

Histogram

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

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

eg:- [10, 20, 25, 30, 35, 40]

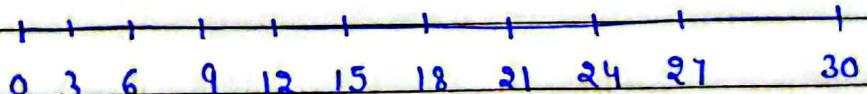
$$\text{min} = 10$$

$$\text{max} = 40$$

Now, I want to divide data into 10 groups.

$$\Rightarrow \text{bins} = 10$$

$$\rightarrow \text{bin size} = \frac{40 - 10}{10} = 3 \quad \left(\frac{\text{max} - \text{min}}{\text{bins}} \right)$$



Here, bins = 10 but the value is ranging b/w 0 to 30
It should be ranging b/w 0 to 40.

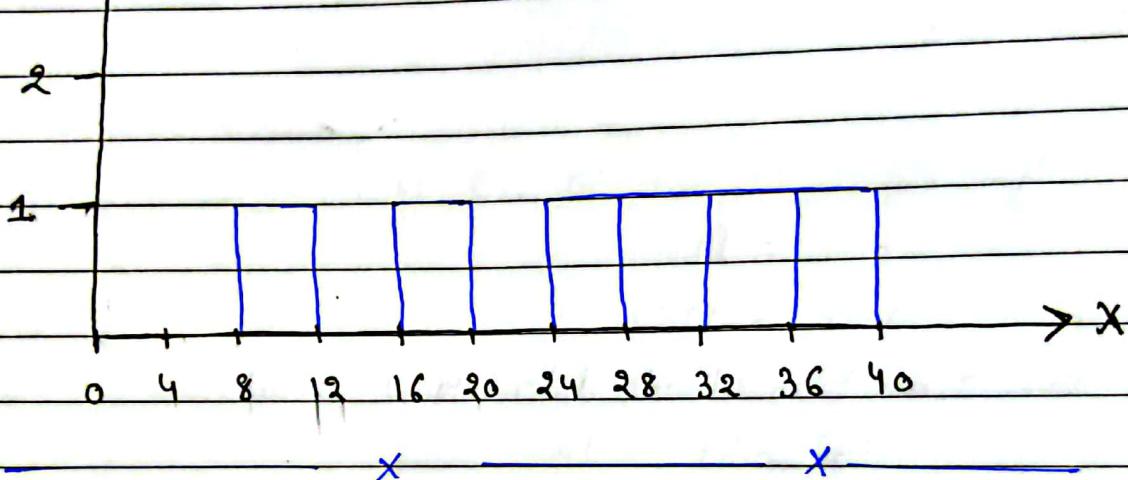
$[10, 20, 25, 30, 35, 40]$

$$\min = 10, \max = 10$$

$$\text{bins} = 10$$

$$\text{bin size} = \frac{40}{10} = 4$$

frequency / counts



now,

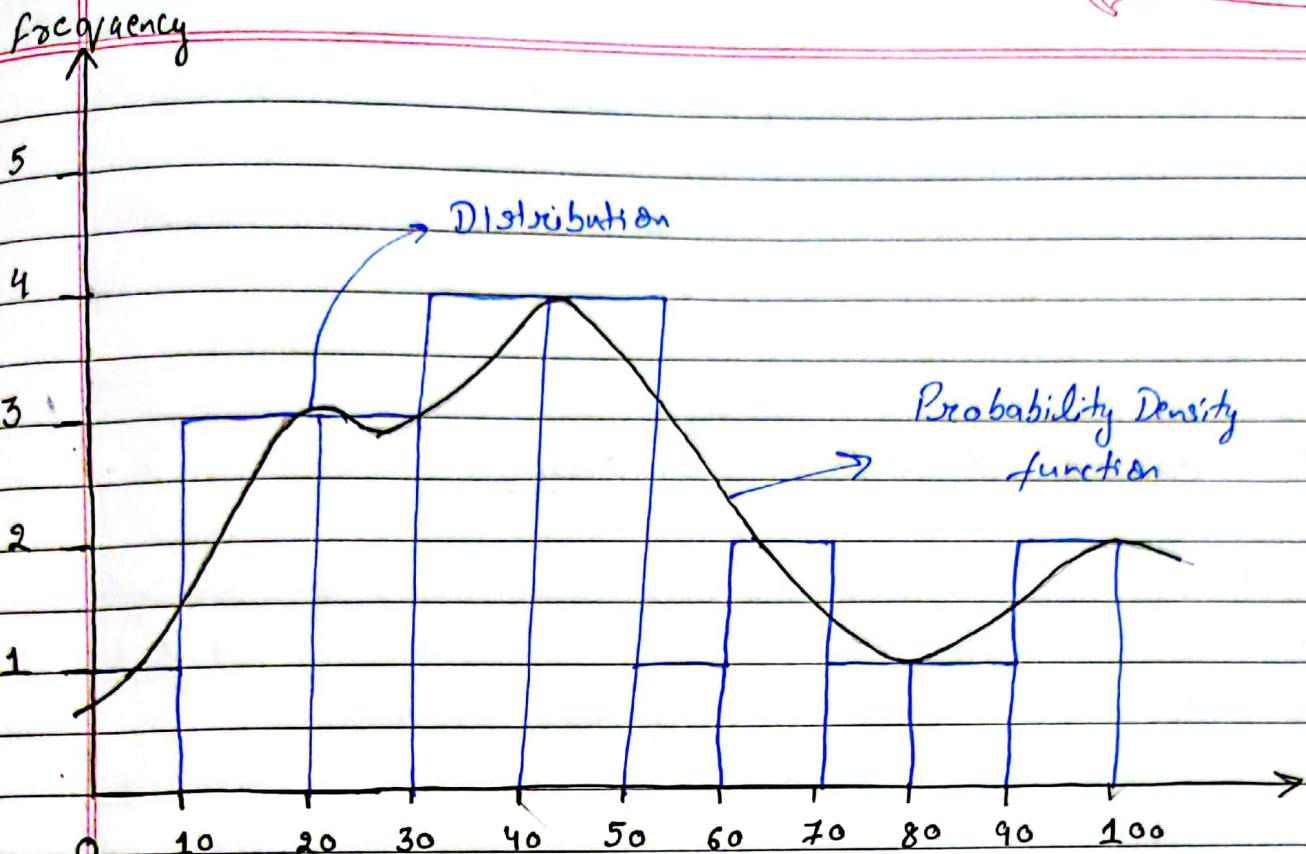
Age = $[10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100]$

