Day-2

2/4/22

Agenda

1) Histogram
2) Measure of central thendency
3) Measure of Dispersion
4) Percentile & Quartiles
5) 5 Number Summary (box plot)

1) Histogram

Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 8

61, 65, 62, 72

90, 95, 100 }

i.) Sort the numbers

Ages are already sorted.

ii.) Bins → No. of groups

iii.) Bin size → size of bins.

if suppose we have in below size

[ 10, 20, 25, 30, 35, 40 ]  min = 10
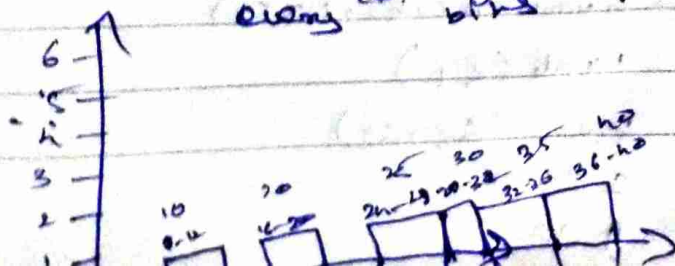
we want 10 bins        max = 40

bins = 10

∴ so we have to ge the 10 group.

of bins. if now take the

max (40) and divided it by bins (n)

40/10 = 4  ∴ so this is size of

every 10. bins.

Ex: 12

Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41,
42, 43, 50, 51, 65, 68, 78, 90, 95, 100}

**Histogram**

i.) Sort the values (Asc)
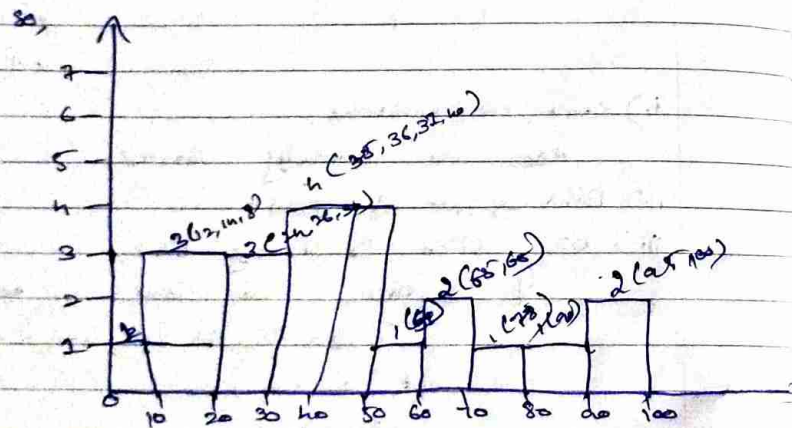   * Alred sorted

ii.) get the number of bins
   bins = 10

iii.) get the bin size based on max value
   & bins
   $\frac{100}{10} = 10$
   bin size = 10



0 - 10 = 1 value (10)
10 - 20 = 3 value (12, 14, 18)
20 - 30 = 3 value (24, 26, 30)
30 - 40 = 4 value (35, 36, 37, 40)
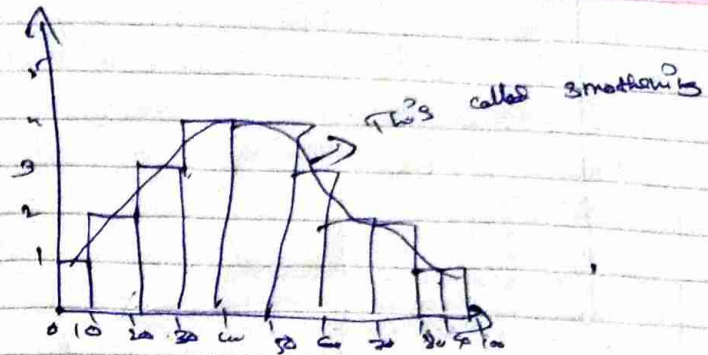40 - 50 = 4 value (41, 42, 43, 50)
50 - 60 = 1 value (51)
60 - 70 = 2 value (65, 68)
70 - 80 = 1 (78)
80 - 90 = 1 (90)



