

statistics

- ① → Histograms
- ② → Measure of Central tendency
- ③ → Measure of Dispersion
- ④ → Percentiles and Quartiles
- ⑤ → 5 Number Summary (Box plot).

Histogram:

Ages = {10, 12, 14, 18, 24, 26, 30, 34, 35, 36, 37, 40, 41, 43, 45, 50, 51, 68, 78, 90, 95, 100}

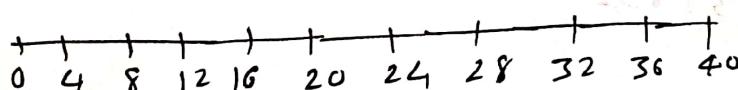
- ① Sort the Numbers
- ② Bins → No. of groups
- ③ Bins size → Size of Bin

$$\boxed{[10, 20, 25, 30, 35, 40]}$$

bin = 10
Min = 10
Max = 40

$$\frac{40-10}{10} = \frac{30}{10} = 3$$

$$\frac{40}{10} = 4.$$



Find the $\boxed{\text{bin size} = 10}$ of Ages.

$$\text{bin size} = \frac{100}{10} \Rightarrow \frac{\text{total Ages}}{\text{bin}}$$

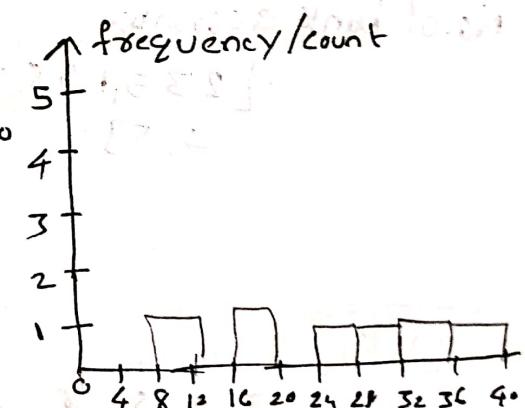
$$\therefore \frac{100}{10} = \frac{100}{10} = 10,$$

bin size are

decided by US//

or wish//

(M)



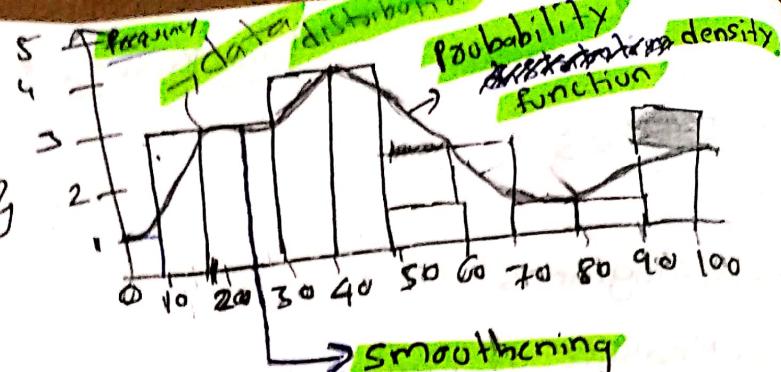
let bins = 20

$$\text{bin size} = \frac{100}{20} = 5$$

$$\boxed{\text{bin size} = 5}$$

Start from 0

Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 55, 68, 78, 90, 95, 100}



Group = bin.

Group size = bin size

If we have negative (-ve) values our bins also get or have (-ve) values.

Weight = [30] 35, 38, 42, 46, 58, 59, 63, 63, 68, 75, 77, 80, 90, 95]

{ Probability density function }

{ Continuous values }

bins = 10

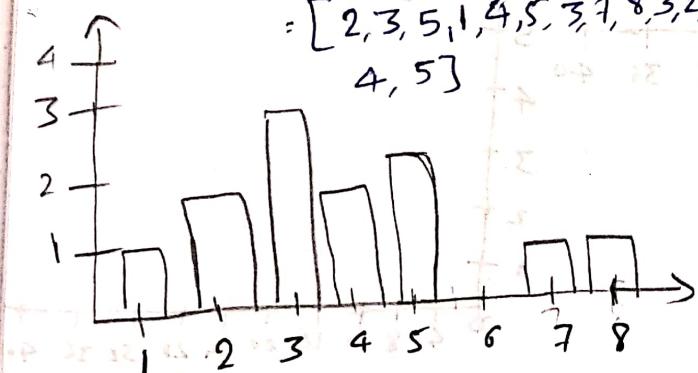
$$\text{bin size} = \frac{95 - 30}{10} = \frac{65}{10} = 6.5$$

(2) Measures of C

Discrete Continuous:

No. of bank accounts:

$$= [2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5]$$



In discrete continuous we are able to find gaps.