$$\min = 10$$

$$\max = 100$$

$$\text{bins} = 10$$

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

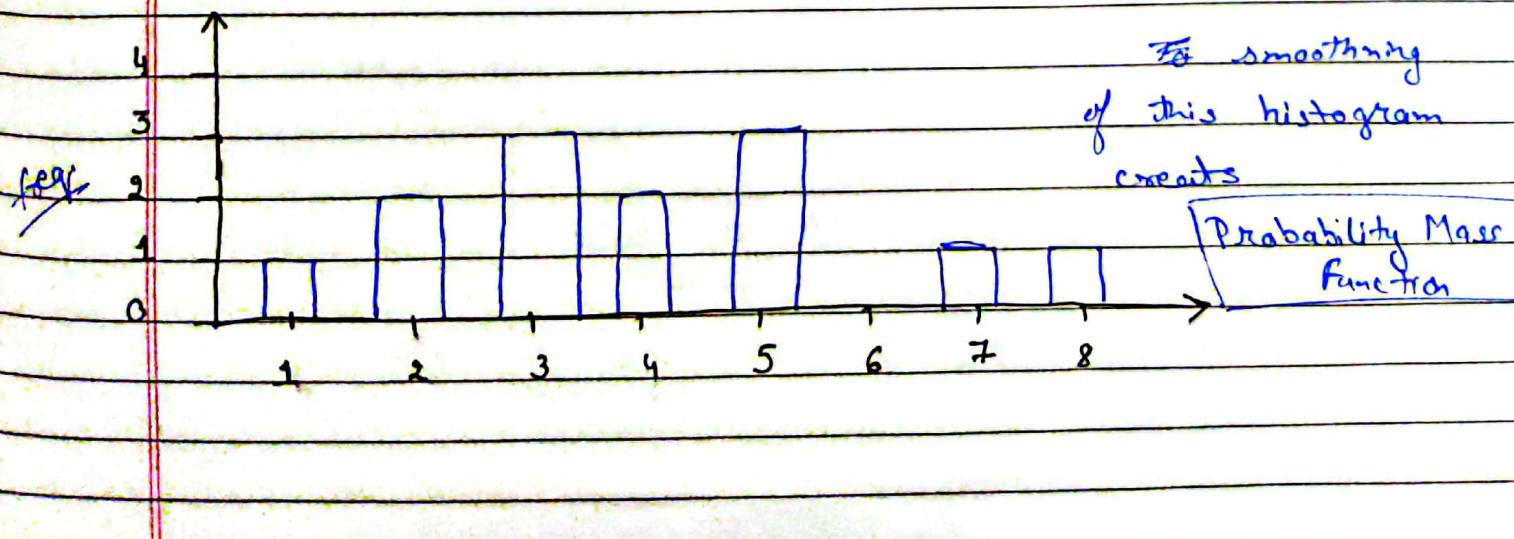


- Smoothing of Histogram creates Probability Density Function
- Through the smoothing of histogram we get distribution of data.

Measures of Central

② Discrete Random Variable

Num of bank accounts :- [2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5]



pdf: probability density function → for continuous variable

pmf: probability mass function → for discrete variable

X — X — X



Measure of Central Tendency

① Mean

② Median

③ Mode

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

① Mean → $X = \{1, 2, 3, 4, 5\}$

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

→ Mean should be defined based on 2 factors

Population (N)

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

Sample (n)

$$N \geq n$$

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

Examples

→ Population Age = [24, 23, 2, 1, 28, 27]
 ↓
 $N = 6$

$$\text{Population mean } (\mu) = \frac{24+23+2+1+28+27}{6} = 17.5$$

Sample Age = [24, 2, 1, 27]
 ↓
 $n = 4$

$$\text{Sample mean } (\bar{x}) = \frac{24+27+1+27}{4} = 13.5$$

* ① $N \geq n$

② $\boxed{\begin{array}{l} \mu > \bar{x} \\ \text{or} \\ \bar{x} > \mu \end{array}}$

→ Practical Application (Feature Engineering)

Age	Salary	Family Size
-	-	-
-	-	-
NaN	-	-
-	NaN	-
-	-	NaN

If we drop whole row with NaN
 there will be loss of info

So we can find ~~mean~~ mean of
 each column and replace
 NaN values with that

② Median :-

eg:- $[1, 2, 3, 4, 5]$

$$\downarrow$$

$$\bar{x} = 3$$

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



$$\bar{x} = 19.16$$

outlier

Steps to find out median

① Sort the Numbers

