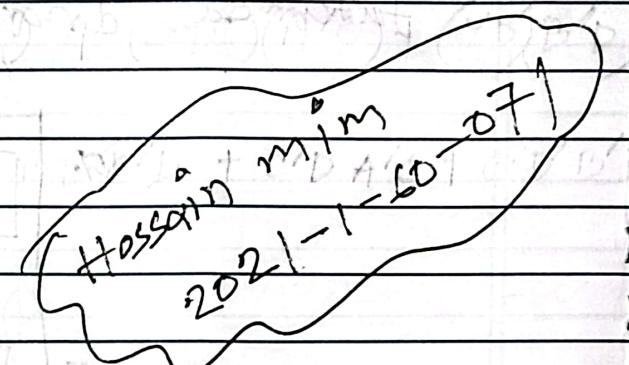


MIDOL

1) • central tendency

- trimmed mean
- mode
- mean absolute deviation
- standard deviation
- variance



2) Box Plot — 2 or \leftrightarrow SD

3) • STANDARD ERROR (Example) Explain

- BOOTSTRAP METHOD

- Example of Bootstrap

- Algorithm of Bootstrap method.

4) Hypothesis testing



Chi-square test

CSE 303

Dmim sir

MIDOL

Dmim sir

Date:

Page:

Hossain Mim
1-2021-1-60-071

Q-P1

60, 70, 80, 75, 65, 70, 80, 70, 65, 65

(Ans) sorted Order:— 60, 65, 65, 65, 70, 70, 70, 75, 80, 80

Here,

$n = 10$

$$\text{mean}(\bar{x}) = \frac{\sum x_i}{n}$$

$$= \frac{700}{10}$$

$$= 70 \quad (\text{Ans})$$

$n = 10$

$$\text{median} : \left(\frac{n}{2} \right)^{\text{th}} \text{ value} + \left(\frac{n}{2} + 1 \right)^{\text{th}} \text{ value}$$

$$= \frac{5^{\text{th}} \text{ value} + 6^{\text{th}} \text{ value}}{2}$$

$$= \frac{70 + 70}{2}$$

$$= 70 \quad (\text{Ans})$$

Q-P2

60, 70, 80, 75, 65, 70, 100, 70, 65, 65

$$\text{mean}(\bar{x}) : \frac{720}{10}$$

$$= 72 \quad (\text{Ans})$$

For Trimmed mean:—

sorted the dataset → 60, 65, 65, 65, 70, 70, 70, 75, 80, 100

For $P = 1$ so, 60 and 100 will be discarded.

$$\text{mean}(\bar{x}) = \frac{\sum x_{i+1}^n}{n-p}$$

$$= \frac{560}{10-2}$$

$$= \frac{560}{8}$$

$$= 70 \quad (\text{Ans})$$

Weighted mean:— $x_i = 40, 45, 80, 75$ and 10 have

$w_i = 1, 2, 3, 4$ and 5

$$\text{(}\bar{x}_{\text{w}}\text{)} = \frac{\sum_{i=1}^n w_i x_i}{\sum w_i}$$

$$= \frac{(40 \times 1) + (45 \times 2) + (80 \times 3) + (75 \times 4) + (10 \times 5)}{1+2+3+4+5}$$

$$= \frac{720}{15} = 48 \quad (\text{Ans})$$

mean Absolute deviation:—

$$\text{MAD} = \frac{\sum |x_i - \bar{x}|}{n}$$

dataset 1: 30, 40, 50, 60, 70

$$\bar{x} = 50$$

$$\text{MAD:— } |(30-50)| + |(40-50)| + |(50-50)| + |(60-50)| + |(70-50)|$$

$$= \frac{20+10+0+10+20}{5}$$

$$= \frac{60}{5}$$

$$= 12 \quad (\text{Ans})$$

Median Absolute deviation:

$$MAD = \text{median}(|x_1 - m|, |x_2 - m|, \dots, |x_N - m|)$$

dataset :- (2) 0, 25, 50, 75, 100

n = 5 (odd number)

