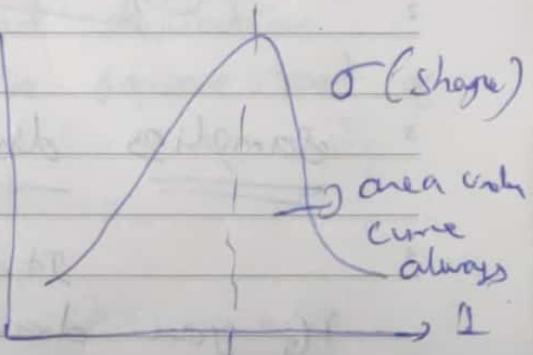
b) continuous



Normal distribution:

equation:

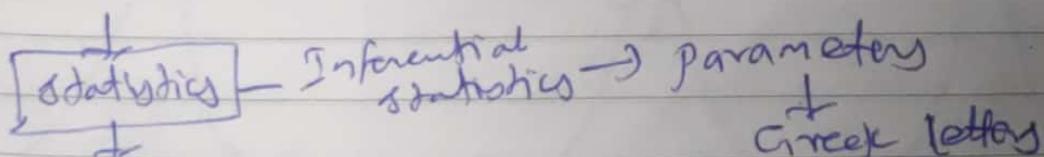
$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



populations? - All data

sample - subset of population (all data)

sample $\xrightarrow{\text{For conclusion}}$ populations



roman characters

\bar{x} = mean

s = standard deviation

μ = mean

σ = standard deviation
in population 2016

29

TUESDAY • MARCH

WK14 • 089-277

MAR - 2016

M	T	W	T	F	S	S
	1	2	3	4	5	
7	8	9	10	11	12	
14	15	16	17	18	19	
21	22	23	24	25	26	
28	29	30	31			

simple Random sample

9

10

11

12

1

2

3

4

5

6

7

population → all student



sampling frame → [collection of all student names]



sample data → 200 students [from the the above collection frame]



[Risqev is Better]

→ sample data from sampling frame

sampling distribution of sample mean:

It is the distribution % that you get if you draw an infinite number of samples from your population and compute the mean of all the collected sample means.

2016

S	M	T	W	T	F	S
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30

APR 2016

MARCH • WEDNESDAY

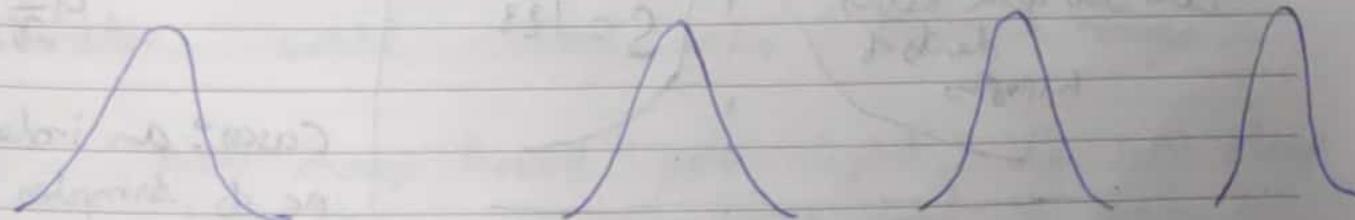
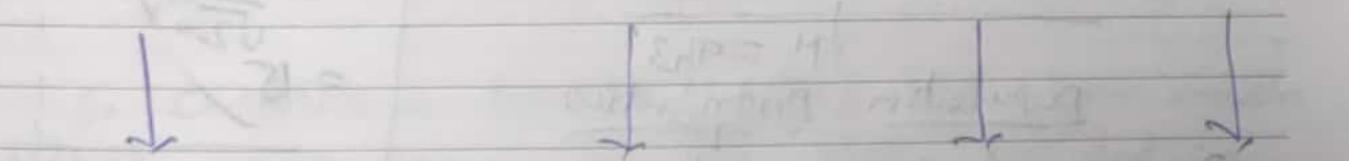
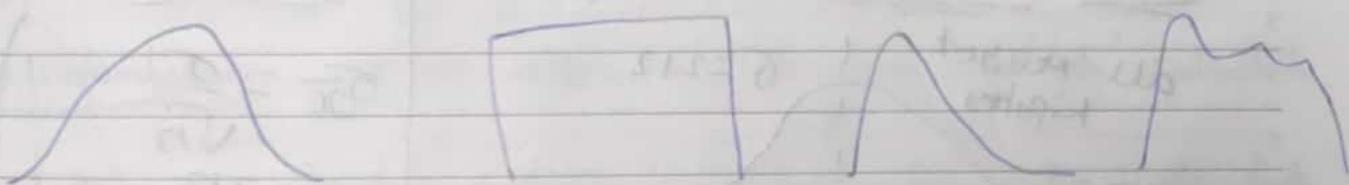
30

WK14 • 090-276

Central Limit theorem:-

- * It says that, provided that the sample size is sufficiently large, the sampling distribution of the sample mean has an approximately normal distribution.
- * The mean of the sampling distribution equals the population mean, and the standard deviation of the sampling distribution equals the standard deviation in the population divided by the square root of the sample size.

