

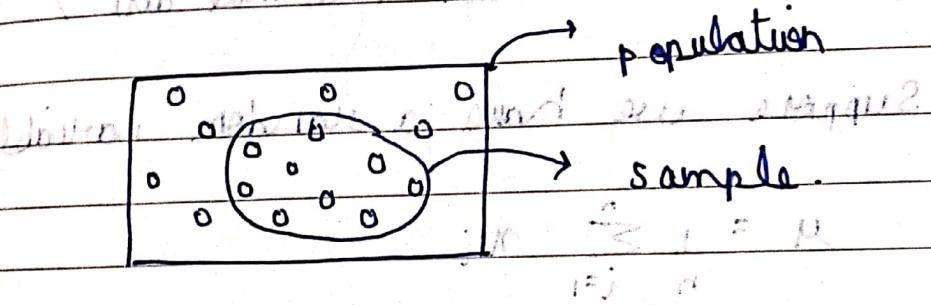
~~Population~~ Consider we want to calculate the average height of all the people in a state.

Total no. of people in state is said to be the population.

$$\text{Population Mean} = \frac{1}{N} \sum_{i=1}^n x_i$$

It refers to the total number of instances present in the data is called its population.

(Actual Average) Population mean



Sample: It refers to the subset of examples present in population.

e.g. Exit polls in an election: The news reports based on a sample of the entire population to determine who will win the election.

This is used to generalize and make an actual prediction on who will win the election.

~~$$\text{Sample Mean} = \frac{1}{n} \sum_{i=1}^n x_i$$~~

~~Population Mean > Sample mean~~

1. Random Variables
- Discrete \rightarrow whole number, not floating
 - Continuous \rightarrow eg. - No. of bank accounts a person has.

Within a range of values, we can have any value.

Within a range of values, we can have any value.

e.g. (10 - 15)

10, 11, 12, 13, 14, 15 $\in (10, 15) \subset \mathbb{R}$

\hookrightarrow decimal, whole number

Height of a person (say 6 ft 2 inch or 6.2 feet)

- Gaussian Distribution (Normal distⁿ)

Suppose we have a random variable x .

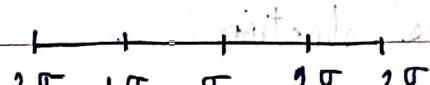
$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$\sigma^2 = \sqrt{\text{Variance}}$, standard deviation

