

Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37,

40, 44, 42, 43, 50, 51, 65, 68, 78, 90, 95,

100}

max = 100

$$\text{Bin Size} = \frac{100}{20} = 5 \quad \text{min} = 10$$

$$\text{Bin Size} = \frac{\text{max}}{\text{bins}}$$

Group = bin size

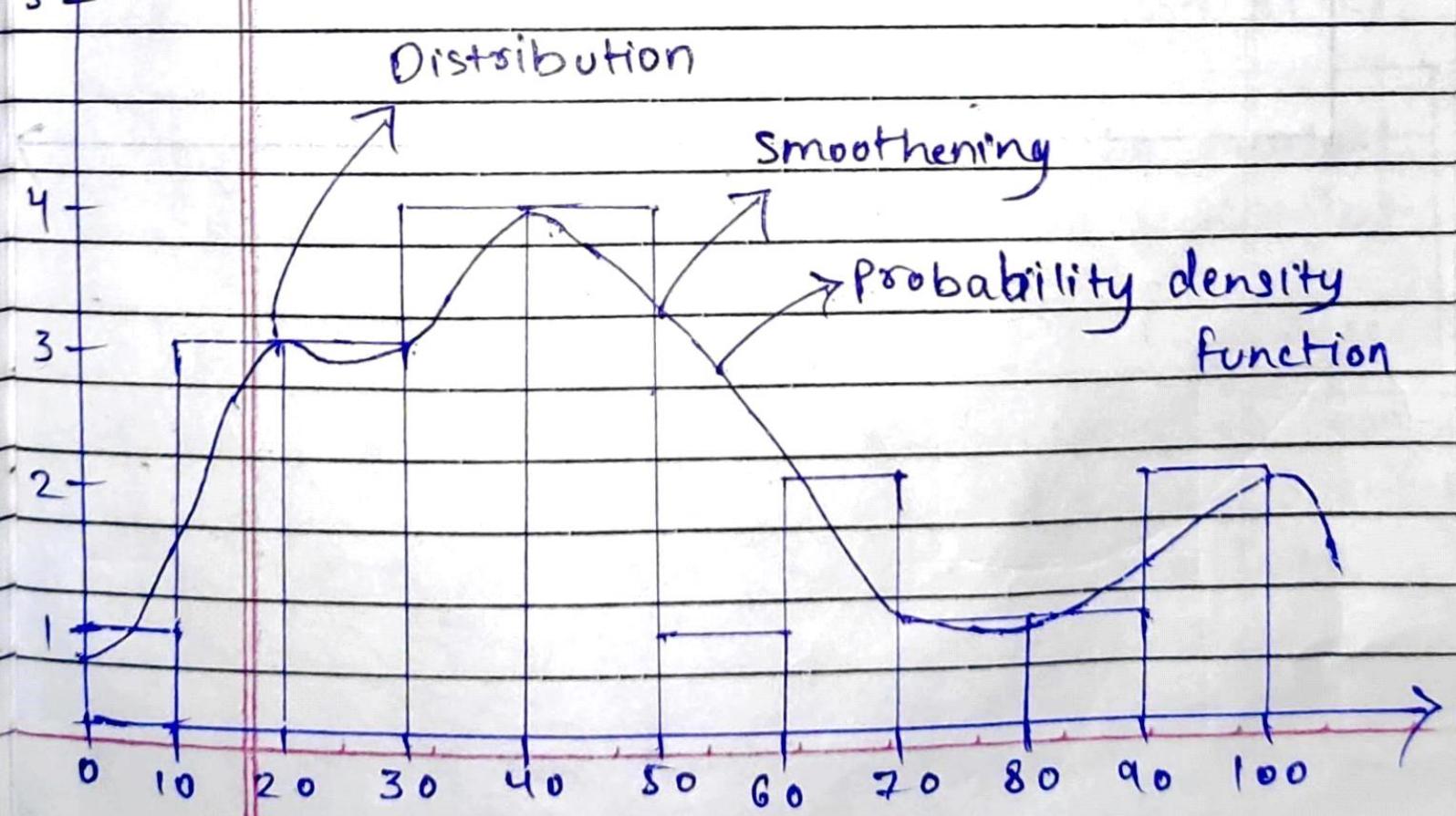
Group size = bin size

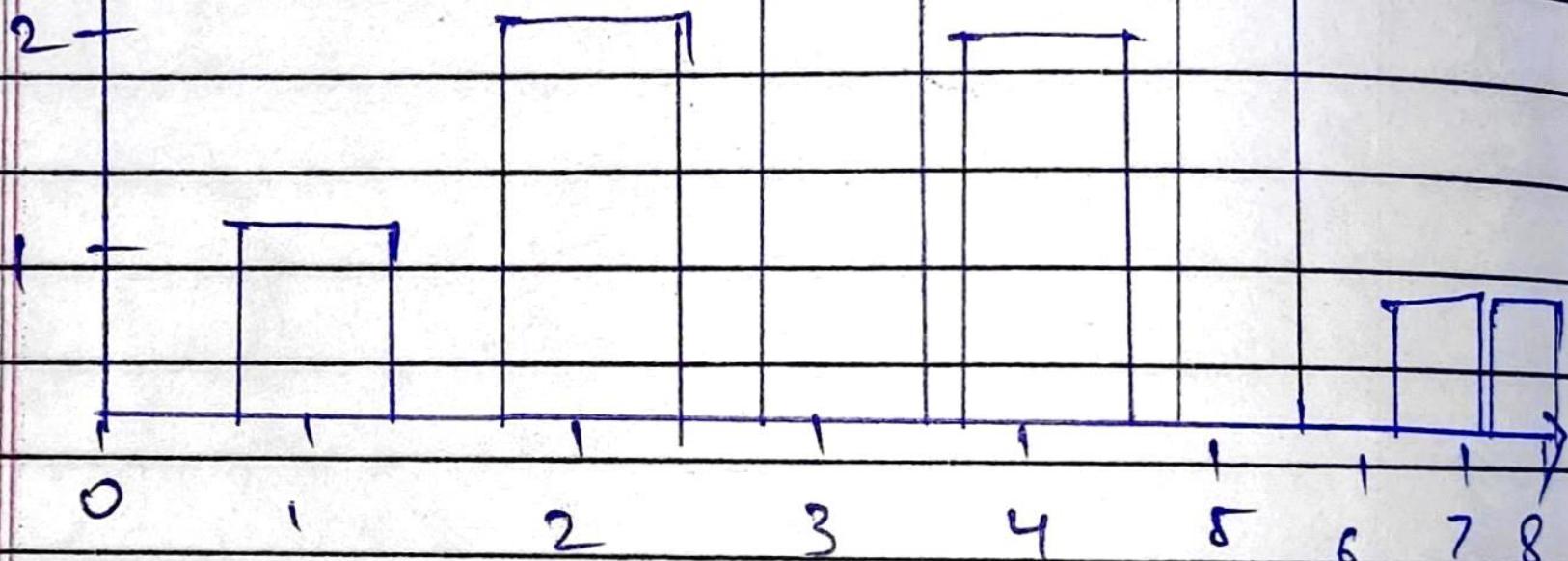
Frequency

Distribution

Smoothening

Probability density function





→ Pdf: Probability density function }

↓ Pmf: Probability mass function }

continuous

(pdf)

↓  
discrete  
(pmf)

① Mean  
 ② Median  
 ③ Mode

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

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

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

Population ( $N$ )       $N > n$       Sample ( $n$ )

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

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

$n=6$        $N > n$

$$\text{Population} = \{24, 23, 2, 128, 27\}$$

Age

$n=4$       sample Age =

~~$\bar{x} \& N$~~

$$\{24, 21, 27\}$$

$\bar{x} \& \mu$

Average

$$\text{sample mean } (\bar{x}) = \frac{24 + 2 + 1 + 27}{4}$$

$$\bar{x} = 13.5$$