population distribution of mean



Sampling distribution of sample mean with
Sample size of $n=30$

2016

31

THURSDAY • MARCH

WK14 • 091-275

M	T	W	T	F
1	2	3	4	5
7	8	9	10	11
14	15	16	17	18
21	22	23	24	25
28	29	30	31	

MAR - 2016

$$\mu_{\bar{x}} = \mu \text{ (population mean)}$$

when \bar{x} is sample mean

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

where

σ = standard deviation of population

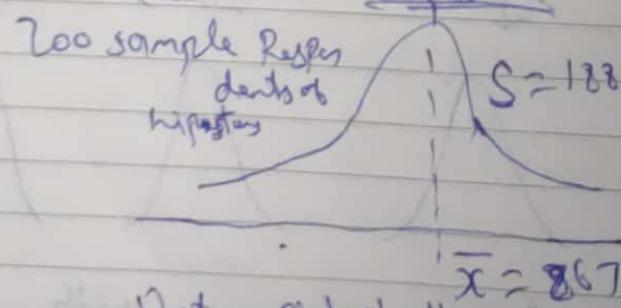
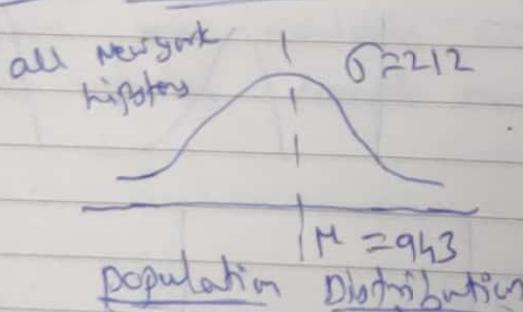
\bar{x} = sample mean

n = sample size

$\sigma_{\bar{x}}$ = standard deviation of sampling distribution

* The larger the variability in the population the larger the variability in the sample means

Three Distributions:



(Theoretical distribution) Sampling Distribution

$$\begin{aligned}\sigma_{\bar{x}} &= \frac{\sigma}{\sqrt{n}} \\ &= \frac{212}{\sqrt{200}} \\ &= 15\end{aligned}$$

$$\mu_{\bar{x}} = 943$$

Cases: an individual
random samples of
200 respondents from
the population of new
york hipsters

2016

S	M	T	W	T	F	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28

MAY - 2016

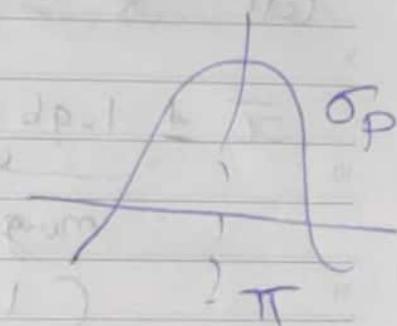
APRIL • FRIDAY

01

WK14 • 092-274

Sampling distribution proportion:

$$\sigma_p = \sqrt{\frac{\pi(1-\pi)}{n}}$$



where σ_p = standard deviation of sampling distributions of sample proportions
 p = sample proportions

Binary categorical variables

Compute proportions

$$\frac{1}{\pi}$$

Confidence Interval:

My understanding: our objective is we need to calculate population parameter which is done by sample data.

So confidence level is nothing but how much confidence we are to say that population parameter will fall into certain range.

Ex: how many hours per night do you sleep less than before you had a baby?

2016

02

SATURDAY • APRIL

WK14 • 093-273

APR-2016

M	T	W	T	F	S	S
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	

lets $\bar{x} = 2.6$, $\sigma = 1.1$ [value from 6o Readys]

$$\bar{x} \pm 1.96 \frac{\sigma}{\sqrt{n}} \text{ where } \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

margin of error

$$(1.96 \cdot 0.112 = 0.22)$$

(C) for mean
with the known
population std dev

$$2.6 \pm 0.22$$

\Rightarrow 95% Confidence Interval
(22.32, 2.88)

$$\bar{x} \pm Z_{95\%} \frac{\sigma}{\sqrt{n}}$$

where $Z_{95\%} \rightarrow$ 2 score
at 95% confidence

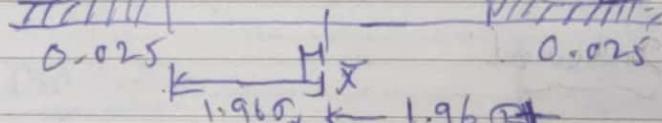
$$\bar{x} \pm 1.96 \frac{\sigma}{\sqrt{n}}$$

where

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

$$\mu_{\bar{x}} = \mu$$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$



margin of error

Margin of error

Margin of error

95% confidence interval

Confidence Interval

03 SUNDAY

(sample)

Compute the Confidence intervals

2016

in 95% of the samples the population value will fall within the confidence interval