$$\text{Ans} \Rightarrow \text{MAD} = \text{median}(|10 - 50|, |25 - 50|, |50 - 50|, |75 - 50|, |100 - 50|)$$

$$= \text{median}(40, 25, 0, 25, 50)$$

$$= \left(\frac{5+1}{2}\right) \left(\frac{6}{2}\right) 3\text{rd value}$$

$$= \text{median}(0, 25, 25, 50, 50)$$

~~MAD
Ans~~

$$n = 5 (\text{odd})$$

$$= \left(\frac{5+1}{2}\right) \text{th value}$$

$$= \left(\frac{5+1}{2}\right)^{\text{th}}$$

$$= 3\text{rd value}$$

$$= 25 \quad (\text{Ans})$$

Variance :- (6)

$$\sigma^2 = \frac{1}{n} \sum |x_i - \bar{x}|^2$$

Example:- The Height : 600 mm, 470 mm, 170 mm, 430 mm, 300 mm

$$\text{mean } (\bar{x}) : \frac{1970}{5} = 394$$

$$\text{variance } (\sigma^2) : \frac{(600 - 394)^2 + (470 - 394)^2 + (170 - 394)^2 + (430 - 394)^2 + (300 - 394)^2}{5}$$

$$= 21704 \quad (\text{Ans})$$

$$\text{standard deviation } (\sigma) = \sqrt{21704}$$

$$= 147.32 \quad (\text{Ans})$$

sample variance (s^2) :-

dataset 1 :- 30, 40, 50, 60, 70

$$n = 5$$

$$\bar{x} = \frac{250}{5} = 50$$

$$s^2 = \frac{(30-50)^2 + (40-50)^2 + (50-50)^2 + (60-50)^2 + (70-50)^2}{5-1}$$

$$= \frac{1000}{4}$$

$$= 250 \quad (\text{Ans})$$

dataset 2 :- 0, 25, 50, 75, 100

$$\bar{x} = \frac{25+100}{2} = 62.5$$

$$= 62.5$$

$$s^2 = \frac{(0-62.5)^2 + (25-62.5)^2 + (50-62.5)^2 + (75-62.5)^2 + (100-62.5)^2}{5-1}$$

$$= \frac{6250}{4}$$

$$= 1562.5$$

standard deviation : $s = \sqrt{s^2}$

$$= \sqrt{250}$$

$$= 15.82 \quad (\text{Ans})$$

standard deviation :-

$$s = \sqrt{1562.5}$$

$$= 39.5m$$

(Ans)

Practice Examples

dataset: — ৮, ৯, ১০, ১০, (১০), ১১, ১১, ১১, ১২, ১৩

mean (\bar{x}): $\frac{10.5}{10}$

= 10.5 (Ans)

 $n = 10$ (Even number)

median = $\frac{(\text{2nd value} + \text{3rd value})}{2}$
= $\frac{10 + 11}{2}$

= 10.5 (Ans)

mode: (10, 11) Bi-modal dataset (Ans)

Range: (max - min)
= (13 - 8)

= 5 (Ans)

$S^2 = \frac{(8-10.5)^2 + (9-10.5)^2 + (10-10.5)^2 + (10-10.5)^2 + (10-10.5)^2 + (11-10.5)^2 + (11-10.5)^2 + (12-10.5)^2 + (12-10.5)^2 + (13-10.5)^2}{10}$

= 0.85 (Ans)

$s^2 = \sqrt{1.85}$

= 1.36 (Ans)

MID ৬।
SMC ২

sample variance

$S^2 = \frac{(8-10.5)^2 + (9-10.5)^2 + (10-10.5)^2 + (10-10.5)^2 + (10-10.5)^2 + (11-10.5)^2 + (11-10.5)^2 + (12-10.5)^2 + (12-10.5)^2 + (13-10.5)^2}{(10-1)}$

= 2.055

$s = \sqrt{2.055}$

= 1.43 (Ans)

mean Absolute Deviation: — MAD: $(|8-10.5| + |9-10.5| + |10-10.5| + |10-10.5| + |10-10.5|)$

