# Statistics-2

## ① Percentiles & Quartiles

**Percentage:-**   1, 2, 3, 4, 5, 6

% of numbers that are odd $= \dfrac{3}{6} = \dfrac{\text{No of odd numbers}}{\text{Total no. of numbers}}$

$= \dfrac{1}{2} = 50\%$

### Percentiles:-

A Percentile is a value below which a certain percentage of data points lie.

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

> (n = total no. of values.)

Percentile Rank of 10 $= \dfrac{\#\text{ of values below 10}}{n} \times 100$

$= \dfrac{\overset{4}{\cancel{12}}}{\underset{\underset{1}{5}}{\cancel{15}}} \times \overset{20}{\cancel{100}} = 80$ Percentile.

80 Percentile = 80% of the distribution fall below the value 10.

Ⓧ what value exists at 25 Percentile?

$$\boxed{\text{value} = \dfrac{\text{Percentile}}{100} \times (n+1)}$$

$$= \frac{28}{100} \times \cancel{16}^{7}$$

$$x_1$$

$$\boxed{= 4^{th} \text{ element}}$$

$$x = \{ \underset{\downarrow\downarrow\;\downarrow\downarrow}{2,3,3,4,6,6,6,7,8,8,9,9,} \underline{10},11,12 \}$$

\* If we get value as decimal
like 4.5
we can take average of 4th and
5th value.

## Quartiles

$Q_1 \rightarrow$ 25 Percentile

$Q_2 \rightarrow$ median $\rightarrow$ 50 Percentile

$Q_3 \rightarrow$ 75 Percentile.

② <u>5 Number Summary</u>

i, minimum

ii, First Quartile (25 Percentile) $(Q_1)$

iii, Median $(Q_2)$

iv, Third Quartile (75 Percentile)

$(Q_3)$

v, Maximum

<u>Remove the outliers</u>

$$X = \{ 1,2,2,2,3,3,4,5,5,5,6,6,6,6,7,8,8,9,\underline{29} \}$$

[ lower Fence $\longleftrightarrow$ upper fence]

lower fence = $Q_1 - 1.5(IQR)$

Higher fence = $Q_3 + 1.5(IQR)$

Inter Quartile Range

$$= \boxed{Q_3 - Q_1}$$

$$X = \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 29\}$$

$Q_1 = 25$ Percentile $= \dfrac{25}{100} \times (20) = 5^{th}$ value $\boxed{3}$

$Q_3 = 75$ Percentile $= \dfrac{75}{100} \times 20 = 15^{th}$ value $= \boxed{7}$

$$IQR = Q_3 - Q_1$$

$$= 7 - 3$$

$$\boxed{= 4}$$

lower fence $= Q_1 - 1.5(IQR)$

$$= 3 - 1.5(4)$$

$$= 3 - 6$$

$$\boxed{= -3}$$

Higher Fence $= Q_3 + 1.5(IQR)$

$$= 7 + 1.5(4)$$

$$= 7 + 6$$

$$\boxed{= 13}$$

$$[-3, 13]$$

Hence, we can consider $\boxed{29}$ as an outlier for the above data.

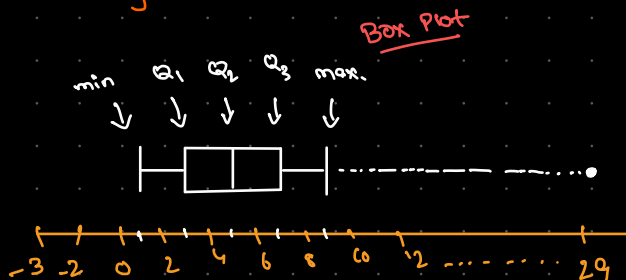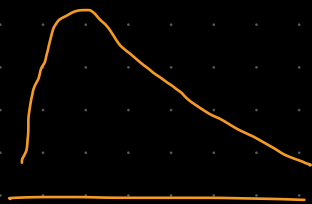③ Box Plot = [to visualize outliers]
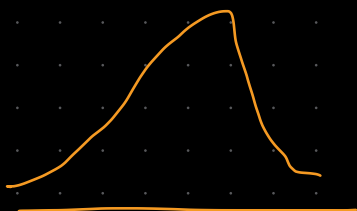
i, minimum value = 1
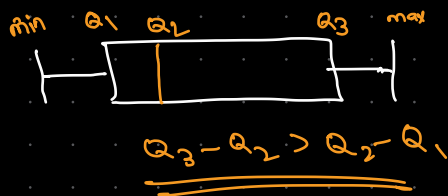
ii, $Q_1 = 3$