M	T	W	T	F	S
3	4	5	6	7	
10	11	12	13	14	
17	18	19	20	21	
24	25	26	27	28	
31					

MAY - 2016

APRIL • MONDAY

04

WK15 • 095-271

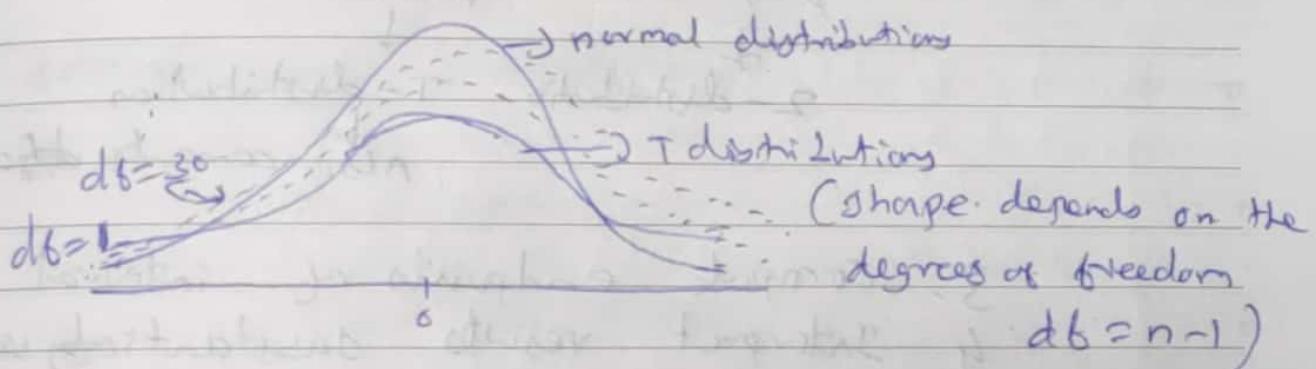
CI for mean with unknown standard deviation:

Assumptions:

1. randomization

2. approximately normal distribution

the wary of extreme outliers



eg:

how many hours per night do you sleep less than before you had a baby?

$$\bar{x} = 2.6 \quad s = 0.9$$

$$\bar{x} \pm t_{95\%} (se)$$

where $t_{95\%} =$ t score at 95% confidence

$$se = \frac{s}{\sqrt{n}}$$

1/2

$$CI \text{ for proportion: } p \pm Z_{95\%} (se)$$

$$\text{where } se = \sqrt{\frac{p(1-p)}{n}}$$

p = population mean

2016

M	T	W	T	F	S	S
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	

Constructing a confidence interval

9

- which confidence level?
- proportion (σ) mean?

11

proportion mean

+ +

z-distribution t-distribution

also compute $t_{df} = t_{df \text{ final}}$

12

- compute endpoints of interval

2 1 - α/2

- Interpret results substantively very

3

Choosing

the sample size

sample size

[Mean]



- magnitude of desired margin of error

(smaller margin of error → Larger sample size)

- confidence level

(larger confidence level → larger sample size)

- variability

(larger the std deviation → longer your sample size)

$$n = \frac{\sigma^2 z^2}{m^2}$$

where $\sigma \rightarrow \text{std dev of population}$ $z \rightarrow z \text{ score respective to confidence interval}$ $m \rightarrow \text{margin of error you want to maintain}$

S	M	T	W	F	S
1	2	3	4	5	6
8	9	10	11	12	13
15	16	17	18	19	20
22	23	24	25	26	27
29	30	31			

MAY - 2016

APRIL • WEDNESDAY

06

WK15 • 097-269

sample size [proportions]

$$n = \frac{p(1-p)z^2}{m^2}$$

sample size computations

ideal world

super large sample

real world → we don't have enough money
to compute super large
sample for regular fir

Sample size comput
tations

* significance Test:

e.g.

1. More than half of all certified divers
in America have more than 35 hours of
diving experience.

→ H_1 (dealing with proportions)

2. Mean no of hours of diving experience
of all certified divers in America
is more than 35?

→ H_1 (dealing with mean)

2016

07

THURSDAY • APRIL

WK15 • 098-268

APR - 2016

M	T	W	T	F	S
4	5	6	7	8	9
11	12	13	14	15	16
18	19	20	21	22	23
25	26	27	28	29	30

Example 1

How many hours of
experience

$p(> 35 \text{ hours}) = 0.57$

$n = 500$

Example 2

$\bar{x} = 35.5 \quad s = 8$

proportion \rightarrow (1) proportion or mean \leftarrow mean

$H_0: \pi = \pi_0 \geq 0.5 \quad (2) \text{ Formulate}$

$H_a: \pi \neq \pi_0$

hypothesis

$H_0: M = M_0 \geq 35$

$H_a: M \neq M_0$

 $\pi > \pi_0 \rightarrow$ Right Tailed Test $\leftarrow M > M_0$ $\pi < \pi_0 \rightarrow$ Left Tailed Test $\leftarrow M \neq M_0$

2 Randomization (3) check assumptions

 $n\pi \geq 15$ and $n(1-\pi) \geq 15$

$n(1-\pi) \geq 15$

Randomization
Population distribution
approximately normal
(if one tailed test
with small n)

