

09

TUESDAY
OCTOBER
2018

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

STATISTICS

09



10 Descriptive State

Inferential State

11

- ↳ Measure of central tendency
- ↳ Measure of dispersion

z - test

t - test

ANOVA test {F test}

CHI SQUARE.

12

"Summarising the data"

Hypothesis testing
[P value]

01

- Histogram
- Poly
- Cdf

Confidence interval

02

- Probability
- Permutation
- Mean, Median, Mode
- Standard deviation
- Variance

03

04

05

* Gaussian distribution

* Binomial "

06

* Bernoulli

* Parity "

* Standard Normal "

Notes

* Transformation & Standardization

* Q-Q plot

#

M	T	W	T	F	S
NOVEMBER 2018	1	2	3	4	5
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	



WEDNESDAY
OCTOBER
2018

10

Statistics - Collect → Organise → Analyse.

09

for better decision making

10

Descriptive Stats -

11

organising & summarizing data

Inferential Stats -

12

draw a conclusion.

01 * Population ^(N) and Sample ⁽ⁿ⁾

Sample → part of entire population.

02

Sampling Techniques -

03

① Simple Random Sampling

04

Every member has an equal chance of being selected.

05

② Stratified Sampling.

06

Population is split into non-overlapping groups

E.g. Male & Female

Notes

→ Based on age group.

③ Systematic Sampling

→ Pick up every n^{th} person

11

THURSDAY
OCTOBER
2018

S	M	T	W	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

Convenience

(4) Convenience Sampling

09

↳ only people having domain knowledge or interest.

10

11

Variable -

12

Quantitative

01

Age, ht.

02

Discrete

Continuous

03

o no of bank acc

o height

04

o no of children

(decimal v.)

o amt of rainfall

Qualitative / Categorical

Gender, Blood Grp
[Categories]

05

Variable Measurement Scale -

06

o Nominal → Categorical data

07

o Ordinal → order of data matters not value

08

o Interval → order & value matter, No 0.

09

o Ratio

Notes

frequency distribution -

Rose - 3

Lily - 4

⋮ ⋮

→ chart -



NOVEMBER	M	T	W	T	F	S	S
2018							
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			



FRIDAY
OCTOBER
2018

12

① Bar Graph → plot for discrete values.

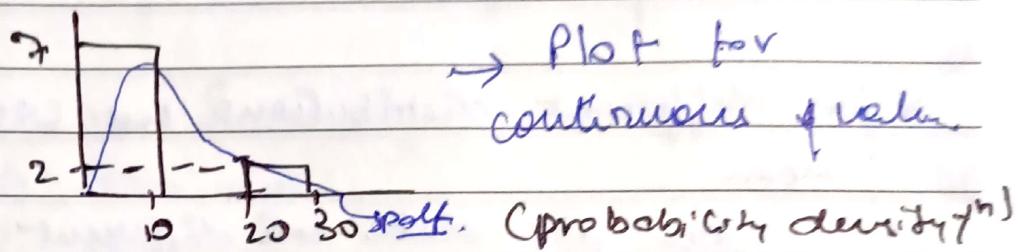
09

② Histogram →

10

$$B.C.n = 10$$

11



12

Pdf: Smoothening of histogram.

01

03 #Central Tendency. - measure used to determine the center of distribution

04

Mean - Average

Population - N

Sample - n

05

$$\mu = \sum_{i=1}^N \frac{x_i}{N}$$

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n}$$

06

Median - due to outliers mean changes a lot

Notes

- Median is considered

o Sort o find middle no long of

Mode - most frequent no.

(most with categorical variable)

Miss value can be replaced with most freq.

13

SATURDAY
OCTOBER
2018

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

Measure of dispersion -

09

o Variance - σ^2 o Standard Dev. - σ

10

↓
Spread of data.

11

* for different distributions we can have same mean -

12

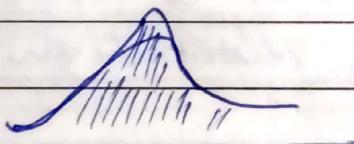
How are data sets different from each other - acc. of to their spread.