$|11-10.5| + |11-10.5| + |11-10.5| + |12-10.5| + |13-10.5|$

$= \frac{5+3+6+11}{10} = \frac{11}{10} = 1.1$ (Ans)

median Absolute Deviation: —

MAD: median $(|8-10.5|, |9-10.5|, |10-10.5|, |10-10.5|, |10-10.5|, |11-10.5|)$ $|11-10.5|, |11-10.5|, |12-10.5|, |13-10.5|$ $n=10$
 $\frac{n}{2} = 5$
 $\frac{n+1}{2} = 5.5$
 0.5×0.5
0.5 (Ans)

= median { 2.5, 1.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 1.5, 2.5 }

= median { 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 1.5, 1.5, 2.5, 2.5 }

Data: 0, 5.6, 8.7, 14.1, 14.1, 15, 17.2, 19.2, 19.3, 24.1, 24.7

Max: 24.7
Min: 0

$$IQR: Q_3 - Q_1 \Rightarrow (19.225 - 11.3) \Rightarrow 11.3$$

Date:

Page:

$$Q_1 = 25\% \times 11$$

= 2.75 + value

$$= 5.6 + (8.7 - 5.6) \times 0.25$$

$$= 7.925$$

$$Q_3 = 75\% \times 11$$

= 8.25 + value

$$= 19.2 + (19.3 - 19.2) \times 0.25$$

$$= 19.225$$

$$Q_2 = 50\% \times 11$$

= 5.5 + value

$$= 14.1 + (15 - 14.1) \times 0.5$$

$$= 14.55$$

$$\text{Upper Extreme: } Q_3 + IQR \times 1.5$$

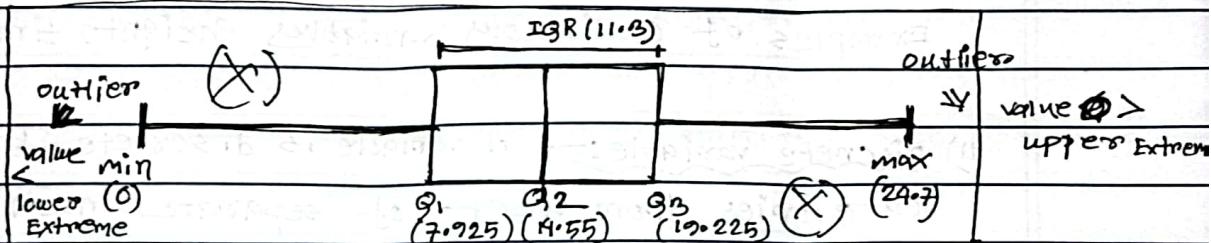
$$= 19.225 + 11.3 \times 1.5$$

$$= 36.175$$

$$\text{Lower Extreme: } Q_1 - IQR \times 1.5$$

$$= 7.925 - 11.3 \times 1.5$$

$$= -9.025$$

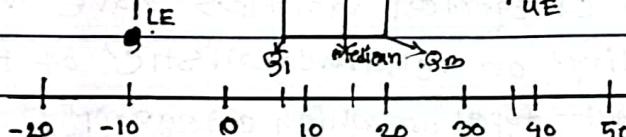


Box Plot

IQR

Hossain Mim

2021-1-60-071



Outlier

5.6, 24.7

LE

Q₁ Median Q₃

UE

50

40

30

20

10

0

-10

-20

$$Q_1 = 25\% \times 10$$

= 2.5 + value

$$= 10 + (15 - 10) \times 0.5$$

$$= 12.5$$

$$UE = (Q_3 + 1.5 \times IQR)$$

$$= 7.5$$

$$LE = (Q_1 - 1.5 \times IQR)$$

$$= -2.5$$

$$Q_2 = 50\% \times 10$$

= 5 + value

$$= 2.5$$

$$IQR = (37.5 - 12.5)$$

$$= 25$$

$$Q_3 = 75\% \times 10$$

= 7.5 + value

$$= 35 + (40 - 35) \times 0.5$$

variable: A variable is any characteristic, number or quantity that can be measured. A variable may also be called a data item.

