

STATISTICS FOR DATA SCIENCE Day 2

Agenda

- (i) Histogram
- (ii) Measure of Central Tendency
- (iii) Measure of Dispersion
- (iv) Percentile and Quartiles
- (v) 5 Number Summary (Box Plot)

(i) HISTOGRAM

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

How to create histogram?

Step 1: Sort the Numbers

Step 2: Bins \rightarrow No. of groups

Step 3: Bin size \rightarrow Size of Bins

[10, 20, 25, 30, 35, 40]

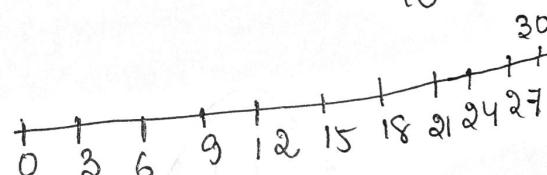
$$\text{min} = 10, \text{max} = 40$$

$$\text{Bins} = 10$$

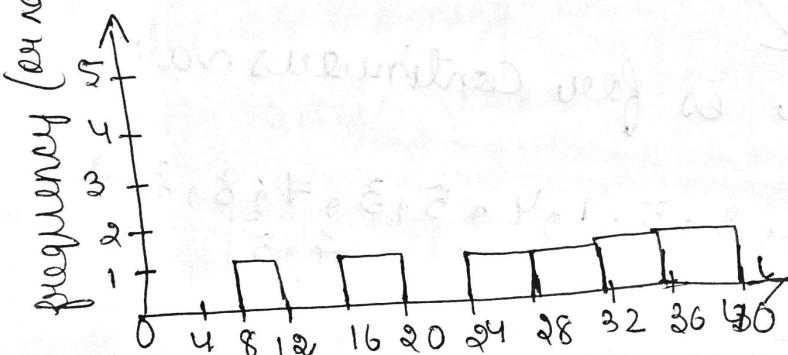
$$\text{bin size} = \frac{40-10}{10} = \frac{30}{10} = 3$$

[10, 20, 25, 30, 35, 40]

$$\text{min} = 10, \text{max} = 40$$



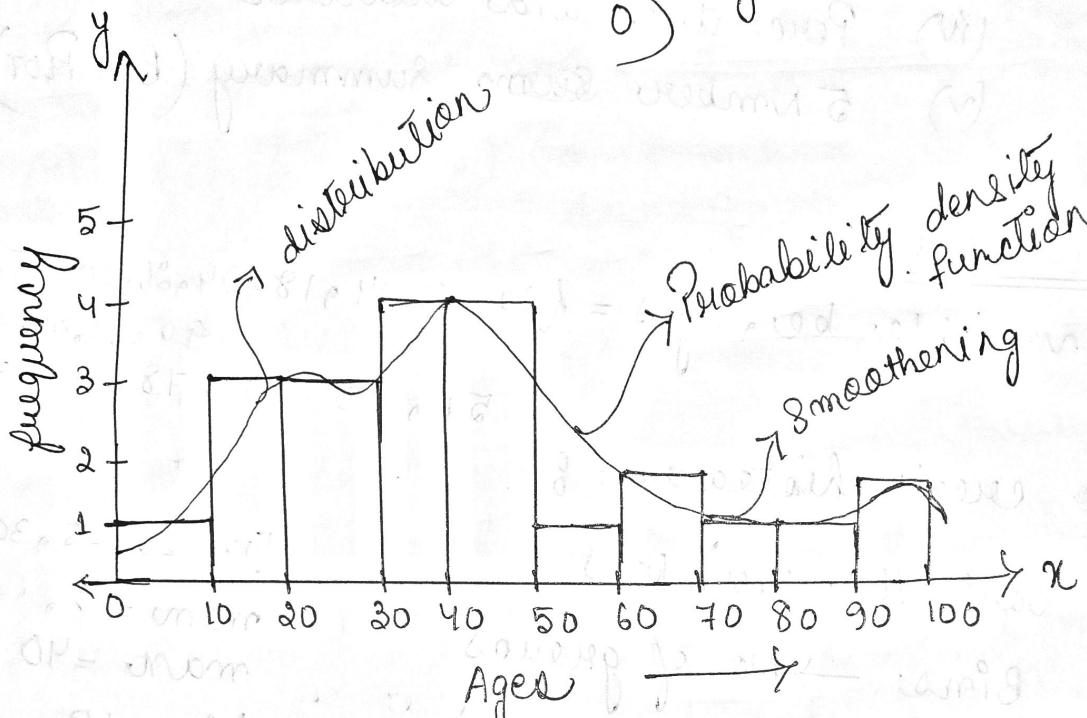
\therefore we got 10 bins.