iii, median $= Q_2 = 5$

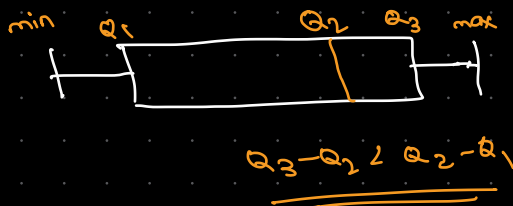iv, $Q_3 = 7$

v, maximum value = 9

Box Plot

min  $Q_1$  $Q_2$  $Q_3$  max.
↓    ↓    ↓    ↓    ↓

-3 -2  0  2  4  6  8  10  12 ..... 29

mean > median > mode

$$Q_3 - Q_2 > Q_2 - Q_1$$



mode > median > mean

$$Q_3 - Q_2 < Q_2 - Q_1$$

Assignment :

$$y = \{-13, -12, -6, -5, 3, 4, 5, 6, 7, 7, 8, 10, 10, 11, 24, 55\}$$

$$Q_1 \ (25) =$$

$$— = \frac{25}{100} \times (16 + 1)$$

$$\boxed{Q_1 = -1}$$

$$= \frac{\overset{1}{\cancel{25}}}{\underset{4}{\cancel{100}}} \times 17 \qquad Q_3 = \frac{\overset{3}{\cancel{25}}}{\underset{4}{\cancel{100}}} \times (17)$$

$$= 4.25 \qquad\qquad\qquad = 12.75 \ \text{element}$$
$$\underline{\text{element}}$$

$$\underline{Q_2 \ (50)} \qquad \frac{6+7}{2} = \frac{13}{2}$$

$$\boxed{Q_2 = 6.5} \qquad\qquad = 6.5 \qquad\qquad \boxed{Q_3 = 10}$$

$$IQR = Q_3 - Q_1$$
$$= 10 - (-1) = \boxed{11}$$

Higher Fence =

Lower Fence
$$= Q_1 - 1.5(IQR)$$
$$= -1 - 1.5(11)$$
$$= -1 - 16.5$$

$$Q_3 + 1.5(IQR)$$
$$= 10 + 16.5$$
$$\boxed{= 26.5}$$

$$= \boxed{-17.5}$$

$$[-17.5, \ 26.5]$$

## Boxplot

① min = -13  ② $Q_1 = -1$  ③ $Q_2 = 6.5$

④ $Q_3 = 10$  ⑤ max = 24

$$\boxed{\text{Outlier} = 55}$$



min        Q1        Q2   Q3              max

-14  -11  -8  -5  -2  1  3  6  9  12  15  18  21  24  27

②

$Z = \{1, 2, 4, 6, 7, 12, 18, 34, 77, 66, 108, 99, 14\}$

Sorted $= \{1, 2, 4, 6, 7, 12, 14, 18, 34, 66, 77, 99, 108\}$

$Q_1 = \dfrac{25}{100} \times 14 = 3.5$ element

$$\boxed{Q1 = 5}$$

$$\boxed{Q_2 = 14}$$

$Q_3 = \dfrac{\cancel{75}^{3}}{\cancel{100}} \times \cancel{14}7 = 10.5$ element

$$\underset{2}{14}$$

$$\boxed{Q_3 = 71.5}$$

$$IQR = Q_3 - Q_1$$

$$\boxed{= 66.5}$$

lower fence $= Q_1 - 1.5(IQR)$

$= 5 - 1.5(66.5)$

$= \underline{\underline{-94.75}}$

High fence $= Q_3 + 1.5(IQR)$

$= 71.5 + 99.75$

$= \underline{171.25}$

$$[-94.75, \ 171.25]$$

# ③ Covariance and Correlation

| X | Y |
|---|---|
| 2 | 3 |
| 4 | 5 |
| 6 | 2 |
| 8 | 9 |

[Relationship between X and Y]