$$\mu > \bar{x} \quad \mu > \bar{x}$$

$$\bar{x} > \mu \quad \bar{x} > \mu$$

`np.nan`  $\Rightarrow$  Null Values.

### Practical Application (Feature Engineering)

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

NAN  $\rightarrow$  Not a Number

Age	salary (k)
24	45
28	50
29	NAN
31	75
36	80
NAN	NAN

Average (mean)

$$\text{Age} = 29.6$$

$$\text{salary} = 62$$

$$\text{outliers} \leftarrow [80] \leftarrow [200]$$

$$(\text{Age}) \text{ mean} = 38 \quad \text{salary} = 85$$

## ② \* Median

outliers

$$\begin{aligned} \{1, 2, 3, 4, 5\} &= \{1, 2, 3, 4, 5, 100\} \\ \bar{x} = 3 &\rightarrow \bar{x} = \frac{1+2+3+4+5+100}{6} = \frac{115}{6} = 19.16 \end{aligned}$$

## Steps to Find Median:

- (1) Sort the Numbers
- (2) find the central number:-

(i) if the no. of elements are even we find the average of central elements.

(ii) if the no. of elements are odd we find the central elements.

### Examples :-

$$\{0, 1, 2, 3, 4, \boxed{5, 6}, 7, 8, \boxed{100, 120}\}$$

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

no impact in median

\* if there are no outliers we use mean

\* if there is outliers then  $\boxed{\text{median} = 5}$  we use median.