How far is the element distributed

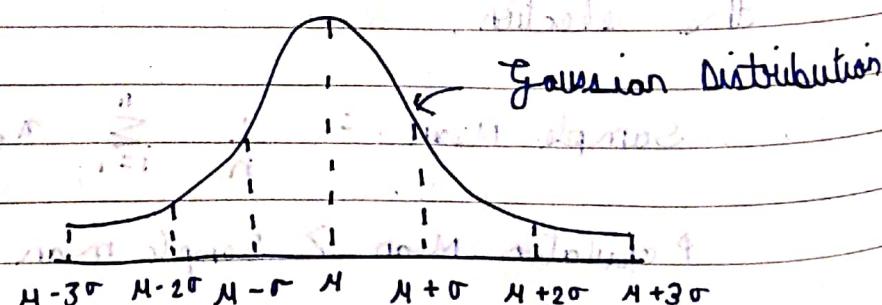
From mean μ until at distance σ



Now how will we get how it is distributed

mean after no. of standard deviation

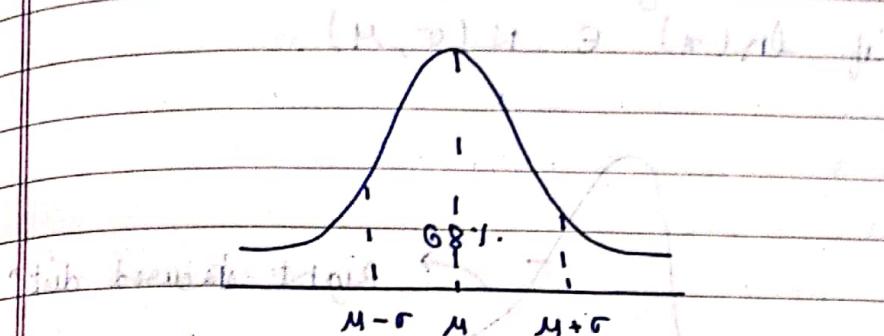
standard deviation



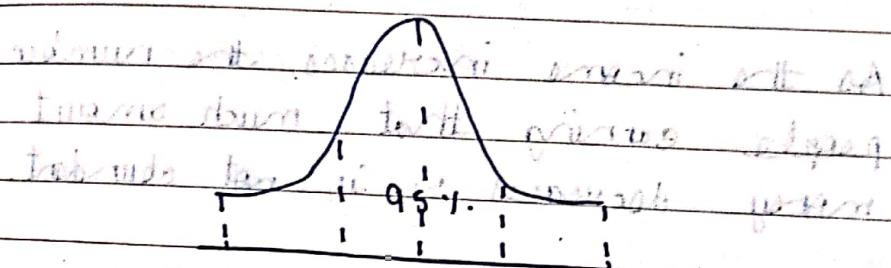
Empirical Formula:

$$P(\mu - \sigma \leq n \leq \mu + \sigma) \approx 68\%$$

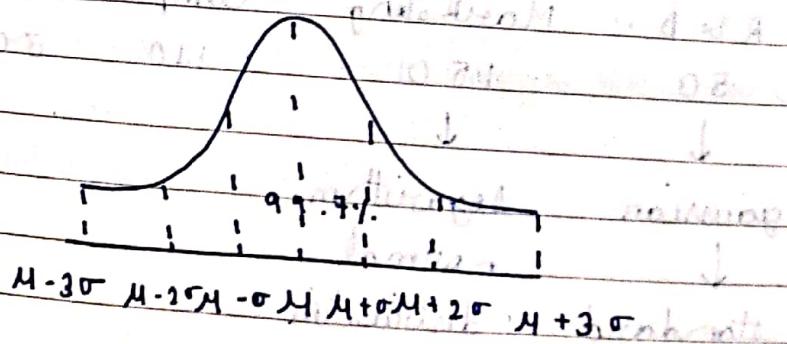
$n \rightarrow X$: Normal dist.



$$P(\mu - 2\sigma \leq n \leq \mu + 2\sigma) \approx 95\%$$



$$P(\mu - 3\sigma \leq n \leq \mu + 3\sigma) \approx 99.7\%$$

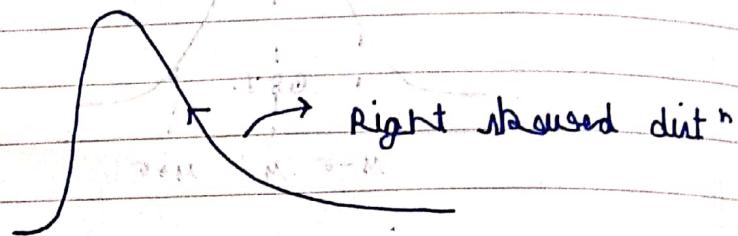


Log Distribution: It is also a form of normal distribution.

$X \sim \text{Log Normal dist}$ if $\ln(x)$ is normally distributed.

$\ln(x_1), \ln(x_2), \ln(x_3)$ follows
normal distⁿ

$\therefore x \sim \text{log normal}$
if $\ln(x) \in N(\mu, \sigma^2)$.



e.g. Income of the people.

As the income increases, the number of people earning that much amount of money decreases & is not abundant.

Feedback reviews are also in form of log distribution, then the number of people reviewing with long reviews will be very less.

e.g. Marketing Camp statistics Profit

50	150	210	300	100
↓	↓			

gaussian logarithm
↓ ↓

standard distribution

normal ↓

($\sigma = 1, \mu = 0$) standard

\therefore We scaled it normal

down by : (because $\log(\text{value})$ follows Gaussian distribution).

$$\bar{x} = \frac{x - \mu}{\sigma}$$

comparable

Covariance in Statistics.