$X \uparrow \quad Y \uparrow$ ✓

$X \uparrow \quad Y \downarrow$    Size of   ↓ house   Price

$X \downarrow \quad Y \uparrow$

$X \downarrow \quad Y \downarrow$ ✓



$X \uparrow \quad Y \uparrow$
$X \downarrow \quad Y \downarrow$



$X \uparrow \quad Y \downarrow$
$X \downarrow \quad Y \uparrow$

## Covariance

$$cov(x,y) = \sum_{i=1}^{n} \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

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

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

⇓

$$\boxed{var(x) \Longleftarrow cov(x,x)}$$

$cov(x,y)$

| $x \uparrow \, y \uparrow$ |
| $x \downarrow \, y \downarrow$ |

+ve covariance

| $x \downarrow \, y \uparrow$ |
| $x \uparrow \, y \downarrow$ |

−ve covariance

**Ex:**

| $x$ | $y$ |
|-----|-----|
| 2 | 3 |
| 4 | 5 |
| 6 | 7 |
| $\bar{x}=4$ | $\bar{y}=5$ |

X and Y are having a positive covariance

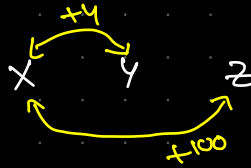$$cov(x,y) = \sum_{i=1}^{n} \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$= \frac{[(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)]}{2}$$

$$= \frac{4 + 0 + 4}{2} = \boxed{4} \quad \text{+ve covariance}$$

<u>Advantages</u>

① Relationship between
    x and Y

+Y

x ⟶ Y     Z

+100

<u>Disadvantages</u>

1) covariance doesnot have a
   specific limit value

$-\infty \longleftrightarrow +\infty$

(It will not gives the strength
   of the relation between
   variables)

(It only gives the direction)

---

**\* Pearson Correlation coefficient** $\left[-1 \text{ to } +1\right]$

$-1 \text{ to } 1$

$\left(\begin{array}{l}\text{row of } x \\ \text{and } y\end{array}\right) \longleftarrow \cdots \quad \int_{x,y} = \dfrac{cov(x,y)}{\sigma_x \cdot \sigma_y} \rightarrow (\text{std deviation of} \atop x \text{ and } Y)$

x ⟶ y     z

$-1 \text{ to } +1$

\* The more the value towards +1, the more +ve correlated it is.
\* The more the value towards -1, the more -ve correlated it is.

|   |   |     |          |
|---|---|-----|----------|
| x | y | 0.4 |          |
| x | z | 0.7 | → strong |

<u>Disadv:</u>
    only linear relationship will be captured by Pearsons
                       correlation coefficient.

---

**\* Spearman Rank correlation**

$r_s = \dfrac{cov(R(x), R(y))}{\sigma_{R(x)} * \sigma_{R(y)}}$

| x | y | R(x) | R(y) |
|---|---|------|------|
| 5 | 6 | 3    | 1    |
| 7 | 4 | 2    | 2    |
| 8 | 3 | 1    | 3    |
| 1 | 1 | 5    | 4    |
| 2 | 2 | 4    | 5    |

## Feature Selection

| +ve | +ve | +ve | Uo impact | -ve | O/P : |
|---|---|---|---|---|---|
| Size of house | No. of rooms | Location | No. of People Staying | is Haunted | Price. |

---

* Probability Distribution Function

* Probability Density Function

* Probability Mass Function.