② Find the central number given two conditions.

③ If num of elements are even we find average of central elements.

④ If num of elements are odd we find ~~average~~ central element.

eg:-

$[1, 2, 3, 4, 5, 6, 7, 8, \underbrace{100}, \underbrace{200}]$



$$\text{median} = \frac{5+6}{2} = \frac{11}{2} = 5.5$$

outliers

$[1, 2, 3, 4, 5, 6, 7, 8, 100]$



$$\text{median} = 5$$

③ Mode :- Most frequent occurring element

$$[1, 2, 2, 3, 3, 3, 4, 4, 5] \rightarrow [3]$$

$$[1, 2, 2, 2, 3, 3, 3, 4, 4, 5] \rightarrow [2, 3]$$

Practical Example

Types of Flowers

Lily
Sunflower
Rose
NaN
Rose
Sunflower
Rose
NaN



Here we can use mode value

* Measure of Dispersion

① Variance (σ^2) → spread of data

② Standard Deviation (σ)

Variance

Population Variance (σ^2)

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

Sample Variance (S^2)

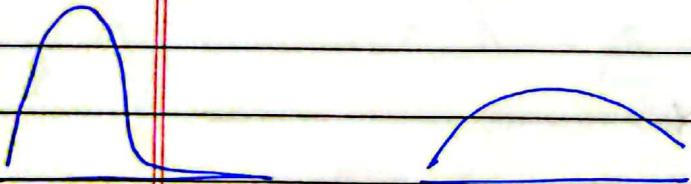
$$S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

$x_i - \mu \rightarrow$ Distance from Mean

① $[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]$

② $[1, 2, 3, 4, 50, 60, 70, 100]$

Variance 1 < Variance 2



As Variance \uparrow Spread \uparrow

Q(2) Standard Deviation (σ)

Standard Deviation ($\sqrt{\sigma^2}$)

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

$$\boxed{\sigma^2 = 2.0, \mu = 3}$$

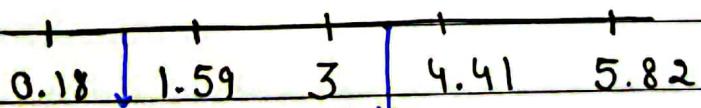
$$\mu = 3$$

$$\sigma^2 = 2$$

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

2 std. to left $\leftarrow -2$ $\rightarrow +2$ \rightarrow 2 std. Δ to the right

1 std. to left $\leftarrow -1$ $\rightarrow +1$ \rightarrow 1 std. to the right



(1 falls in 2nd std.dev to left) \rightarrow (so 4 falls in 1st std.dev to right)

std. ~~dev~~ tells that how many std. away a number falls away
the mean.