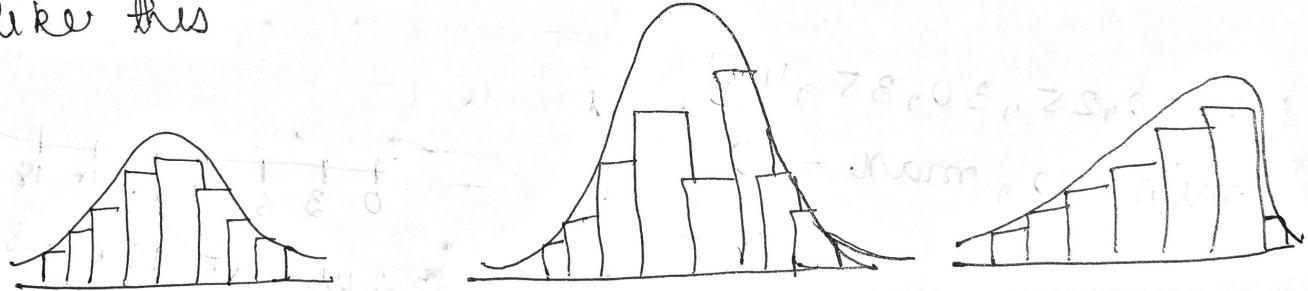
e.g. Let, the bin = 10 for this data

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

min = 10, max = 100
bin size = $\frac{100-0}{10} = 10$ (we are starting from 0
that's why we are subtracting by 0)



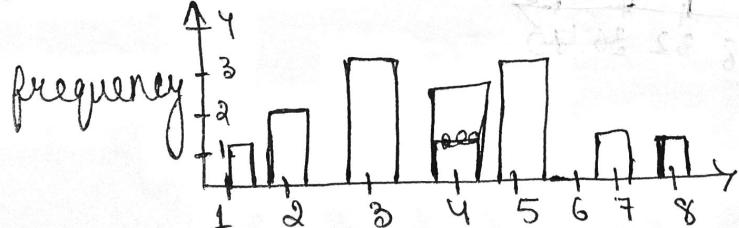
Smoothing will help to form the distribution like this



↓
This one is for continuous value

*Discrete Continuous

No. of Banks accounts = [2, 3, 5, 9, 4, 5, 3, 7, 8, 3, 1, 4, 5]