$\alpha = 0.05$

(4) determine α

$\alpha = 0.05$

$Z = \frac{\hat{p} - \pi_0}{SE_0}$

(5) compute
statistic

$t = \frac{\bar{x} - M_0}{SE}$

$SE_0 = \sqrt{\frac{\pi_0(1-\pi_0)}{n}}$

$SE = \frac{s}{\sqrt{n}}$

$Z = (0.57 - 0.5)/$

$t = \frac{(35.5 - 35)}{8/\sqrt{500}}$

$\sqrt{\frac{0.5(0.5)}{500}}$

$2016 \quad Z = 3.13$

$t = 1.40$

S	M	T	W	T	F	S
5	6	7	8	9	10	11
12	13	14	15	16	17	18
19	20	21	22	23	24	25
26	27	28	29	30	31	

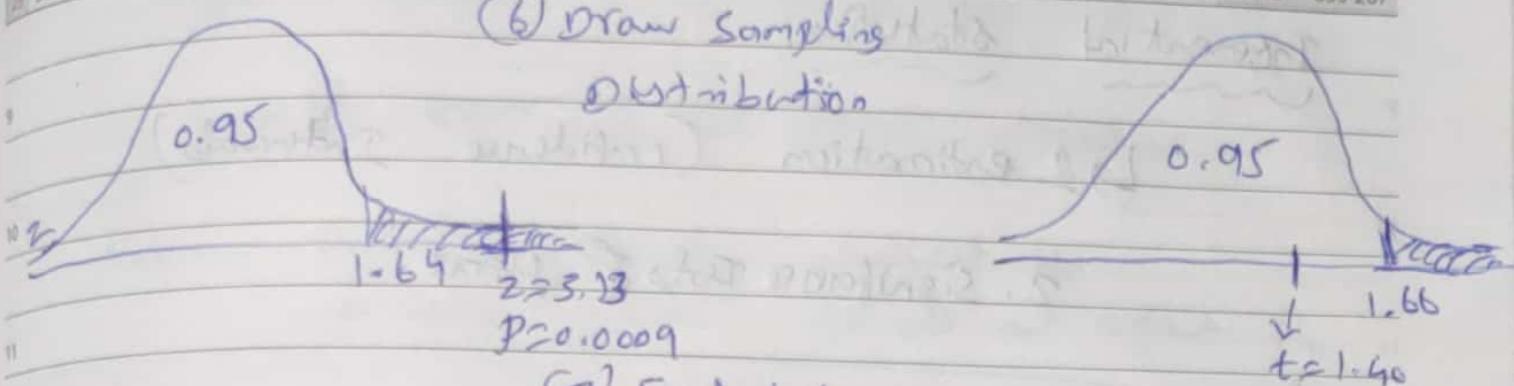
MAY - 2016

APRIL • FRIDAY

08

WK15 • 099-267

(6) Draw Sampling Distribution



(7) Find location of

test statistic in sub-area

(8) Reject H_0 ?

No

because t value

falls in Rejection Area

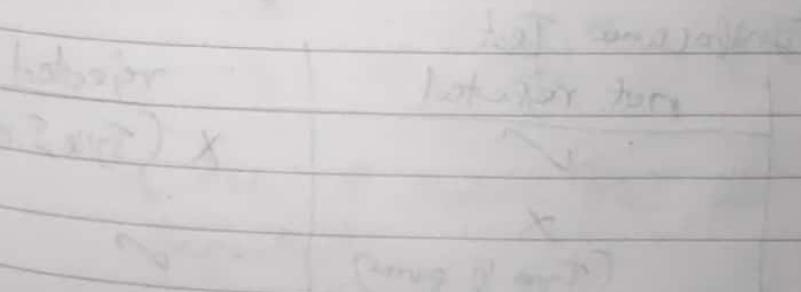
t value falls in Rejection Area

more than half of all certified drivers in America have more than 35 hours of daily experience

(9) Interpret findings

we cannot

conclude that there mean of no of driving exp is greater than 35 hours



2016

09

SATURDAY • APRIL

WK15 • 100-266

APR-2016

M	T	W	T	F	S	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30					

Inferential statistics

b) 1. estimation (confidence intervals)

2. Significance tests & point

P-value in a two tailed significance test ≤ 0.05 = 95% confidence interval does NOT contain H_0 value

()

P-value in a two tailed significance test > 0.05 = 95% confidence interval will contain the H_0 value

Type I and Type II error

defendant innocent
" guilty

	set free	convicted
defendant innocent	✓	✗
" guilty	✗	✓

10 SUNDAY

Significance Test	
H_0 True	not rejected
H_0 False	x (Type I error)

S	M	T	W	T	F	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28

MAY - 2016

APRIL • MONDAY

11

WK16 • 102-264

Type I error : \rightarrow False positive

Type II error : \rightarrow False negative

if $p\text{-value} \leq \alpha \rightarrow$ reject H_0

$p\text{-value} > \alpha \rightarrow$ do not reject H_0

chi-squared Test:

$$\chi^2 = \sum_i \frac{(\text{observed}_i - \text{expected}_i)^2}{\text{expected}_i}$$

how to choose degree of freedom for the distribution?