PDF: Probability density function

PMF: Probability Mass function

To smooth these discrete continuous

We use { Probability Mass function }



pdf: Probability density function
 PMF: Probability Mass Function

② → Measure of Central tendency.

- ① Mean
- ② Median
- ③ Mode

definition

A Measure of CT is a single value that attempts to describe a set of data identifying the central position

Central Position

Mean

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

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

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

Population (N)

Sample (n)

Population mean (μ)

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

$$N \gg n$$

Sample mean (\bar{x})

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

$$n \geq N$$

never happens

Population

$$\text{Age} = \{24, 23, 21, 28, 27\}$$

$$\text{Population (N)} = 6$$

$$\text{Population mean} (\mu) = \frac{24 + 23 + 21 + 28 + 27}{6}$$

here

$$N > n$$

$$\text{Sample Age} = \{24, 21, 27\}$$

$$\text{Sample (n)} = 4$$

$$\text{Sample mean} (\bar{x}) = \frac{24 + 21 + 27}{4}$$

$$\begin{aligned} \mu &\geq \bar{x} \\ \bar{x} &\geq \mu \end{aligned}$$

$$\bar{x} = 23.5$$

Practical application of Feature Engineering

Used in ~~Machine Learning~~

NAN = Null Values
Not A Number

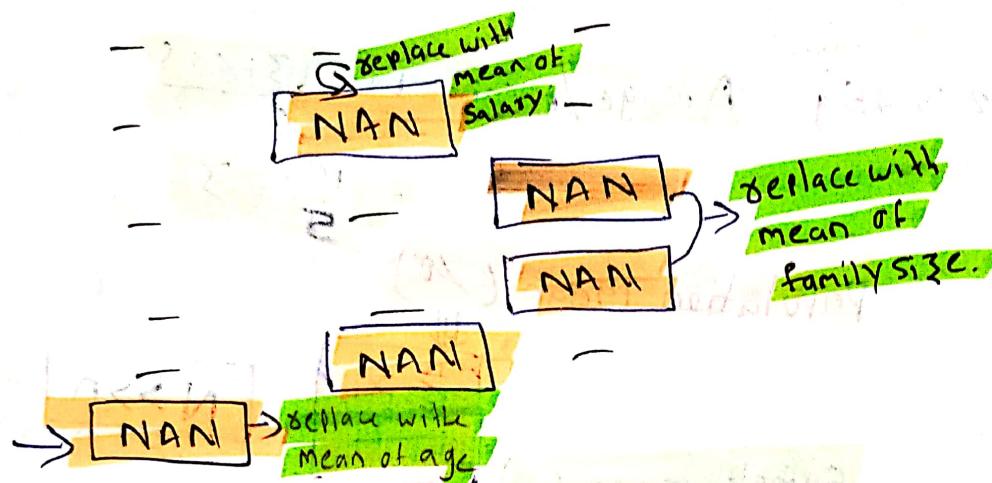
I have features like



Age	Salary	Family size
-	-	-
-	-	-
NAN	-	-

if we drop this table

We loss of info



if we drop the table we will loss the info so that, we have to the total out come.

instead of this we have to find

Mean of Age and replace the NAN value with the Mean of Age so that you cannot loss the info

Mean is talking about the

Central Information of data

<u>Ex:</u>	<u>Age</u>	<u>Salary</u>	<u>Mean of Age</u>	<u>Mean of Salary</u>
	24	45	$24 + 28 + 21 + 31 + 36 \over 5$	$45 + 50 + 60 + 75 + 80 \over 5$
	28	50	$= 29.6$	$= 62$
	29	NAN		
	31	60		
	32	NAN		
	31	75		
	32	80		
	32	NAN		
	80	NAN		
	200	⇒ outliers		
<u>outliers</u>				

$$\text{mean of Age after adding outliers} = \frac{24 + 28 + 21 + 31 + 36 + 80}{6} = 38.11$$

$$\text{mean of Salary after adding outliers} = \frac{45 + 50 + 60 + 75 + 80 + 200}{6} = 85.11$$

$$\text{before outliers mean of Age} = 29.6$$

$$\text{mean of Age}$$

$$\text{after outliers mean of age} = 38$$

$$\text{before outliers mean of salary} = 62$$

$$\text{after outliers mean of salary} = 85$$

there a huge difference by adding outliers to the mean of Age, Salary.