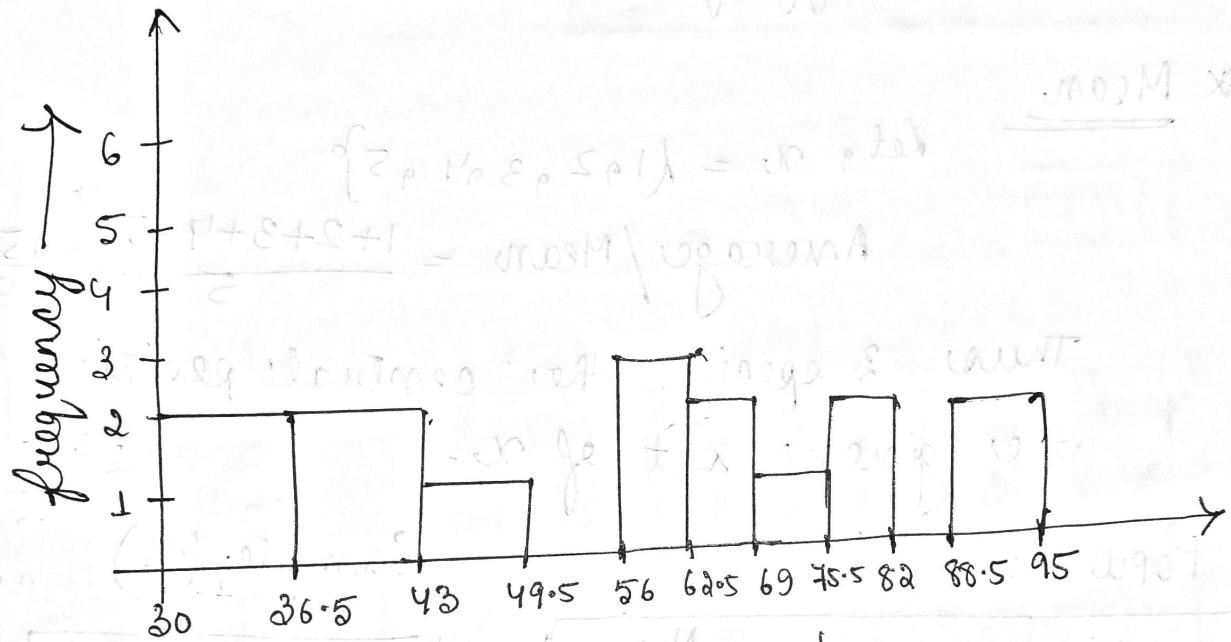


$$\text{E.g. weight} = \{30, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, 80, 90, 95\}$$

bins = 10.

If we want to start from 30, then

$$\text{bin size} = \frac{95 - 30}{10} = \frac{65}{10} = 6.5$$



$$\text{Normal Algme} \\ \frac{\sum x^2}{n} = \bar{x}^2 \\ \sum x^2 = \bar{x}^2 n$$

$$\text{Weights} \rightarrow \text{normal algme}$$

- * For smoothing the continuous histogram we will use Probability density function (pdf)
- * For smoothing the discrete continuous histogram we will use Probability mass function (pmf)

(ii) Measure of Central Tendency

- (a) Mean (b) Median (c) Mode

* A Measure of Central Tendency is a single value that attempts to describe the set of data identifying the central position.

* Mean

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

Average / Mean = $\frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$

Thus 3 specify the central position of the given set of x .

Population (N)

Sample (n)

$$\text{Population mean } (\mu) = \sum_{i=1}^N x_i / N$$

$$\text{Sample mean } \bar{x} = \sum_{i=1}^n x_i / n$$

Example

and off course $N \neq n$

Population Ages = $\{24, 23, 2, 1, 28, 27\}$

$$\therefore N = 6$$

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

$$= \frac{105}{6} = 17.5$$

Sample Age = {24, 29, 31, 27} [Let us pick randomly 4 values from Population Age]

$$\text{Sample mean } (\bar{x}) = \frac{24+29+31+27}{4}$$

$$= \cancel{24} + \cancel{29} + \cancel{31} + \cancel{27} = 100$$

$$= \frac{54}{4} = 13.5$$

PRACTICAL APPLICATION (Feature Engineering)

Let, the three features be

Age	Salary	Family Size
-	-	-
-	-	-
-	-	-
NAN	-	-
-	NAN	-
-	-	-
-	-	NAN
-	NAN	-
NAN	-	-

- * If we drop the NAN here then we will loss some of our information.
- * we will use mean here in place of NAN value.