*

Percentiles and Quartiles

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

$$\% \text{ of even nos} = \frac{\text{no. of even numbers}}{\text{Total elements}} \times 100$$

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

Percentiles :- A percentile is a value below which a certain percentage of observations lie.

eg:-

99 percentile \rightarrow It means the person has got better marks than 99% of entire students

Ascending order

Dataset :- $\underbrace{2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12}_{1 \quad \quad \quad 16 \quad \quad \quad 20}$

a) What is the percentile rank of 10

$$\text{Percentile Rank of } x = \frac{\text{no. of values below } x}{\text{sample size } (n)}$$

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

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

\downarrow
4.5% of entire values
is less than 8

Q. What is the value that exist at 25 percentile?

$$\boxed{\text{Value} = \frac{\text{Percentile} \times n+1}{100}} \rightarrow \text{for odd}$$

$$\rightarrow \text{Value} = \frac{25 \times 20}{100}$$

= 5th index



5 will be the output

now 95 percentile value

$$\boxed{\cancel{\text{Value} = \frac{95 \times 20}{100}}} =$$

$$\text{value} = \frac{95 \times 20}{100} = 19.95 \text{ index} \\ = 12 \rightarrow \text{output}$$

* 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

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



outlier

Here, we will try to create a fence
and value should range b/w lower fence and
higher fence value.

[Lower fence \longleftrightarrow Higher fence]

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

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

\downarrow
Inter Quartile Range

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

\downarrow
1.5 std dev. to the mean

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

$$= \frac{25}{100} \times 21 = 5.25 \text{ index}$$

Take average of 5" & 6" index value

$$\text{Output} = [3]$$

$$Q_3 = \frac{75}{100} \times 21 = 15.75 = \frac{8+7}{2} = [7.5]$$

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

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

[1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27]

$$\textcircled{1} \quad \text{minimum} = 1$$

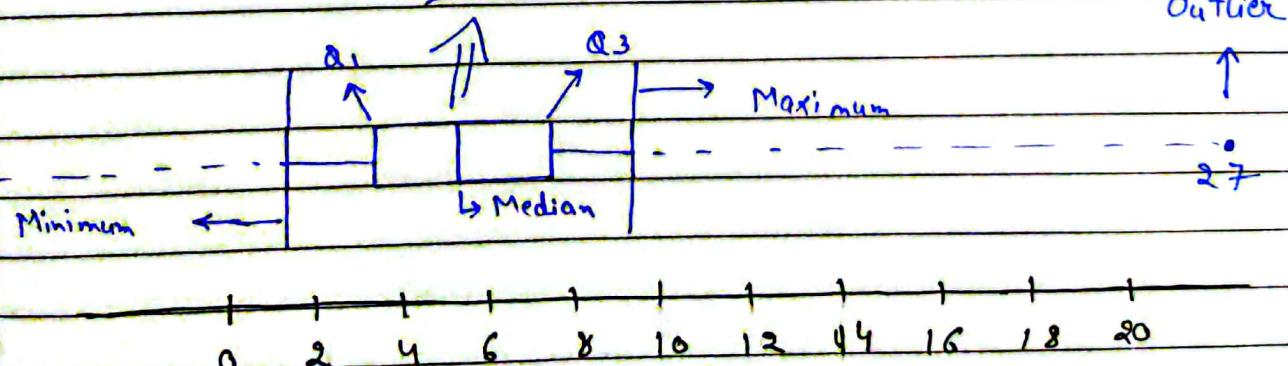
$$\textcircled{4} \quad Q_3 = 7.5$$

$$\textcircled{2} \quad Q_1 = 3$$

$$\textcircled{5} \quad \text{Maximum} = 9$$

$$\textcircled{3} \quad \text{Median} = 5$$

Box Plot



To Treat Outliers