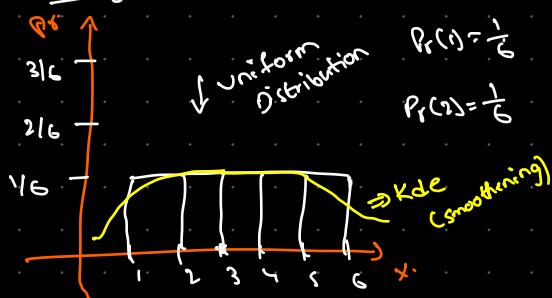→) **Probability Distribution Function**

    i) Probability Density Fn

    ii) Probability Mass Fn
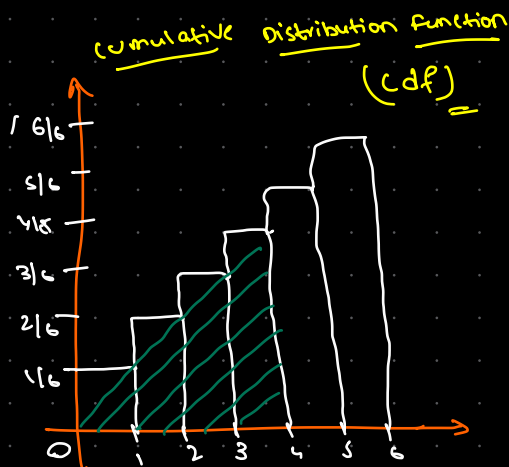
    iii) Cumulative Distributive Fn.

① **Probability Mass Function (PMF)**

    ✗ when dataset has Discrete Random Variable

Ex: Rolling a Dice $\{1,2,3,4,5,6\}$



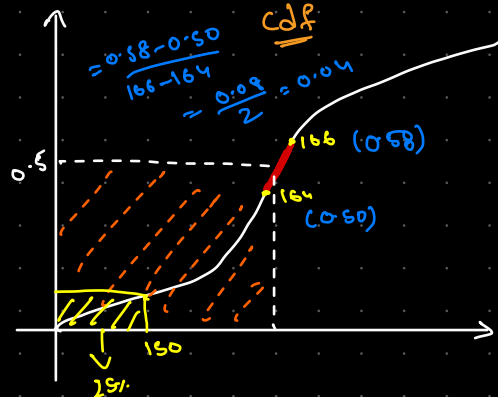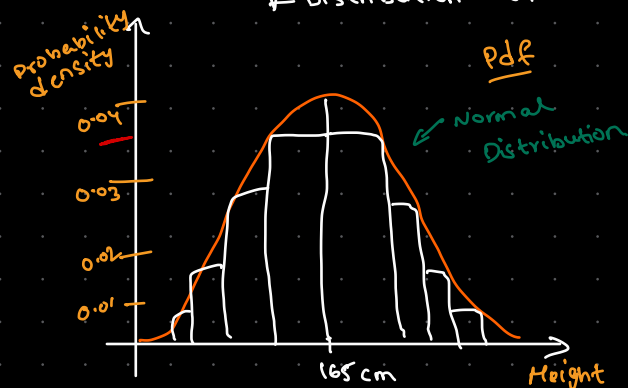↓ uniform Distribution

$P_r(1) = \frac{1}{6}$

$P_r(2) = \frac{1}{6}$

⇒ kde (smoothening)

$P_r(1 \text{ or } 2) = \frac{1}{6} + \frac{1}{6} = \frac{2}{6}$

$= \frac{1}{3}$

**Cumulative Distribution Function (cdf)**



$P_r(x \leq 4) = P_r(x=1) + P_r(x=2)$

$+ P_r(x=3) + P_r(x=4)$

$= \frac{4}{6}.$

## ② Probability Density Function

*Distribution of continous Random Variable.*



**Pdf**

← Normal Distribution

Probability density

0·04

0·03

0·02

0·01

165 cm          Height



**Cdf**

$= \dfrac{0.58 - 0.50}{166 - 164}$

$= \dfrac{0.08}{2} = 0.04$

166 (0.58)

164 (0.50)

0.5

150

25%

* During the reduce in the size of histogram bins, the slope of the curve gets reduced