Age	Salary
24	45
28	50
29	NAN
NAN	60
31	75

Age	Salary
NAN	80
NAN	NAN

$$\text{Age mean} = \frac{24+28+29+31+36}{5}$$

Replacing NAN value by 29.6

If we add one more value 80 to the set then the NAN value will change to 30
(here 80 is outlier)

$$\text{Salary mean} = \frac{45+50+60+75+80}{5}$$

Replacing NAN value by 62 on salary set

If we add one more value 200 to the set then the NAN value will change to 85
(here 200 is outlier)

Median

$$\begin{aligned} &\{1, 2, 3, 4, 5\} \\ &\bar{x} = \frac{1+2+3+4+5}{5} \\ &= 3 \end{aligned}$$

$$\begin{aligned} &\{1, 2, 3, 4, 5, 100\} \\ &\bar{x} = \frac{1+2+3+4+5+100}{6} \\ &= 19.16 \end{aligned}$$

NAN \rightarrow Null
see
Not number

By adding one outlier, the mean is changing $\bar{x} = 3 \rightarrow \bar{x} = 19.16$, which is huge.

so we will use median here

Steps to find out median

(i) Sort the numbers

(ii) Find the central number

(a) If the no. of elements are even we find the avg. of central elements

(b) If the no. of elements are odd we find the central element.

For e.g. 1, 2, 3, 4, 5, 6, 7, 8, 100, 120

∴ no. of elements is even [100 & 120 are outliers]
 $\therefore \text{median} = \frac{5+6}{2} = 5.5$ [last]

(since mean \neq median)

$$\therefore \text{mean} = 25.6$$

* when there is no outlier we use mean
for NAN value

* when there is outlier we use median
for NAN value

MODE : Most frequent occurring element

e.g. $\{1, 2, 2, 3, 3, 3, 4, 5\}$ $\{1, 2, 2, 2, 3, 3, 3, 4\}$

↓ ↓

mode = 3 mode = 2 & 3.

Practical Application

Dataset

Types of flower

Lily

Sunflower

Rose

NAN

Rose

Sunflower

Rose

NAN

* Mode we will use in categorical dataset

since here mode is Rose, so will replace NAN value by Rose

Measure of Dispersion (i) Variance (σ^2)

(ii) Standard deviation (σ)

* Variance (σ^2) → Spread of Data

Population Variation (σ^2)

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

Sample Variance (s^2)

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

$x_i - \mu \rightarrow$ distance from mean with every data point

~~1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13~~ ~~1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13~~

Assignment ①: why sample variance is divided by $n-1$?

* Standard deviation ($\sqrt{s^2}$)

{1, 2, 3, 4, 5} calculated using witness (vii)

~~1, 2, 3, 4, 5~~ \downarrow

$$\text{mean} = 3$$

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

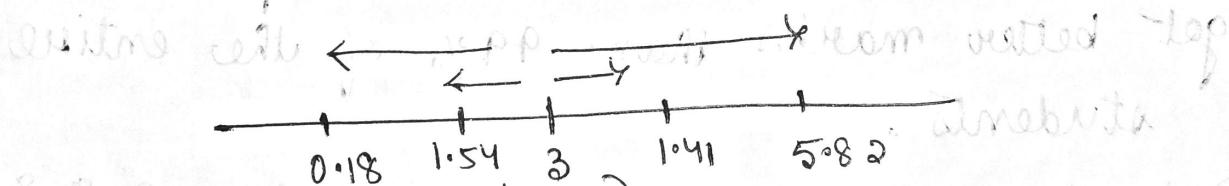
$$\text{and variance } (s^2) =$$

$$\text{standard deviation} = \sqrt{\frac{4+1+1+4}{5}} = \sqrt{2} \approx 1.41$$

$$= \sqrt{6^2/5} = \sqrt{2}$$

$$= 1.41$$

Calculated with test sheet no. 16 = witness P.