$$\text{D.o.F} = (\text{No. of Rows} - 1) * (\text{No. of columns} - 1)$$

	Fruit	Flowers	Mixed	
Early Renaissance	11	5	1	17
late "	8	6	8	22
Baroque	3	10	12	25
	22	21	21	64

$$\Rightarrow (3-1) * (3-1) = (2) * (2) = 4 \text{ D.o.F}$$

if the P value is < 0.05 the two variables dependent

\rightarrow dependent

\rightarrow independent

2016

S	M	T	F	S
1	2	3	4	5
9	10	11	12	13
16	17	18	19	20
23	24	25	26	27
30	31			

Jan
2019

APRIL • WEDNESDAY

13

WK16 • 104-262

Machine Learning: (Once written, they can keep on learning on its own and evolve as new data gets introduced with no human intervention required)

1642 - Mechanical Adder

1801 - First data storage of data

1847 - Boolean Logic

1890 - Mechanical system for statistical calculations

1950 - Turing Test

1952 - The perceptron

1967 - Pattern Recognition

1979 - Stanford car driving research at 20 km/h

1981 - Explanation Based Learning

1990's - Machine Learning Applications

2000's - Adaptive programming

Artificial Intelligence:

It is the study of how to train the computers so that computers can do things which at present humans can do better.

Machine Learning:

ML is said to learn from experience.

14

THURSDAY • APRIL

WK16 • 105-261

APR - 2016

M	T	W	T	F	S	S
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	

Artificial Intelligence

Machine Learning

- * ↑ chance of success not accuracy

↑ accuracy, and doesn't come about success

- * The goal is to simulate natural Intelligence to solve complex problem.

The goal is to learn from data on certain task to maximize the performance of machine on the task

- * AI is decision making

ML allows system to learn new things from data.

- * AI will go for finding the optimal solution

ML will go for only solution but that whether it is optimal or not.

- * AI leads to Intelligence or wisdom

ML leads to Knowledge.

Limitations of ML:

- * Performance
- * Computation

2016

S	M	T	W	T	F	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

Feb 9
MON 2016

APRIL • FRIDAY

15

WK16 • 106-260

Machine Learning process

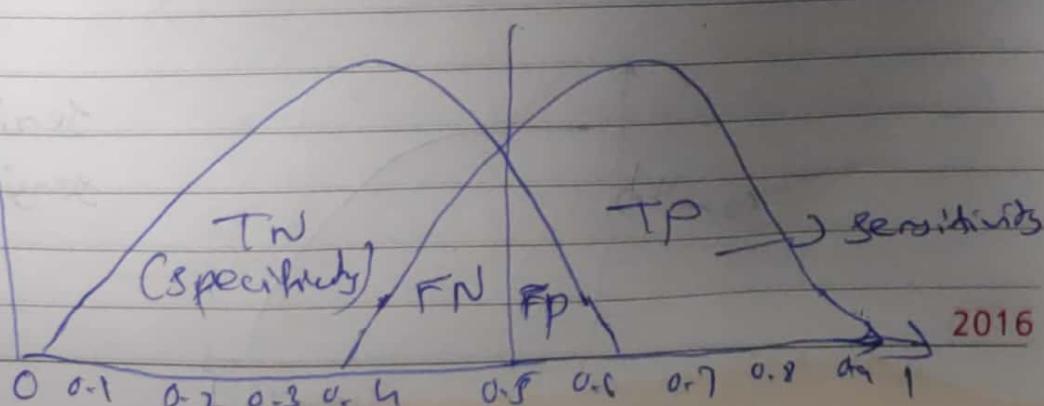
collecting data → pre processing and preparing data → Exploring data → choosing a model → Training a the model → evaluating the model → Improving the performance of model.

Confusion matrix:

It's a technique used for summarizing the performance of the classification algorithm.

e.g.: to predict cancer Yes/No

		predicted:		Actual:		Total
		No	Yes	Yes	No	
Actual:		TN = 50	FP = 10	TP = 100	FN = 5	165
NO						
Actual						
yes						
		55	110	110	165	
		Type II error	Type I error			



16

SATURDAY • APRIL

WK16 • 107-259

M	T	W	T	F	S
4	5	6	7	8	9
11	12	13	14	15	16
18	19	20	21	22	23
25	26	27	28	29	30

Pads APR-2016

sensitivity (or) Recall

- the proportion of patients that were identified correctly to have the disease.
- (i.e) True positive upon the total number of patients who actually have the disease is called sensitivity or Recall.