So to avoid these difference we use

Median

Median

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

$$\frac{1+2+3+4+5}{5} = 3$$

$$\boxed{\bar{x} = 3, \bar{x} = 19.16}$$

outlier

$$\{1, 2, 3, 4, 5, \boxed{100}\}$$

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

$$= \frac{115}{6} = 19.16$$

Just by adding one outlier the mean is changing a huge.

in order to prevent this median is used.

⇒ Steps to find out Median

- ① ⇒ Sort the Numbers
- ② ⇒ Find the central numbers

if the no of elements are even we find the average of central elements

if the n.o of Elements are odd we find the central elements

Sorted no of elements are even.

$$\{1, 2, 3, 4, \boxed{5, 6}, 7, 8, 9, 10, 11\}$$

$$\text{median} = \frac{5+6}{2} = 5.5$$

if we are adding "0"

$$\{0, 1, 2, 3, 4, \boxed{5}, 6, 7, 8, 9, 10, 11\}$$

now no of elements are odd

then we take the central

$$\text{Median} = 5 \quad \text{Element} = 5$$

~~Note:-~~ When ever we have outlier we have to use median.

{ No outliers \rightarrow Mean }
with outliers \rightarrow Median }

③ Mode - { The Most frequently occurring elements }

Ex { 1, 2, 2, [3, 3, 3], 4, 5 }

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

Dataset:-

Types of flowers

Lily

Sunflower

Rose

NAN

Rose

Rose

Sunflower

Rose

NAN

Rose

in recent version of python it is giving you errors, but in previous version of python it use to give you

[2, 3]
Mode

Mode most of time we use with categorical value variable

③ Measure of Dispersion

- (a) Variance (σ^2) ← Spread of data
(b) Standard deviation (σ)

Variance

Population variance (σ^2)

Sample variance (s^2)

Population variance (σ^2)

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

Sample variance (s^2)

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

Ex:

First One

||

Variance ↑↑
Spread ↑↑

second one

||

$$\{1, 2, 3, 4, 5, 7, 7, 8, 9, 10\}, \{1, 2, 3, 4, 50, 60, 70, 100\}$$

||

Variance

Ex: 2

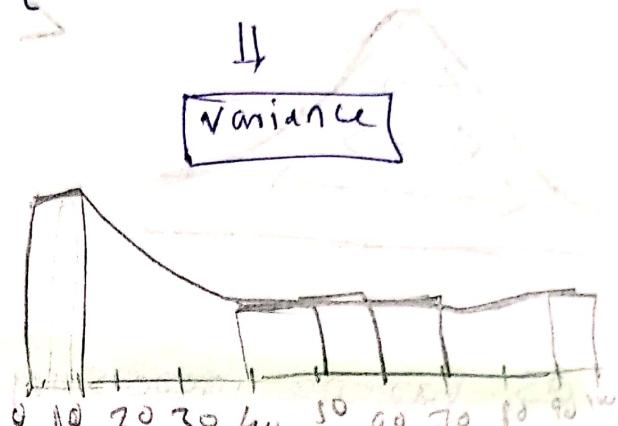
$$\textcircled{a} \{1, 2, 3, 4, 5\}$$

$$\mu = \frac{1+2+3+4+5}{5}$$

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

$$\mu = 3$$

Variance ↑↑
Spread ↑↑
Variance ↑↑



Variance ↑↑ Spread ↑↑

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

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

$$\textcircled{b} \{1, 2, 3, 4, 5, 6, 80\}$$

$$\mu = \frac{1+2+3+4+5+6+80}{7} = \frac{101}{7} = 14.4$$

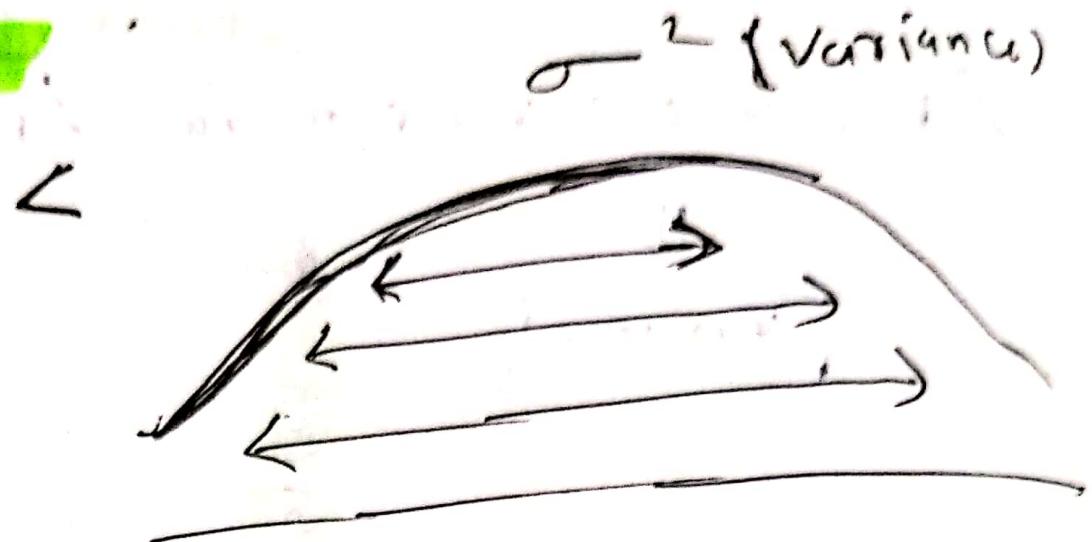
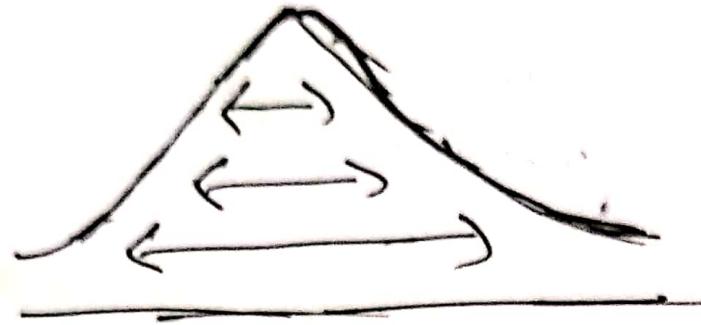
$$\sigma^2 = \frac{[(1-14.4)^2 + (2-14.4)^2 + (3-14.4)^2 + (4-14.4)^2 + (5-14.4)^2 + (6-14.4)^2 + (80-14.4)^2]}{7}$$

$$\sigma^2 = 749.10$$

||

by graph

σ^2 (variance)



AS Variance increasing spread also increasing

So Variance ↑↑ spread ↑↑

∴ Variable is ~~nothing~~ nothing but which indicating the spread of the data.

⑥ b

Standard deviation: (Ans)

$$(\sqrt{\sigma^2}) = \sigma$$

Standard deviation = σ

$$\therefore \{1, 2, 3, 4, 5\} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

$\boxed{\mu = 3}$

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

$$\sigma^2 = \frac{4+1+0+1+4}{5} = 2.11$$

$$\sigma^2 = 2.11$$

$$\sigma = \sqrt{2.11} = 1.41.$$

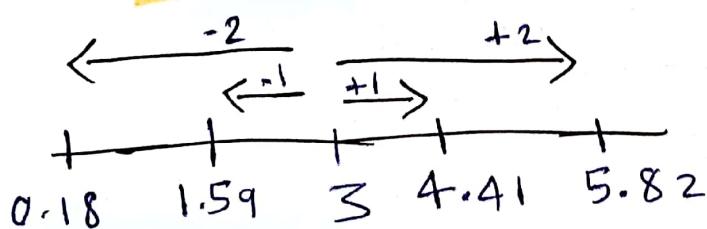
$\boxed{\mu = 3}$

When we are moving one step backward (-1) which means

Subtracting σ (1.41)

$$3 - 1.41 = 1.59$$

$$1.59 - 1.41 = 0.18$$



When we are moving one step forward (+1) which means adding σ (1.41)

$$3 + 1.41 = 4.41$$

$$4.41 + 1.41 = 5.82$$

Standard deviation tells that:-

Suppose we are taking value "4" how many standard deviation will fall away from mean.

like \Rightarrow how many standard deviation "4" is falling away from mean is within (+1) range

\Rightarrow how many standard deviation "3" is falling away from mean is "0"