Size Price }

1200 sgm

1800 sgm

3

100 gm

200

Can we quantify a relation between them?

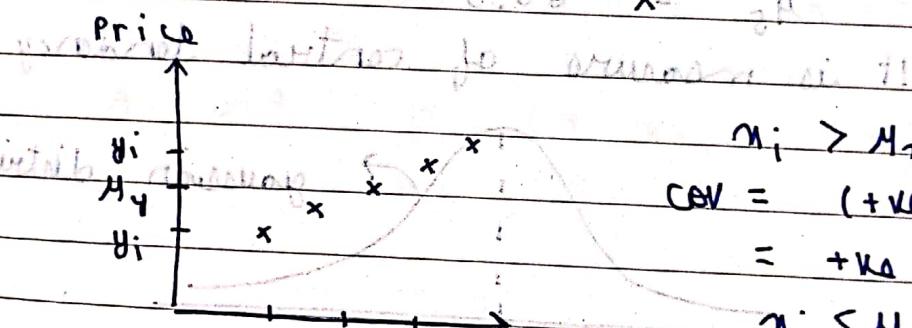
$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n [(x_i - \bar{x}_x) * (y_i - \bar{y}_y)]$$

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

$$\therefore \text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_x)(y_i - \bar{y}_y)$$

$$\therefore \text{Cov}(x, y) = \text{Variance}(x)$$

If $x \uparrow y \uparrow$ or $x \uparrow y \downarrow$



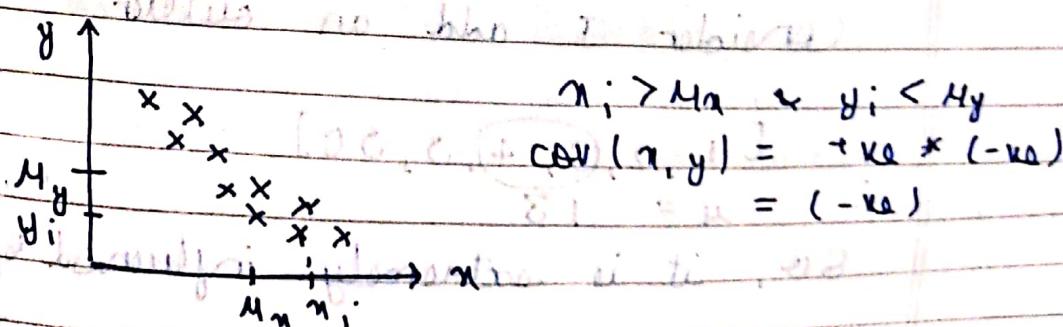
$$\begin{aligned} \text{Cov} &= (+ve) * (+ve) \\ &= +ve. \end{aligned}$$

$$\begin{aligned} x_i &< M_x \\ \text{Cov} &= (-ve) * (-ve) \\ &= +ve \end{aligned}$$

But,

Covariance will not always be positive.

Consider $x \uparrow, y \downarrow$:



$$x_i > M_x \text{ and } y_i < M_y$$

$$\begin{aligned} \text{Cov}(x, y) &= +ve * (-ve) \\ &= -ve \end{aligned}$$

if $x_i < M_x$ and $y_i > M_y$

$$\text{Cov}(x, y) = -ve + (+ve) = +ve$$

Covariance $X \uparrow, Y \uparrow \Rightarrow +va$ (how much?)
 $X \downarrow, Y \uparrow \Rightarrow -va$

We need to measure how much we read it.
 wanted outcome \rightarrow less variance

Mean, Median and Mode.