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

Actual
predicted
1 TP
0 FP
0 FN
1 TN

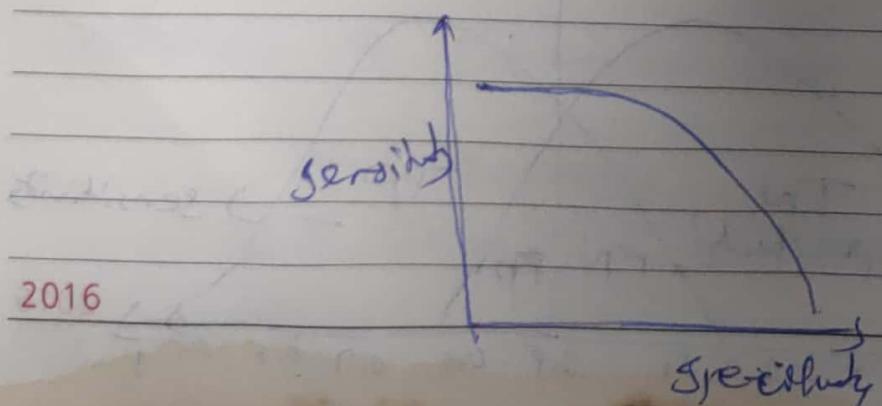
specificity

- the proportion of Patients that were identified correctly not to have the disease
- i.e) True negative upon the total number of patients who do not have the disease.

$$\text{Specificity} = \frac{TN}{TN + FP}$$

Actual
1 TP
0 FP
0 FN
1 TN

17 SUNDAY



sensitivity ↑ specificity ↓
sensitivity ↓ specificity ↑

2016

M	T	W	T	F	S
3	4	5	6	7	
10	11	12	13	14	
17	18	19	20	21	
24	25	26	27	28	
31					

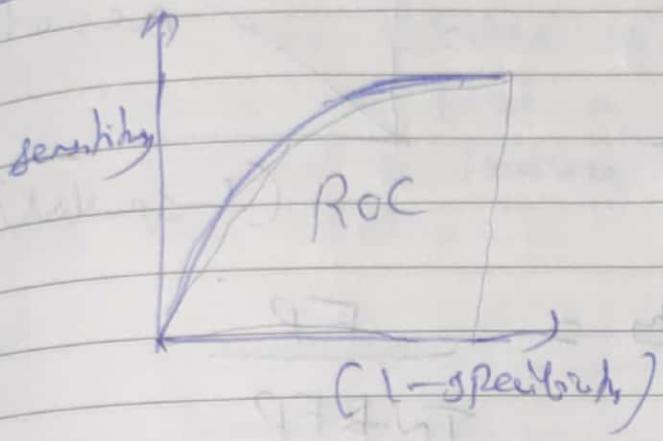
MAY - 2016

APRIL • MONDAY

WK 17 • 109-257

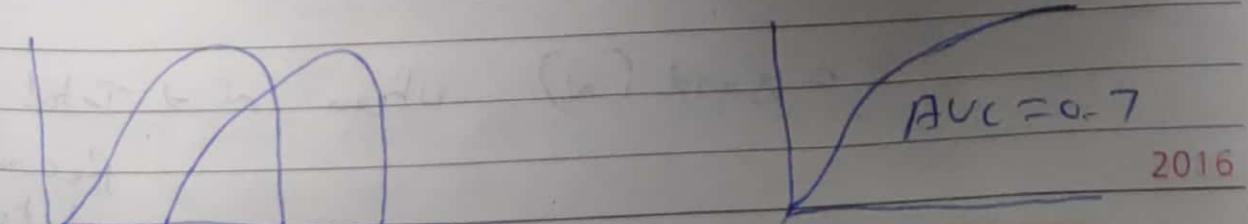
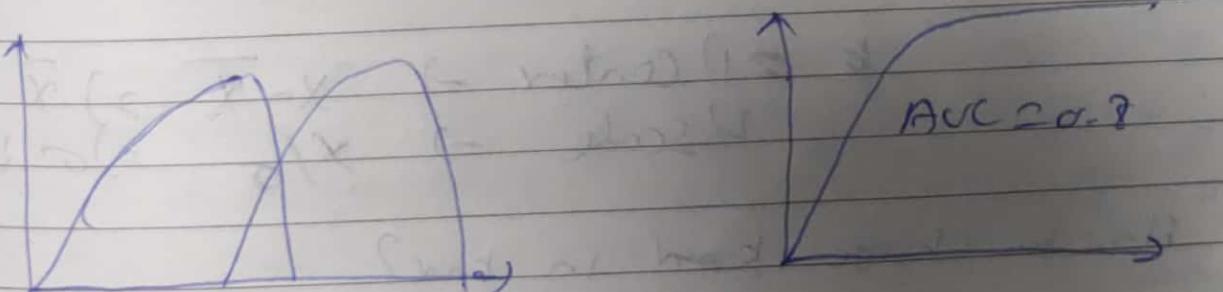
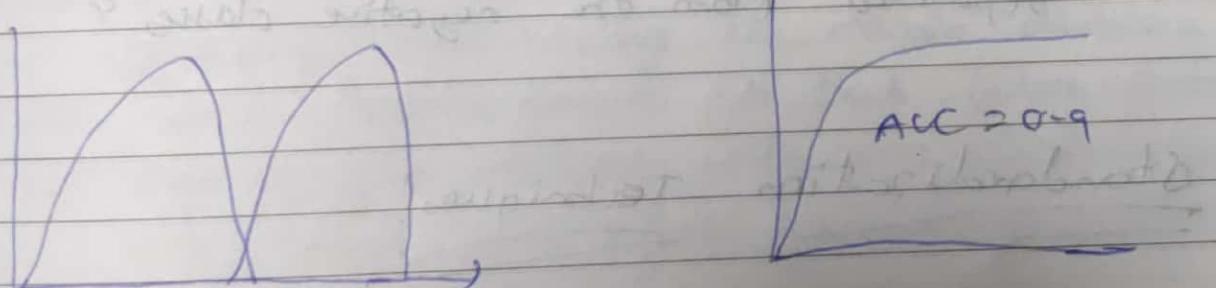
18