Variance talk about Spread

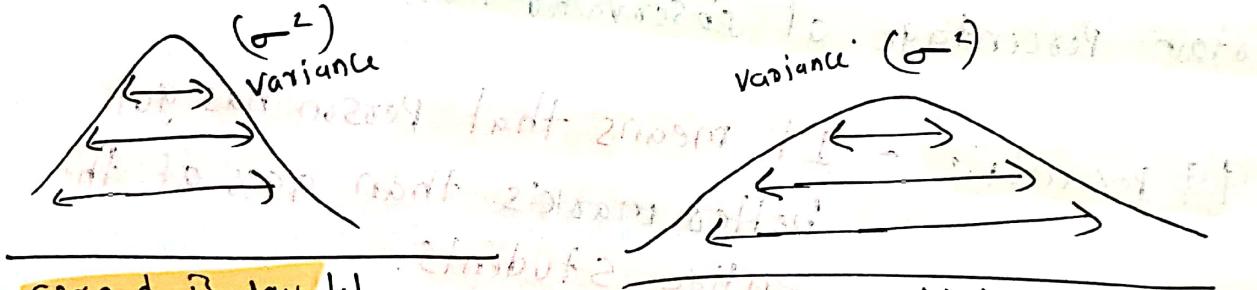
* *
*.

How the data is spread

Variance $\uparrow \uparrow$

Spread $\uparrow \uparrow$

Spread of data is calculated by variance



spread is low $\downarrow \downarrow$

variance is low $\downarrow \downarrow$

spread is high $\uparrow \uparrow$

variance is high $\uparrow \uparrow$

Standard deviation says that $\{1, 2, \dots, 1000\}$

in side these numbers, if i take any number i will

group of

be able to find how many standard deviation is away from the mean.

* Percentiles & Quartiles:

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

Percentage of even numbers

$$= \frac{\text{No. of even numbers}}{\text{total no. of numbers}}$$

$$= 4/8 = 0.5 = 50\%$$

Percentiles: Gate, IAT, IELTS, SAT, GRE, JEE, NEET

Definition:

A Percentile is a value below which a certain percentage of observations lies.

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

Ascending order:

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

what is the Percentile rank of 10

Percentile Rank of $x = \frac{\# \text{ No. of Value below } x}{n}$

$$= \frac{16}{20} = 80 \text{ Percentile}$$

what is the value that exist at 25 percentiles.

for even

$$\text{Value} = \frac{\text{Percentiles}}{100} * n + \frac{1}{2}$$

$$= \frac{25}{100} \times \frac{21}{21} + 1 = 5^{\text{th}} \text{ index}, \\ \text{output} = 5$$

for odd

$$\text{Value} = \frac{\text{Percentiles}}{100} * n$$

$$= \frac{95}{100} \times 21 = 19.95^{\text{th}} \text{ index} \\ \text{output} = 19$$

* 5 number Summary

(x) outliers

① Minimum

② First Quartile (25 Percentile) (Q1)

③ Median

④ Third Quartile (75 Percentile) (Q3)

⑤ Maximum

these 5 summary

used to

remove

the

outliers,

$$[2, 4] = 2.5$$

$$[2, 4, 6, 8] = 4.5$$

Ex: $\{1, 2, 2, 2, 2, \boxed{3}, 3, 4, 5, 5, 5, 6, 6, 6, \boxed{7}, 8, 8, 9, 22\}$

[lower fence] \longleftrightarrow [higher fence]

$$\text{Lower fence} = Q_1 - 1.5(\text{IQR})$$

$$\text{IQR} = Q_3 - Q_1$$

IQR \Rightarrow Inter Quartile Range (IQR)

$$\text{Lower fence} = Q_1 - 1.5(\text{IQR})$$

$$\text{Higher fence} = Q_3 + 1.5(\text{IQR})$$

$$Q_1 = \frac{25}{100} * (n+1) = \frac{25}{100} * 21 = 5.25 \Rightarrow \text{index}$$

$$5.25 \Rightarrow \text{index} = 3_{11}$$

$$\text{We don't have } 5.25 \text{ so we take average or } \frac{3+3}{2} = \frac{6}{2} = 3_{11}$$

$$Q_3 = \frac{75}{100} * 21 = 15.75 \quad \text{Index} = \frac{8+7}{2} = 7.5$$

$$Q_3 = 7.5$$

$$\text{Lower fence} = 7.5 + 1.5(4.5) =$$

$$3 - (1.5)(4.5) = \boxed{-3.65}$$

$$\text{Higher fence} = 7.5 + (1.5)(4.5) = \boxed{14.85}$$

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

① **minimum = 1**

② **$Q_1 = 3$**

③ **median = 5**

④ **$Q_3 = 7.5$**

⑤ **maximum = 9**