(3) Mode: { Most frequent occurring elements }

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

3

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

Dataset } categorical variable }

Types of Flowers:

Lily

[Rose] ← Most frequent elements.

sunflower      NAN = Rose.

NAN ← Rose

[Rose]

sunflower

[Rose]

NAN ← Rose

\* Measure of Dispersion:

① Variance ( $\sigma^2$ )

② standard deviation. ( $\sigma$ )

Variance

Population Variance ( $\sigma^2$ )

Example:



①

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



N

$\mu = \text{Variance (Population mean)} = \frac{3}{5}$

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

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

5

$$+ (80 - 14.5)^2 \Big]$$

7

$$\sigma^2 = \frac{719.10}{7}$$

$$\frac{(\sigma^2)}{2} < \frac{(\sigma^2)}{719.10}$$

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

$$\left\{ 1, 2, 3, 4, 5 \right\} = M = \underline{\underline{3}}$$

$$\sigma^2 = \frac{2}{7}$$

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

$$0.1 \cdot P1F = 0.1$$

No. of even Numbers

$$\frac{2}{8} = 0.25$$

Total no. of Numbers

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

Percentiles = Grade (CAT, JEE, SAT)

↓

Percentiles.

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

$$\frac{8}{20} = 45 \text{ percentile}$$

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

$$\frac{14}{20} = 70 \text{ percentile.}$$

\* what is the value that consists at 25 percentile.

$$\text{Value} = \frac{\text{Percentile}}{100} \times n$$

100

= 19.95

### \* 5 number summary:

① minimum

② First Quartile (25 percentile) (Q1)

③ Median

④ Third Quartile (75 percentile) (Q3)

⑤ Maximum

remove the outliers.



Box Plot

$\{1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6,$   
 $7, 8, 8, 9, 12\}$

Create a set:

lower fence  $\leftarrow \rightarrow$  Higher fence

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

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

IQR  $\Rightarrow$  Inter Quartile Range (IQR)

75 25

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

$$Q_1 = \frac{25 * (n+1)}{100} = \frac{25 * 21}{100} = 5.25 \Rightarrow \text{index} = 3$$

$$Q_3 = \frac{75 * 21}{100} = 18.75 \Rightarrow \text{index} = \frac{8+7}{2}$$

$$= [7.5]$$

$$\text{Lower fence} = 3 - 1.5(4) = 3 - 6 = -3$$

$$\text{Lower fence} = [-3]$$

$$\text{lower fence} = 3 - (1.5)(4.5) = -3.65$$

$\boxed{-3.65}$

$$\text{Higher fence} = 7.5 + (1.5)(4.5) = 14.25$$

$\boxed{14.25}$

lower fence  $\longleftrightarrow$  Higher fence

$3.65 \longleftrightarrow 14.25$

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

outliers  $\rightarrow$   $\boxed{9, 10}$

① Minimum = 1

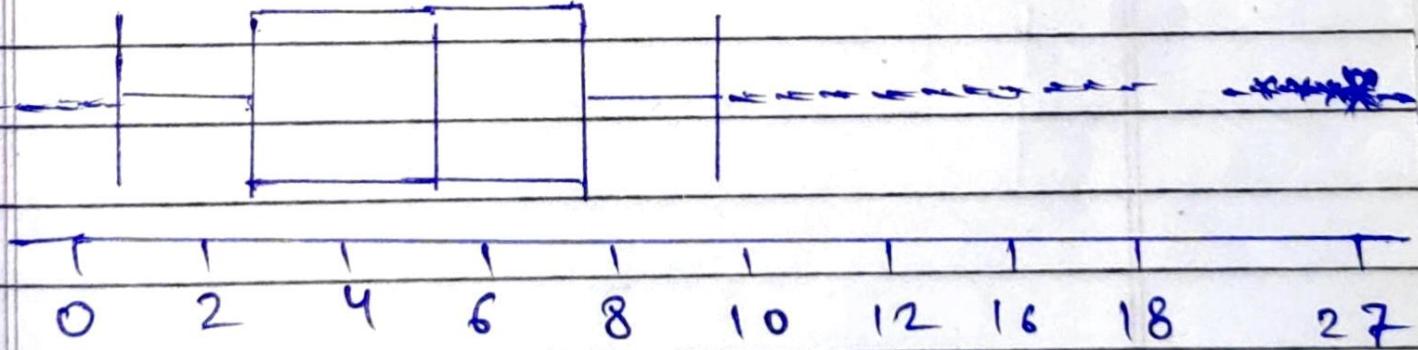
② Q<sub>1</sub> = 3

③ Median = 5

④ Q<sub>3</sub> = 7.5

⑤ Maximum = 9.

Box Plot



To treat Outliers