Age, sex, eye colour, country of birth etc are examples of variable.

Numerical variables:— Numerical variable have values that describe a measurable quantity as a number, like 'how many' or 'How much'.

i) continuous variable:— A continuous variable is a numeric variable. Observations can take any value between a set certain set of real numbers.

Examples of continuous variables height, time, age etc.

ii) discrete variable:— A variable is discrete if it possible categories from a set of separate numbers.

Examples of discrete variables, numbers of children in a family, Number of registered cars, Number of boy in the school university.

categorical variables:— Categorical variables have values that describe a 'quality' or 'characteristic' of the data unit like 'what type' or 'which category'.

i) Nominal variables:— A nominal variable is made up of various categories which has no order.

Example of nominal variables include sex, eye colour, religion etc.

ii) ordinal variable:— An ordinal variable is a categorical variable. Observation can take a value that can be logically ordered or rank.

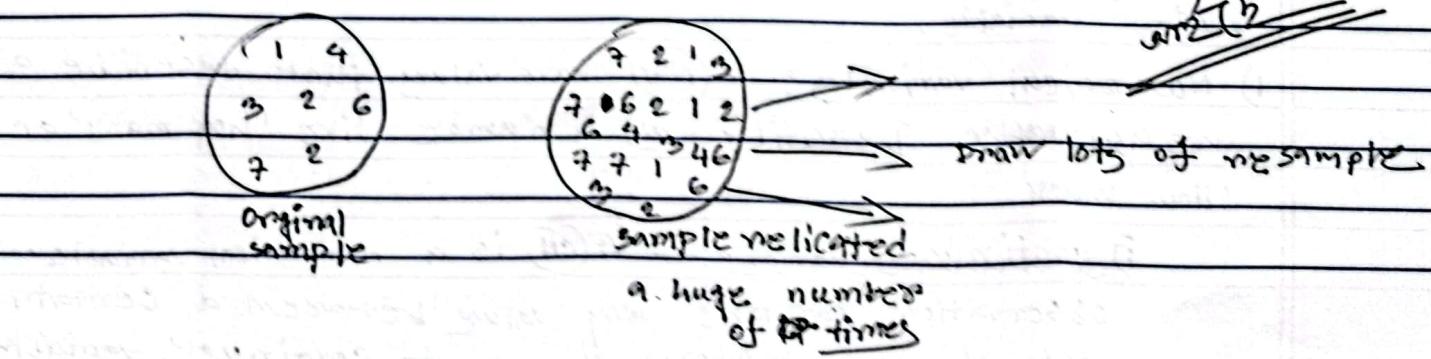
Examples of ordinal variables include

High \rightarrow minimum
 $1021-1-60-071$

There are mainly seven steps involved into it :-

- i) Decide on the objectives :— These objectives may usually require significant data collection and analysis.
- ii) what to measure and how to measure :— measurement generally refers to the assigning of the numbers to indicate different values of variables
- iii) Data collection :— once you know what type of data you need for your statistical study then you can determine whether your data can be gathered.
- iv) Data cleaning :— This is another crucial step in data analysis pipeline is to improved data quality for your existing data.
- v) Summarizing :— Exploratory data analysis helps to understand the data better.
- vi) Data modeling :—
- vii) Optimize and Repeat :— The data analysis is a repeatable process and sometime tends to continuous improvements , both to the business .

The Bootstrap:— one easy and effective way to estimate the sampling distribution of a statistic of a model parameters , is to draw additional samples , with replacement from the sample itself and recalculate the statistic or model each resample . This procedure is called bootstrap and it does not necessarily involve any assumption about the data or sample statistic ~~being~~ being normally distributed .



~~বোর্ড কলেজ এবং মাসিক~~
~~বোর্ড কলেজ এবং মাসিক~~

The algorithm for a bootstrap sampling of the mean
is as follows for a sample of size n :-

1) Draw a sample value, record, replace it

2) Repeat n times

3) Record the mean of the n resampled values

4) Repeat steps (1-3) R times

5) Use the R results to:-

i) calculate their standard deviation

ii) produce a histogram or boxplot

iii) Find a confidence interval.