* For every probability density that we are calculating in the PDF is coming from the gradient of the CDF. *
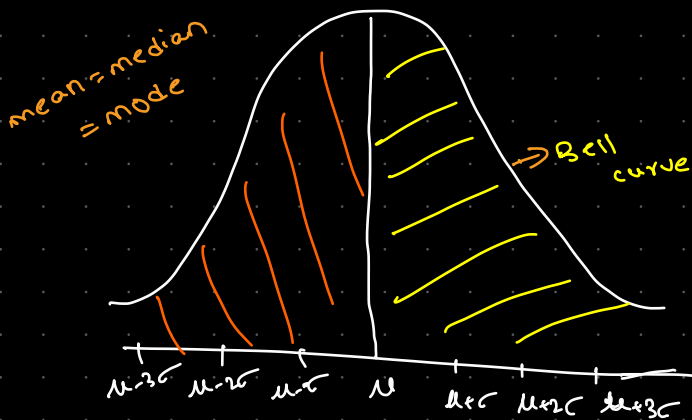
## * Probability Density ⇒ Gradient of cumulative curve.



Probability Density

0.02



0·02

(Getting derivative of gradient of 0.02)

# Different types of Distribution

1. Normal Distribution → PDF
2. Standard Normal Distribution → PDF
3. Log Normal Distribution → PDF
4. Power Law Distribution → PDF
5. Bernoulli Distribution → PMF
6. Binomial Distribution → PMF
7. Poisson Distribution → PMF
8. Uniform Distribution { Discrete → PDF / continous → PMF
9. Exponential Distribution → PDF
10. CHI SQUARE Distribution → PDF
11. F. Distribution → PDF

## 1 Normal / Gaussian Distribution



mean = median = mode

→ Bell curve

$x$ : continous Random Variable

$\mu-3\sigma$  $\mu-2\sigma$  $\mu-\sigma$  $\mu$  $\mu+\sigma$  $\mu+2\sigma$  $\mu+3\sigma$
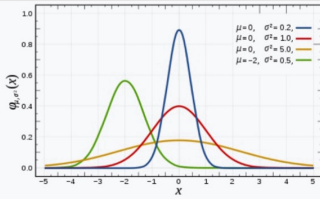
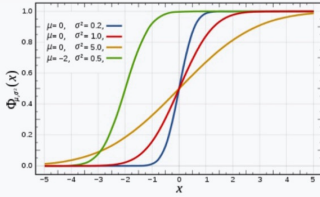Ex: Height, weight, age ..    (IRIS Dataset)

## Normal distribution

### Probability density function



The red curve is the *standard normal distribution*

### Cumulative distribution function
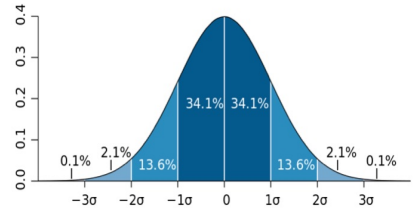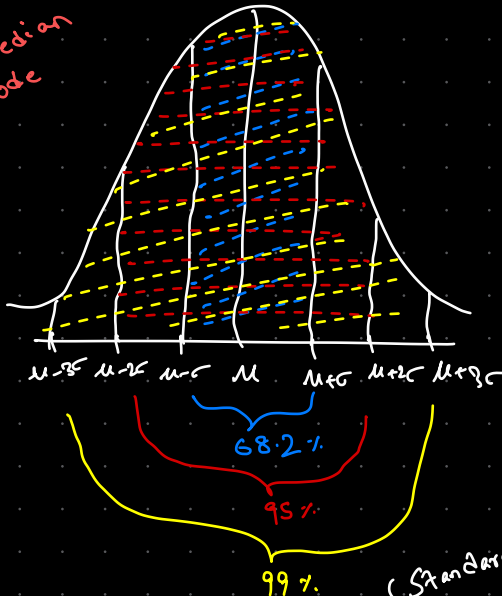


| Notation | $\mathcal{N}(\mu, \sigma^2)$ |
|---|---|

$$X \underset{\sim}{\sim} N(\mu, \sigma^2)$$

Support Parameters     $\mu = $ mean

$\sigma^2 = $ variance

$$PDF = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

## Empirical Rule   (68 - 95 - 99.7% Rule)

(3 - Sigma Rule)