ROC curve: Receiver operating characteristic



AUC → Area Under the curve

The AUC is the Area under the ROC curve. This score gives us a good idea of how well the model performs.



2016

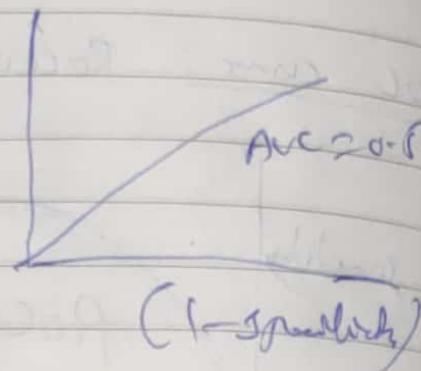
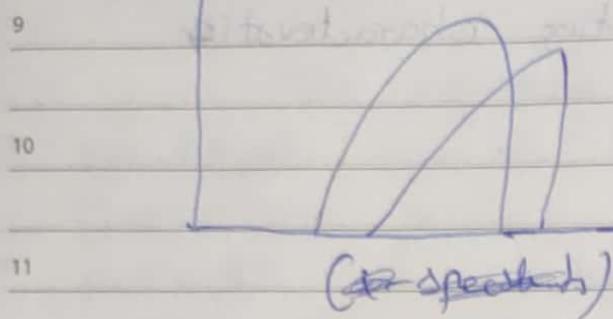
19

TUESDAY • APRIL

WK17 • 110-256

APR - 2016

M	T	W	T	F	S
1	2	3	4	5	6
7	8	9	10	11	12
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30



$$1 - \text{specificity} = \frac{FP}{TN + FP}$$

sensitivity \rightarrow True Positive Rate

(1 - specificity) \rightarrow False positive Rate

The AUC ROC score indicates how well the probabilities from the positive classes are separated from the negative classes.

Standardization Techniques:

$$k \approx 1) \text{Center} \rightarrow x - \bar{x} \quad 2) \hat{x} \text{ b mean}$$

$$4) \text{Scale} \rightarrow x/\sigma \quad 2) \sigma \text{ b std dev}$$

How to choose k and σ ?

$K = \sqrt{n}$ where $n \rightarrow \text{Total no of records in dataset}$

2016

S	M	T	W	T	F	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

Feb 9
MON 2016

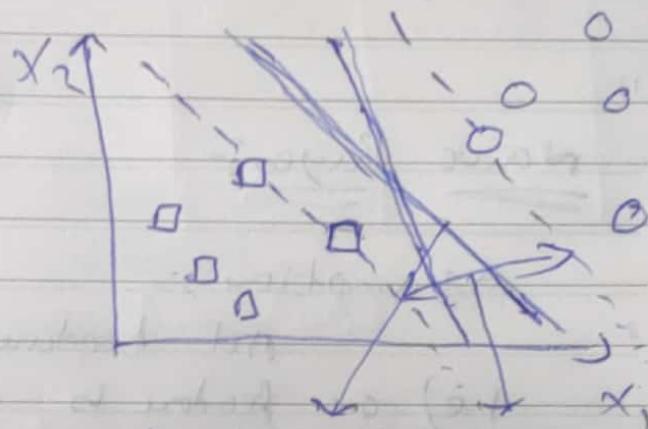
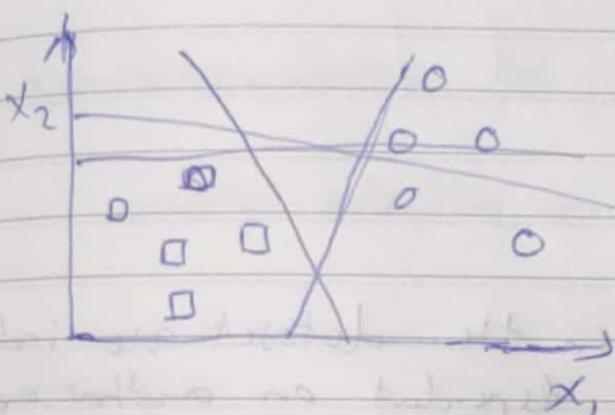
APRIL • WEDNESDAY

20

WK17 • 111-255

SVM - support vector machines:

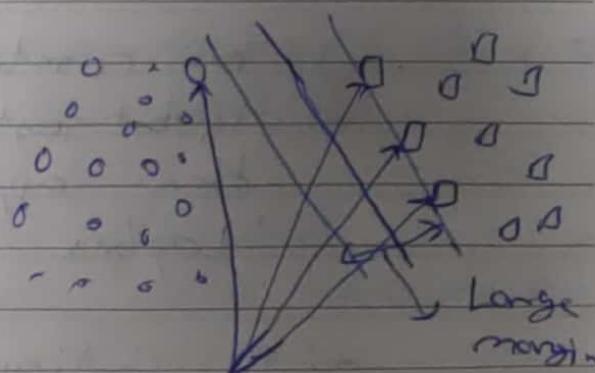
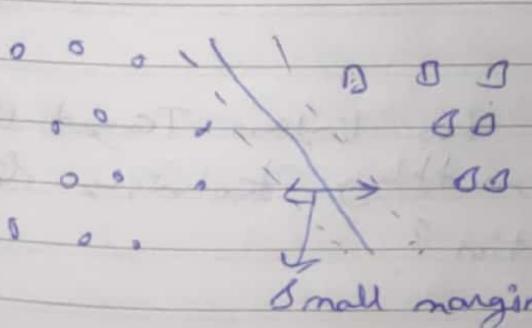
objective: The objective of support vector machine algorithm is to find a hyperplane in an N -dimensional space (N = # of features) that distinctly classifies the data points.



optimal hyperplane maximum margin

maximum margin \rightarrow maximum distance b/w data points of both classes.

By maximizing the margin distance provides some reinforcement so that future data points can be classified with more confidence.



Support vectors are data points that are closer to the hyperplane and influence the position and orientation of the hyperplane. 2016

21

THURSDAY • APRIL

WK17 • 112-254

APR-2016

M	T	W	T	F
4	5	6	7	1
11	12	13	14	2
18	19	20	21	3
25	26	27	28	29

- in Logistic Regression we take linear function output and squash the value within the range of $[0, 1]$ using the sigmoid function.
- in SVM we squash the linear function output values to $[-1, 1]$

Naive Bayes:-

Assumption:-

- All features in the dataset are independent i.e) one feature is not dependent on another one.

$$p(C|x) = \frac{p(x|C)p(C)}{p(x)}$$

↳ Likelihood ↳ class prior
 ↳ probability

↳ posterior probability ↳ Predicted prior probability

6. Tips to improve Naive Bayes Model:

- No correlation
- If test data has missing Target value then apply smoothing techniques like "Logistic Correction".
- If continuous variable do not have normal distribution, should we Transformation or diff methods to convert it in Normal distribution.

M	T	W	T	F	S
1	2	3	4	5	6
8	9	10	11	12	13
15	16	17	18	19	20
22	23	24	25	26	27
29	30	31			

MAY - 2016

APRIL • FRIDAY

22

WK17 • 113-253

application:

- mostly used in Text classification.
- spam filtering / sentiment analysis
- Recommendation system.
- multi-class predictions

Hyperparameter for NB:



Gaussian:-

It is used in classification and it assumes that features follow a normal distribution

multinomial :-

It is used for discrete counts.

e.g.: how often word occurs in the document

Bernoulli:-

If the feature vectors are binary (0's & 1's)

e.g.: Text classification with "word at

"Bag of words" model where

'1' → word occurs in the document

'0' → word does not occur in the document

2016

23

SATURDAY • APRIL

WK17 • 114-252

APR - 2016

M	T	W	T	F
4	5	6	7	8
11	12	13	14	15
18	19	20	21	22
25	26	27	28	29

NB example:

weather vs play

weather	play
Sunny	No
overcast	Yes
Rainy	Yes
Sunny	yes
:	:
n	in

weather	No	Yes
overcast	4	5
Rainy	3	2
Sunny	2	3
Total	5	9

Likelihood Table

overcast	4	$\frac{4}{14} = 0.29$
Rainy	3	$\frac{3}{14} = 0.36$
Sunny	2	$\frac{2}{14} = 0.36$
	$\frac{25}{14}$	$\frac{29}{14} = 0.64$
	$= 0.36$	$= 0.64$

problem: players will play if weather is sunny.
is the statement correct?

$$P(\text{yes} | \text{sunny}) = P(\text{sunny} | \text{yes}) * P(\text{yes})$$

24 SUNDAY

 $P(\text{sunny})$

$$\Rightarrow P(\text{sunny} | \text{yes}) = \frac{3}{9} = 0.33,$$

$$\begin{aligned} P(\text{sunny}) &= \frac{5}{14} = 0.36 \\ P(\text{yes}) &= \frac{9}{14} = 0.64 \end{aligned}$$

2016

S	M	T	W	T	F	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

MAY - 2016

APRIL • MONDAY

25

WK18 • 116-250

$$\begin{aligned}
 P(\text{yes})_{\text{sunny}} &= \frac{0.33(0.64)}{0.36} \\
 &= 0.60 \rightarrow \text{higher probability} \\
 &\text{So statement is correct.}
 \end{aligned}$$

2016