This is called smoothing

* Smoothening concept gives the
   probability density function.
* Kernal does distribution estimation

E + 2:

weight = {30, 35, 38, 42, 46, 58, 59, 62, 63
68, 75, 78, 90, 95}
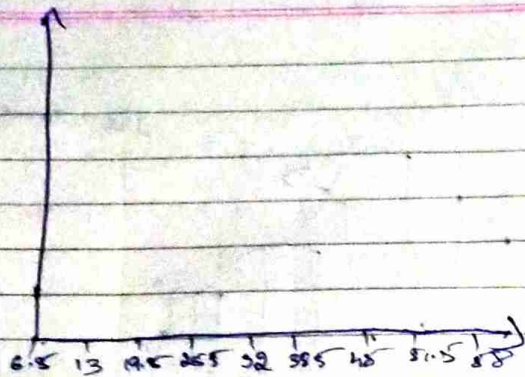
i.) Sort the weight (Asc)

ii.) get the bin
   bin = 10

iii.) If our values start from '0' then
   we can take max Number, and
   divide it by bin, we get bin size
   !. So here we have values
   start from 30 so we have
   to Substract the max max - min
   !. $\frac{95 - 30}{10} = \frac{85}{10} = 6.5$

6.5  13  19.5  26.5  32  38.5  45  51.5  58

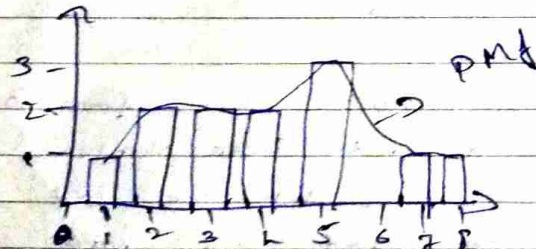**2) Measure of central Tendency**

1) Mean
2) Median
3) Mode

Discrete continues
2

no of Bank account = $\{2,3,5,1,4,5,3,7,8,3,2,4\}$

i) sort the value



pmf

0  1  2  3  4  5  6  7  8

pdf = Probability density function [Continuous var]
pmf = Probability mass function [Discrete var]

---

**2) Measure of central Tendency**

-) Measure of central Tendency is a single value that attempts to describe a set of data Identifying the central Position.

1) Mean (Average)

$$x = \{1,2,3,4,5\} =$$

$$\bar{x} = \frac{\sum x}{n} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

$$\boxed{\bar{x} = 3}$$

| Population (N) mean | N > n | Sample (n) mean |
|---|---|---|
| Mean → $\mu = \frac{\sum_{i=1}^{N} x_i}{N}$ | Mean $\bar{x} = \frac{\sum_{i=1}^{n} x_i}{n}$ | |
| | $\mu \gtrless \bar{x}$ or $\bar{x} \gtrless \mu$ | |

Ex: Population of Age $= \{24,23,2,1,28,27\}$

N = 6

P. Mean $(\mu) = \frac{24+23+2+1+28+27}{6}$

$$= \frac{105}{6}$$

$$\boxed{\mu = 17.5}$$

Ex: Sample of Age $= \{24,2,1,27\}$

$$\bar{x} = \frac{24+2+1+27}{4} = \frac{54}{4} = 13.5$$

$$\boxed{\bar{x} = 13.5}$$

## Practical Application [Feature Engineering]
### (mean)

| Age | Salary | Family size |
|-----|--------|-------------|
| | | |
| NAN | — | — |
| — | NAN | — |
| — | — | NAN |
| — | NAN | — |

→ loss of Data

* So here we have to handle this (NAN values my put the Average of the particular column.

NAN ⇒ Not A number

| Age | Salary |
|-----|--------|
| 24 | 45 |
| 28 | 50 |
| 29 | NAN |
| NAN | 60 |
| 31 | 75 |
| 36 | 80 |
| NAN | NAN |
| Av 29.6 | Av 62 |

⇩

| Age | Salary |
|-----|--------|
| 24 | 45 |
| 28 | 50 |
| 29 | 62 |
| 29.6 | 60 |
| 31 | 75 |
| 36 | 80 |
| 29.6 | 62 |

## 2) Median

$$\{1,2,3,4,5\} = \{1,2,3,4,5,100\}$$
↳ one layer

$$\bar{x} - 3 \qquad \Rightarrow \bar{x} = \frac{1+2+3+4+5+100}{6}$$

$$= \frac{19.16}{1.5}$$

* by adding one outlier we get more difference in mean so here median will come $\boxed{= 19.16}$

### Steps to find out the median.

i) Sort the values
ii) Find the center values
   i) if values are even no we find the values average
   ii) if values are odd no we find the central elements.