mean = median = mode



$\mu - 3\sigma$   $\mu - 2\sigma$   $\mu - \sigma$   $\mu$   $\mu + \sigma$   $\mu + 2\sigma$   $\mu + 3\sigma$
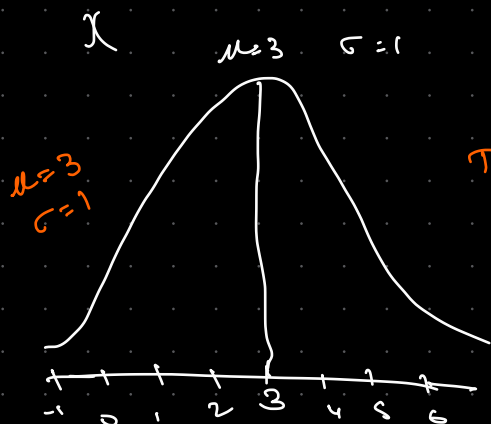
68.2%

95%

99%



For the normal distribution, the values less than one standard deviation away from the mean account for 68.27% of the set; while two standard deviations from the mean account for 95.45%; and three standard deviations account for 99.73%.
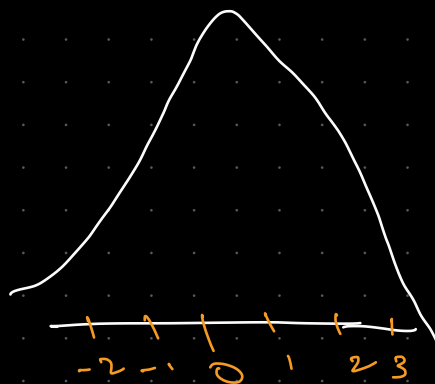
(Standard Deviation and coverage)

# ② Standard Normal Distribution.

The simplest case of normal distribution is known as standard normal distribution or unit normal distribution

This is special case when $\mu = 0$ and $\sigma = 1$

$x$

$\mu = 3 \quad \sigma = 1$

$\mu = 3$
$\sigma = 1$

Transformation
$\Rightarrow$
$\mu = 0, \sigma = 1$

-1  0  1  2  3  4  5  6

-2 -1  0  1  2  3

(Standard normal distribution)

$\downarrow$

$$Z\text{-score} = \frac{x_i - \mu}{\sigma}$$
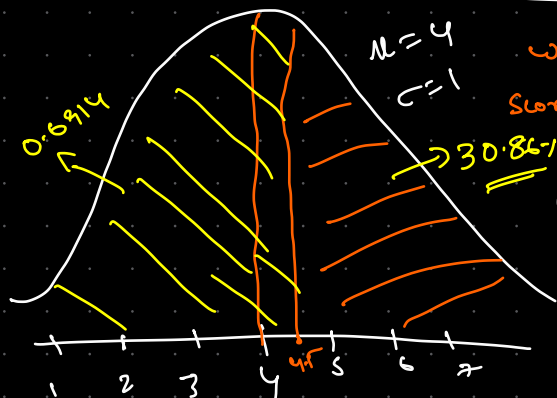
Z-score tells you about a value how many standard deviation away from mean

$\Rightarrow \dfrac{1-3}{1} = -2$

$\dfrac{3-3}{1} = 0 \qquad \dfrac{2-3}{1} = -1$

$\boxed{4.5} \Rightarrow \dfrac{4.5-3}{1} = 1.5$

$x$

0.6914

$\mu = 4$
$\sigma = 1$

$\Rightarrow 30.86\%$

what is the percentage of scores lies above 4.5?

(we need to use Z-table to calculate the area).

1  2  3  4 4.5 5  6  7

$$Z - score = \frac{4.5 - 4}{1}$$

$$= 0.5$$

Area under curve = $1 - 0.6914$

$$= 0.3086 \implies 30.86\%$$