~~সম্পর্ক~~
Confidence Interval:- A CI in statistics is a range of values that are used to estimate an unknown population parameter, such as the mean or proportion. It provides a range of plausible values for the parameter and is calculated from sample data.

For example 95% CI for mean ~~weight~~ of a population might be from 60 kg to 70 kg, indicating that we're 95% confident that the true population mean weight falls within this range.

The width of the interval is influenced by the sample size and the desired level of confidence.

variable:- A variable is any characteristic, number, quantity that can be measured. A variable also be called a data item. ~~for~~ Age, sex, eye colour, country of birth etc are examples of variable.

1) Numerical variables:- (Nv) have values that describe a measurable quantity as a number like 'how many' or 'How much'

i) continuous variable:- A (CN) is a numerical variable.

Observation can take any value between a certain set of real numbers. Examples of continuous variable include weight, time, age.

A confidence interval is a range of values that is likely to contain a population parameter, such as the mean or proportion.

It provides a range of plausible values for the parameters and is calculated from sample data.

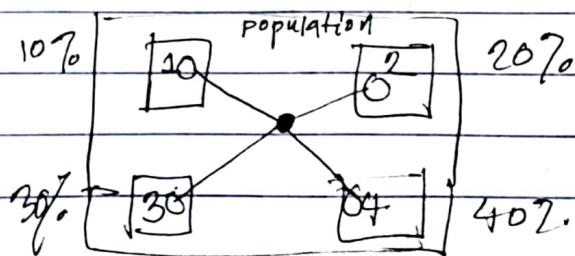
95% ⇒ It means that if we repeat an experiment or survey over and over again, 95% of the time our results will match the result we got from a ~~pop~~ population.

self-selection bias in ~~statistic~~ occurs when individuals or entities self-select into a sample being studied, leading to a non-random or biased representation of the population.

An Example self-selection bias in political polling. Consider a situation where a poll about a certain candidate on online platforms. People who actively seek out and participate in such poll might have stronger opinions or be more engaged in the topic compared to general population.

stratified sampling is a sampling technique in statistics where the population is divided into distinct subgroups, called strata, based on certain characteristics that are relevant to the study.

The strata are homogeneous within themselves but different from each other in some key aspects.



sorted data:— 64, 65, 68, 69, 70, 71, 72, 72, 75, 79, 80, 81, 83, 85

$$25^{\text{th}} \text{ percentile } (Q_1) = 25\% \times 14 \\ = 3.5 \text{ th value} \\ = 68 + (69 - 68) \times 0.5 \\ = 68.5$$

~~MID
Q1 68.5
Q2 72~~

$$50^{\text{th}} \text{ percentile (median)} = 50\% \times 14 \\ = 7^{\text{th}} \text{ value} \\ = 72$$

$$75^{\text{th}} \text{ percentile } (Q_3) = 75\% \times 14 \\ = 10.5 \text{ th value} \\ = 75 + (80 - 75) \times 0.5 \\ = 77.5$$