Ex: $\{1, 2, 3, 4, 5, 6, 7, 8, 100, 120\}$

1) Sort = $\{1, 3, 4, 5, \boxed{6, 7}, 8, 20, 10, 120\}$

2) center values = $\frac{6+7}{2} = \frac{13}{2} = 6.5$

$$\boxed{\text{Median is } 6.5}$$

* If there is No outliers we have to use Mean
* If there is outliers we have to use Median.

2) Mode

* Most Repeated values

$Age = \{1, 2, 2, \boxed{3, 3}, 4, 5\}$

Mode = 3

$Age = \{1, \boxed{2, 2}, \boxed{3, 3}, 4, 5\}$

$\boxed{mode = 2, 3}$

Practical example

Types of flower

Rose
LPlus &rarr; How to handle NAN
Sunflower    In this situation
Plase    Categorical Data with
NAN    we have to replace
Rose    Mode value instead
$\boxed{mode \Rightarrow Rose}$    of NAN value

Rose
LPlus
Sunflower
Rose
Rose
Rose

3) Measure of Dispersion

1) Variance $(\sigma^2)$
2) Standard $(\sigma)$

Variance &rarr; spread of data

| Population variance $(\sigma^2)$ | $\sigma = s$ | Sample variance $(s^2)$ |
|---|---|---|
| $\sigma^2 = \dfrac{\sum\limits_{i=1}^{N}(x_i - \mu)^2}{N}$ | $(n-1)$ | $s^2 = \dfrac{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2}{n-1}$ |

$\therefore x_i - \mu \text{ or } x_i - \bar{x}$

basically

it gives the difference of x value
to mean value
$x \Rightarrow$ values
$\mu, \bar{x} \Rightarrow$ mean $(N, n)$

1) $\{1, 2, 3, 4, 5\}$    2) $\{1, 2, 3, 4, 5, 6...\}$

$\mu = 3$    $n = (n-1)$

$\sigma^2 = \dfrac{[(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2]}{5}$

$= \dfrac{4 + 1 + 0 + 1 + 4}{5} = \dfrac{10}{5} = 2$

$\boxed{\sigma^2 = 2}$

$\sigma^2 = \dfrac{[(1-x)^2 + (2-x)^2 + (3-x)^2 + (4-x)^2 + (5-x)^2 + (6-x)^2]}{7}$

$\sigma^2 = \boxed{7 \cdot 9 \cdot 6}$

## 2) Standard deviation $(\sqrt{\sigma^2})$

$\{1,2,3,4,5\}$

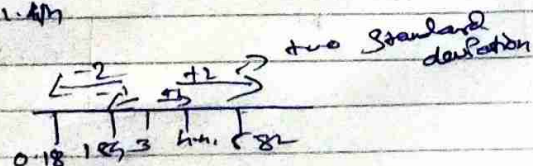$\mu = \frac{15}{5} = 3$

$\sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5}$

$= \frac{4+1+0+1+4}{5} = \frac{10}{5} = 2$

$\sigma^2 = 2$

$\sigma = \sqrt{2} = 1.41$

two Standard deviation

$\xleftarrow{\;\;-2\;\;}\; \xrightarrow{\;\;+2\;\;}$
$\frac{-2}{-2} \quad +2$

0.18  1.59  3  4.41  5.82

**How many** standard deviation away from the mean.

Ex: what is the standard dev of 4
it is 4.1
because standard deviation value is should add$^{on}$ $_{Sub}$ to the mean and see how much time to take to reach the given number (4)

---

## 4) percentile And Quartiles

Percentage = $\{1,2,3,4,5,6,7,8\}$
$\downarrow$

Per $g$ even Number
Pn percentage $\Bigg\} = \dfrac{\text{No. } g \text{ even numbers}}{\text{Total No } g \text{ Number}}$

$= \frac{4}{8} = 0.5 = 50\%$

### Percentile

**Def**

A percentile is a value below which a certain percentage $g$ observations lie.