02

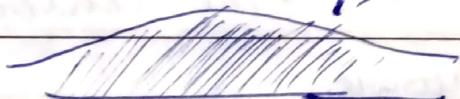
$$\sigma^2 = \frac{\sum_{i=1}^n (m_i - \mu)^2}{N}$$

$$\sigma^2 = \frac{\sum_{i=1}^n (m_i - \bar{x})^2}{n-1}$$

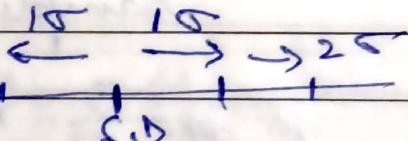
03



more variance,



04



05

V = Spread of data

S.D = range of spread

Unit to tell the distance.

* Percentile & Quartile.

SUNDAY 14

Percentile \Rightarrow $\frac{\text{no. of values below } x}{\text{total no. of values}} \times 100$

NOVEMBER	M	T	W	T	F	S	S	
2018					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			



MONDAY
OCTOBER
2018

15

What value exist for a particular percentile -

09

$$\text{Value} = \frac{\text{Percentile}}{100} \times (n+1)$$

10

↳ get index.

11

Five number Summary -

12

1. Minimum

2. Maximum

01

3. 1st Quartile

4. Median

02

5. 3rd Quartile

03

Removing Outliers -

[Lower fence \leftrightarrow Higher fence]

04

Lower fence - $Q_1 - 1.5(IQR)$

05

Higher fence - $Q_3 + 1.5(IQR)$

06

$IQR = \text{Inter Quartile Range} = Q_3 - Q_1$

$Q_3 \Rightarrow 75\%$. $Q_1 \Rightarrow 25\%$.

Notes

$\{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 5, 6, 6, 6, 7, 8, 8, 9, 27\}$

$$Q_1 \Rightarrow \frac{25}{100} \times (2019+1) \Rightarrow \frac{25 \times 20}{100} \Rightarrow 5 \Rightarrow 3$$

$$Q_3 \Rightarrow \frac{75}{100} \times 20 \Rightarrow 15 \Rightarrow 7 \Rightarrow 3$$

$IQR = 7 - 3 = 4$

$$L \Rightarrow 3 + 1.5 \times 4 = 7$$

$R \Rightarrow 7 + 1.5 \times 4 = 11$



Scanned with OKEN Scanner

16

TUESDAY
OCTOBER
2018

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

09

Minimum = 1

 $Q_1 = 3$

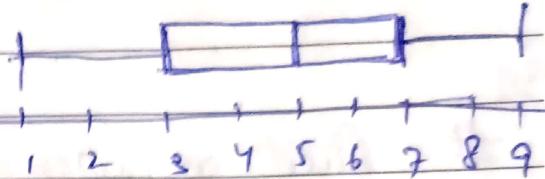
Median = 5

 $Q_3 = 7$

Mode = 9.

} Box Plot.

10



11

12

Distributions -

01

↳ Normal / Gaussian

02

↳ Standard Normal dis.

↳ Z-score

03

↳ Log normal distribution

↳ Bernoulli "

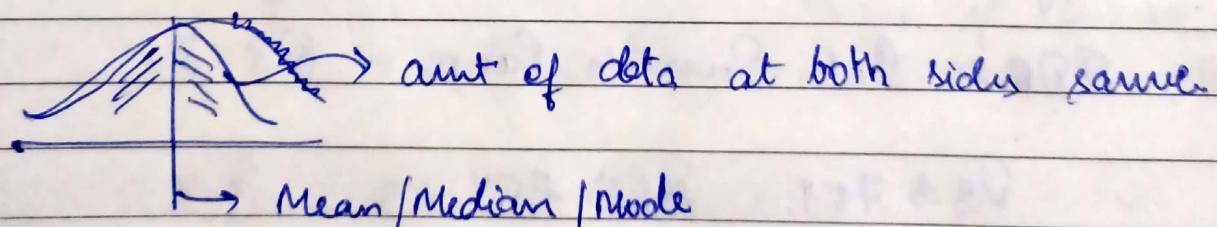
04

↳ Binomial "