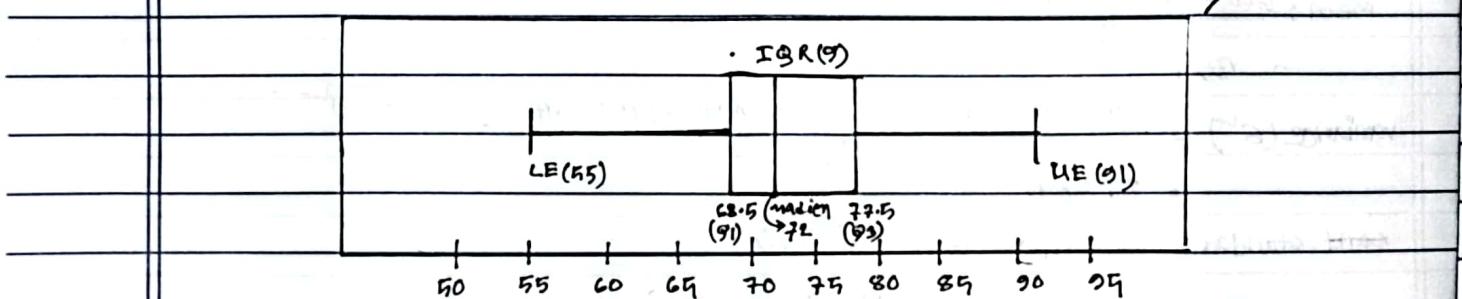
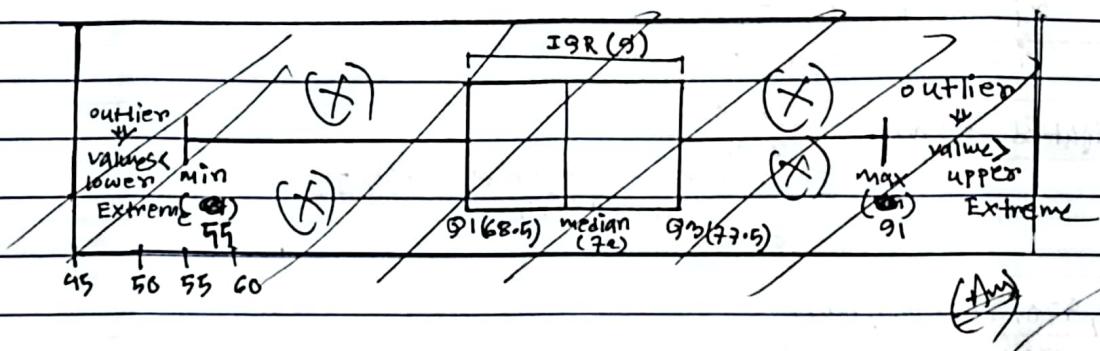
$$\text{IQR} = (77.5 - 68.5) \\ = 9$$

max = 85

min = 64

$$\text{upper Extreme} = \{ Q_3 + (IQR \times 1.5) \} \\ = \{ 77.5 + (9 \times 1.5) \} \\ = 91$$

$$\text{Lower Extreme} = Q_1 - IQR \times 1.5 \\ = \{ 68.5 - (9 \times 1.5) \} \\ = 59$$



we can see

no outliers

~~Ans~~

H_0 : There is no link between gender and political party preference

Date _____

H_1 : There is link between gender and political party preference

Page _____

$$\text{Expected value for male republican} = \frac{240 \times 200}{420}$$

$$= 114.28$$

$$\text{Expected value for female republican} = \frac{1240 \times 220}{420}$$

$$= 125.72$$

$$\text{Expected value for male Democrat} = \frac{130 \times 200}{420}$$

$$= 61.91$$

MIND
ONCE

$$\text{Expected value for female Democrat} = \frac{130 \times 220}{420}$$

$$= 68.09$$

$$\text{Expected value for male Independent} = \frac{50 \times 200}{420}$$

$$= 23.81$$

$$\text{Expected value for female Independent} = \frac{50 \times 220}{420}$$

$$= 26.19$$

NOW,

$$\chi^2 = \frac{(O_i - E_i)^2}{E_i}$$

$$\bullet \chi^2 \text{ are all sum of values} = (1.79 + 1.63 + 1.07 + 0.97 + 1.61 + 1.46)$$

$$\text{Degree degrees of freedom } (2-1) \times (3-1) = 8-4 = 4$$

$$= \frac{(1 \times 2)}{2}$$

we can see that alpha level of 0.05 and 2 degrees of freedom

the critical statistic is 5.91

the obtain statistic is 8.51

as $\chi^2 > 5.91$ so, H_0 rejected (Ans)

standard Error:- The SE is a single metric that sum up the variability the sampling distribution for a statistic

$$SE = \frac{s}{\sqrt{n}} \rightarrow (\text{sample})$$