99 percentile → It means the Person has got better than than 99% g entire Students.

**Dataset**

$A = \{2,2,3,4,5,5,5,6,7,8,8,8,8,9,9,9,10,11,11,12\}$

i.) What is the percentile rank g 10
i.) Sort the data first
ii.) Percentile Rank g $x = \dfrac{\text{No. } g \text{ Values below}(x)}{h}$

$= \frac{16}{20 \; 5} \quad \frac{4}{8} = 0.8 = 80\%$

ii) What is the value that exists at 25 percentile

value = $\frac{\text{percentile}}{100} \times$ → for odd numbers

for even (n+1)

$= \frac{25}{100} + 20$

$= 5^{th}$ index value of the given values.

iii) What is the value that exists at 20 percentile.

Value $= \frac{20}{100} \times 20 + 1$
$\frac{}{5}$

$= \frac{20 + 21}{100}$

$n \cdot 2$

$= \frac{1}{5} + 21 = 21/5 = n \cdot 2 = n^{th}$ index of given

5) number summary $\overset{-}{=} \frac{}{2}$

  1.) Minimum
  2.) First quartile (25 percentile)(Q1)
  3.) Median
  4.) Third quartile (75 percentile)(Q3)
  5.) maximum

Box plot
↑
Remove the outlier

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

→ i) To find outliers
  ↓ [Lower Fence → Higher Fence]

Lower fence → Q1 → 1.5 (IQR) : IQR = Q3 - Q1
  ↓                              75 - 25
  (Q2)                       Inter Quartile
                              Range (IQR)

Higher fence = Q3 + 1.5 (IQR)
                  ↓
                 (75)

$Q_1 = \frac{25}{100} + (n+1)$

$= \frac{25}{100} \times (20 + 1)$
           5.2.1

$= \frac{25}{100} \times 21$

$= \frac{25 \cdot 5 \cdot 2}{}$

$\boxed{Q_1 = 5.2}$ → Index ⇒ $\boxed{\frac{3+3}{2}} = \frac{6}{2} = 3$

$Q_3 = \frac{75}{100} \times (21)$

$= 15.75^{th}$ Index value so 15 & 16$^{th}$ value

$= \frac{7+8}{2} = \frac{15}{2} = 7.5$

$\boxed{Q_3 = 7.5}$

Lower fence ⇒ $3 - (1.5)(4.5) = -3.65$

Higher fence ⇒ $7.5 + 1.5 (4.5) ⇒ 14.25$

$(-3.65 \longleftrightarrow 14.25)$

in between these values only we should
have . other dot is an outliers.
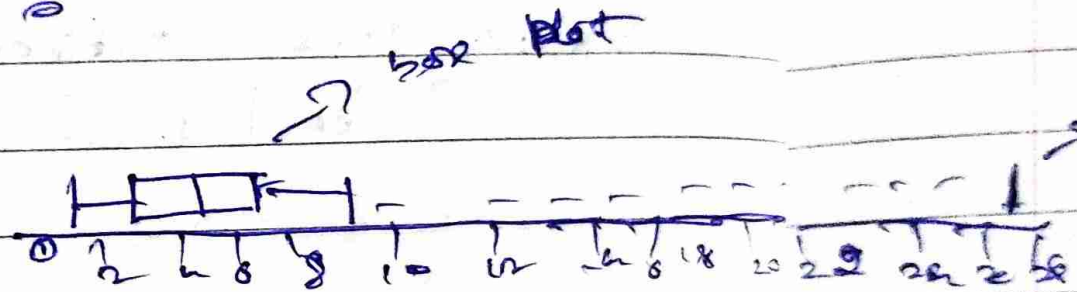
To remove out lier

box plot

1) Minimum — 1

2) $q_1$ — 3

3) media — 5

4) $q_3$ — 7.5

5) maximum — 9

→ outlier

→ outlier