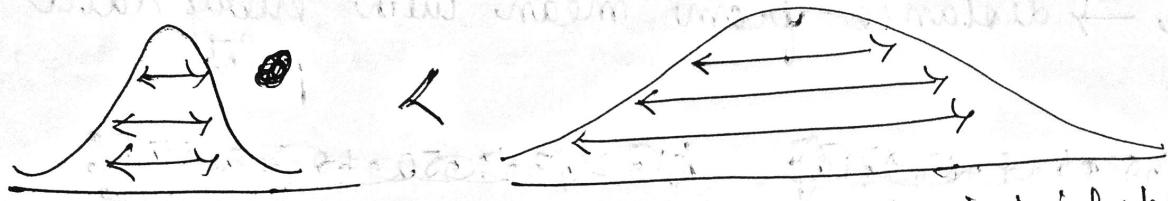


Standard deviation means how many standard deviations away the no. from the mean

So, in the data {1, 2, 3, 4, 5}, the no. 2 is

away by ~~one~~ standard deviation from 3.
-0.1

Witness P = $\frac{1}{5} = 0.2$ for witness



spread is less here

spread is high here

Variance is high

Standard deviation is high

(iv) Percentile And Quartiles

Percentage -

Percentiles : MATHE, CAT, SAT; JEE, NEET \Rightarrow marks are count on basis of Percentile

Defⁿ: A percentile is a value below which a certain percentage of observation i.e.

99 percentile = It means that the person has got better marks than 99% of the entire students.

Dataset : 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

Q. what is the percentile rank of 10

\Rightarrow Percentile rank of $x = \frac{\text{no. of value below } x}{n}$

Percentile rank of 10 = $\frac{16}{20} = 80$ percentile

Percentile rank of 8 = $\frac{9}{20} = 45$ percentile

Q what is the value that exist at 25 percentile

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

$$= \frac{25}{100} \times 21 = 5$$

Q what is the value that exist at 95 percentile

$$\text{Value} = \frac{95}{100} \times 21 = 19.95$$

5 numbers Summary

i) Minimum (25 percentile) (Q1)

ii) First Quartile (25 percentile) (Q1)

iii) Median

iv) Third Quartile (75 percentile) (Q3)

v) Maximum

These 5 numbers Summary used to remove outliers.

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

here 27 is the outlier.

[Lower Fence \leftarrow Higher Fence]

$$\text{Lower Fence} = Q1 - 1.5(\text{IQR})$$

$$\text{Higher Fence} = Q3 + 1.5(\text{IQR})$$

$$\text{IQR} = Q3 - Q1$$

Interquartile Range

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

$$= \frac{25}{100} \times 21 = 5.25 \Rightarrow \text{Index} = \frac{3+3}{2} = 3 \quad (\text{Adding } 5^{\text{th}} \text{ and } 6^{\text{th}} \text{ index value})$$

$$Q_3 = \frac{75}{100} \times 21 = 15.75 \Rightarrow \text{Index} = \frac{8+7}{2}$$

$$= \frac{15.5}{2} = 7.5$$

$$\text{Lower Fence} = Q_1 - 1.5(IQR) = 3 - 1.5(4.5) \\ = -3.65$$

$$\begin{aligned} IQR &= Q_3 - Q_1 \\ &= 7.5 - 3 \\ &= 4.5 \end{aligned}$$

$$\begin{aligned} \text{Higher Fence} &= Q_3 + 1.5(IQR) \\ &= 7.5 + 1.5 \times 4.5 \\ &= 14.25 \end{aligned}$$

Now,

~~1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100~~

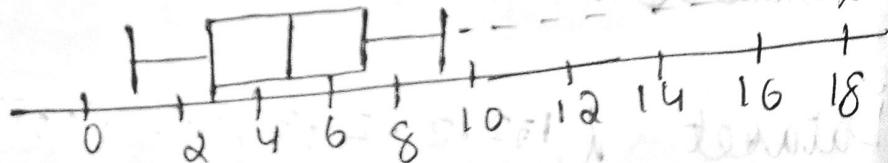
① Minimum = 1

② $Q_1 = 3$

③ Median = 5

④ $Q_3 = 7.5$

⑤ Maximum = 9



To treat outliers

[treat with \rightarrow new values]

18-50-58-59

(Q1) 2.1 - 12.5 and out