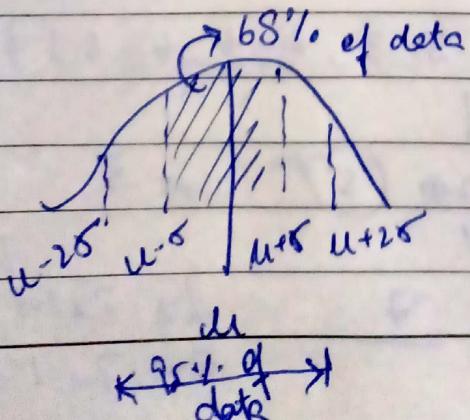
05

① Gaussian / Normal -

06



Notes



Empirical formula -

$$68 - 95 - 99.7$$



NOVEMBER	M	T	W	T	F	S	S	
2018					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			



WEDNESDAY
OCTOBER
2018

17

Example - Height - Normally distributed

09 wt - "

IRIS - "

10

Z score - how much standard deviation is it away from the mean

11 12 Z score $\Rightarrow \frac{x_i - \mu}{\sigma}$

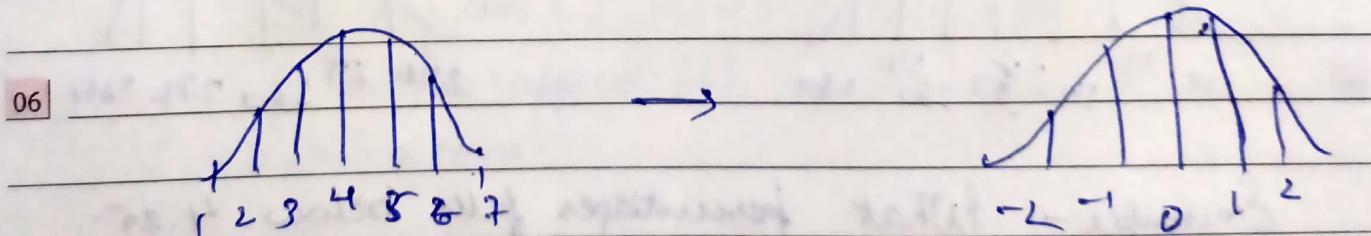
01 Say $\sigma = 1$ $\mu = 4$

02 for $x = 4.75$ $\frac{4.75 - 4}{1} \Rightarrow 0.75$

03 +ve \rightarrow Right direct \quad -ve \rightarrow left direct.

04 Standard Normal distribution -

On applying Zscore to all values of sample the distribⁿ obtained is called SND



Notes

$$\{1, 2, 3, 4, 5, 6, 7\} \rightarrow \{-3, -2, -1, 0, 1, 2, 3\}$$

$$\mu = 0$$

$$\sigma = 1$$

Standardization - Conversion into standard normal distib.

18

THURSDAY
OCTOBER
2018

SEPTEMBER	M	T	W	T	F	S	S
2018							
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	

Normalizatⁿ → convert to (0 to 1) data set.

09

Min Max Scalar

10

 $\rightarrow (0, 1)$.

11 Example -

2021

2022

12 Series avg 250 260
S.D. 10 12

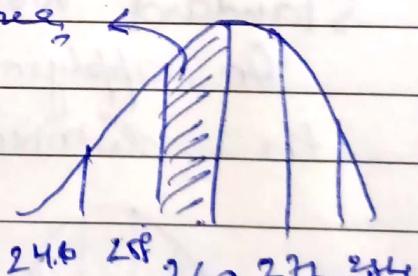
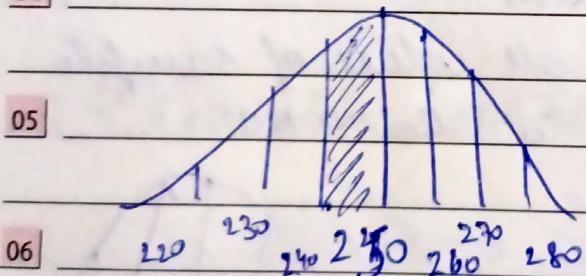
01 Peony score 240 245

02 which score is better

03 Z-scores $\frac{240 - 250}{10} \rightarrow -1$

$$\frac{245 - 260}{12} \rightarrow -1.25$$

more area,



Example - What percentage fall below 4.20.

Notes

* find Z-scores $\frac{x - \mu}{\sigma}$

Look at Z-table to find the area

= % of score fell above/below



Scanned with OKEN Scanner

NOVEMBER	M	T	W	T	F	S	S	
2018					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			



FRIDAY
OCTOBER
2018

19

Example - $\text{Avg IQ} = 100$
 09 $\sigma = 15$

10 lower than 85?

10 Z-score $\rightarrow \frac{85-100}{15} \rightarrow -\frac{15}{15} \rightarrow -1$
 11 \rightarrow look into table,

12 Removing Outliers using python

01 outliers = []

02 def detect_outliers(data):

threshold = 3

03 mean = np.mean(data)

std = np.std(data)

04 for i in data:

Zscore = (i - mean) / std

if np.abs(z-score) > threshold:

outlier.append(i)

return outliers.

Notes $IQR = Q_3 - Q_1$

- ① sort the data
- ② calculate Q_1 & Q_3
- ③ $IQR = Q_3 - Q_1$
- ④ find lower fence
- ⑤ find upper fence

20

SATURDAY
OCTOBER
2018

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

dataset = sorted(dataset)

09

$$\theta_1, \theta_3 = \text{mp. percentile}(\text{dataset}, [25, 75])$$

10

$$IQR = \theta_3 - \theta_1$$

$$\text{lower fence} = q_1 - (1.5 * IQR)$$

11

$$\text{higher fence} = q_3 + (1.5 * IQR)$$

12

Probability -
measure of likelihood.

01

Addition Rule -

✓ Mutual Exclusive - cannot occur at same time

02

$$P(A \text{ or } B) = P(A) + P(B) = 1$$

03

✓ Non - mutual Exclusive

04

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

05

Multiplication Rule -

✓ Independent Event -

one event not dependent on another event

$$P(AB) \text{ AND } P(B) = P(A) * P(B)$$

✓ Dependent Events -

SUNDAY 21 $E \rightarrow$ 3 Red marble 2 white marble

$$R \rightarrow \frac{3}{5}, \frac{2}{4}$$

Name Bayes

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

NOVEMBER	M	T	W	T	F	S	S
2018							
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			



MONDAY
OCTOBER
2018

22

At Permutation -

$${}^m P_r \Rightarrow \frac{m!}{(m-r)!}$$

✓

09

10

Combination -

unique w.r.t. elements (no swaps allowed)

11

$${}^m C_r = \frac{m!}{(m-r)! r!}$$

12

P value :

$P = 0.8$ out of 100 ; 80 times will be this.

02

03 Influential Statistics -

04 \rightarrow Test whether a coin is fair or not by performing test for 100 times.

05

if 50 times H comes then \rightarrow fair.

06

Hypothesis testing -

(1) Null Hypothesis \rightarrow coin is fair

Notes (2) Alternate " \rightarrow coin is unfair

(3) Experiment.

(4) Reject / Accept null hypothesis.

23

TUESDAY
OCTOBER
2018

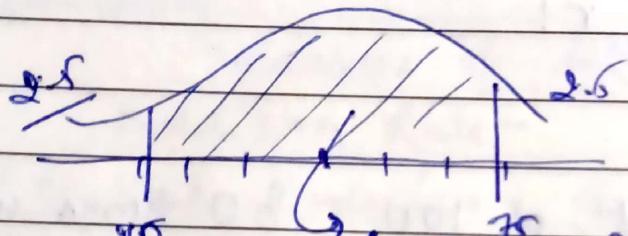
S	M	T	W	T	F	S	5
						1	2
3	4	5	6	7	8	9	
10	11	12	13	14	15	16	
17	18	19	20	21	22	23	
24	25	26	27	28	29	30	

Say \rightarrow 30 times heads.

09 "define how far can the value be away from the mean"

10 \hookrightarrow Significant Value - $\alpha = 0.05$ (5%)

11 $100 - \alpha = 95\% \equiv$ Confidence Interval.



defined by domain
exp. ~

02

03

04

05

06

Notes