Sample of Height: f 168, 170, 150, 160, 182,

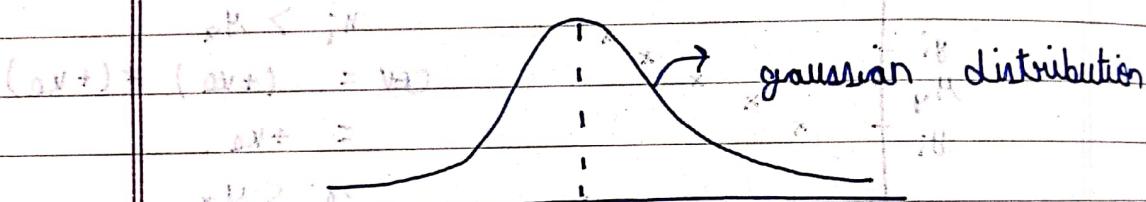
f(168+182) = 142, 175 f(150)

$$\text{Mean} (\mu_s) = \frac{1}{n} \sum_{i=1}^n [m_i] * \frac{1}{h}$$

$$\mu_s = \frac{168 + 170 + 150 + 160 + 182}{5}$$

$$\mu_s = 163.5$$

It is measure of central tendency



$$\text{Mean} = \frac{(m_1) + (m_2) + (m_3) + (m_4) + (m_5)}{5} = 163.5$$

$$\text{with } [1, 2, 3, 4, 5], M = \frac{5(5-1)}{2} = \frac{5 \times 4^2}{8 \times 5} = 3$$

Consider I add an outlier

$$[1, 2, 3, 4, 50]$$

$$\text{Median} = 13$$

So, it is extremely influenced by outliers

$$\text{Median} = \frac{3+4}{2} = 3.5$$

$$\text{Mean} = \frac{(m_1) + (m_2) + (m_3) + (m_4) + (m_5)}{5} = 163.5$$

It lessens impact of outlier

Mode is the most frequent occurring datapoint

Age : 23, 24, 27, 32, 35, 21

How to handle missing value in age?
Mean, Median or Mode?

If there are not many outliers then we can use mean or else use median.

Central Limit Theorem.

Consider a random variable X which may or may not belong to Gaussian dist".

Then probabilties for \bar{x} given X

$$X \sim \text{GD}(\mu, \sigma^2) \quad n \geq 30.$$

$$\bar{X} \sim \text{GD}(\mu, \sigma^2)$$

$$\mu = \sigma^2 \quad S_1 = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_{30}) \Rightarrow \bar{\bar{x}}_1$$

$$\mu = \sigma^2 \quad S_2 = (\bar{x}_{31}, \dots, \bar{x}_{60}) \Rightarrow \bar{\bar{x}}_2$$

$$\vdots$$

$$S_{100} \Rightarrow (\bar{\bar{x}}_{100})$$

If I take all these means and try to plot it on a histogram then

$$\bar{\bar{x}} \approx \text{GD}\left\{\mu, \frac{\sigma^2}{n}\right\}$$

If we take multiple samples, then

$$\bar{\bar{x}} = \text{GD}\left\{\mu, \frac{\sigma^2}{n^2}\right\} \text{ in single bell curve.}$$

triangular form Chebyshev's Inequality

$x \sim \text{GD}(\mu, \sigma)$. P. S. ch. 6.1

Suppose, in above diagram standard deviation

$$P_r(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68\%$$

$$P_r(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95\%$$

$$P_r(\mu - 3\sigma \leq x \leq \mu + 3\sigma) = 99.7\%$$

$y \not\sim \text{GD}$

more? limit last

$$P_r(\mu - k\sigma < x < \mu + k\sigma) \geq 1 - \frac{1}{k^2}$$

higher k odd more narrow & higher K^2

typically $k=2$, get the result from

k specifies range of standard deviation

$$\text{Ansatz } (\mu - 2\sigma, \mu + 2\sigma) = X$$

$$\therefore P_r(\mu - 2\sigma < x < \mu + 2\sigma) \geq 1 - \frac{1}{4}$$

$$\text{Ansatz } (\mu - 2\sigma, \mu + 2\sigma) = 200 \quad k = \frac{2}{2}$$

$$\text{Ansatz } (\mu - 2\sigma, \mu + 2\sigma) = 200 \quad = 1 - \frac{1}{4}$$

$$\text{Ansatz } (\mu - 2\sigma, \mu + 2\sigma) = 200 \quad = \frac{3}{4} = 75\%$$

with lower bound & upper bound of $\mu \pm 2\sigma$

Pearson Correlation Coefficient.

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

with additional condition about σ_x

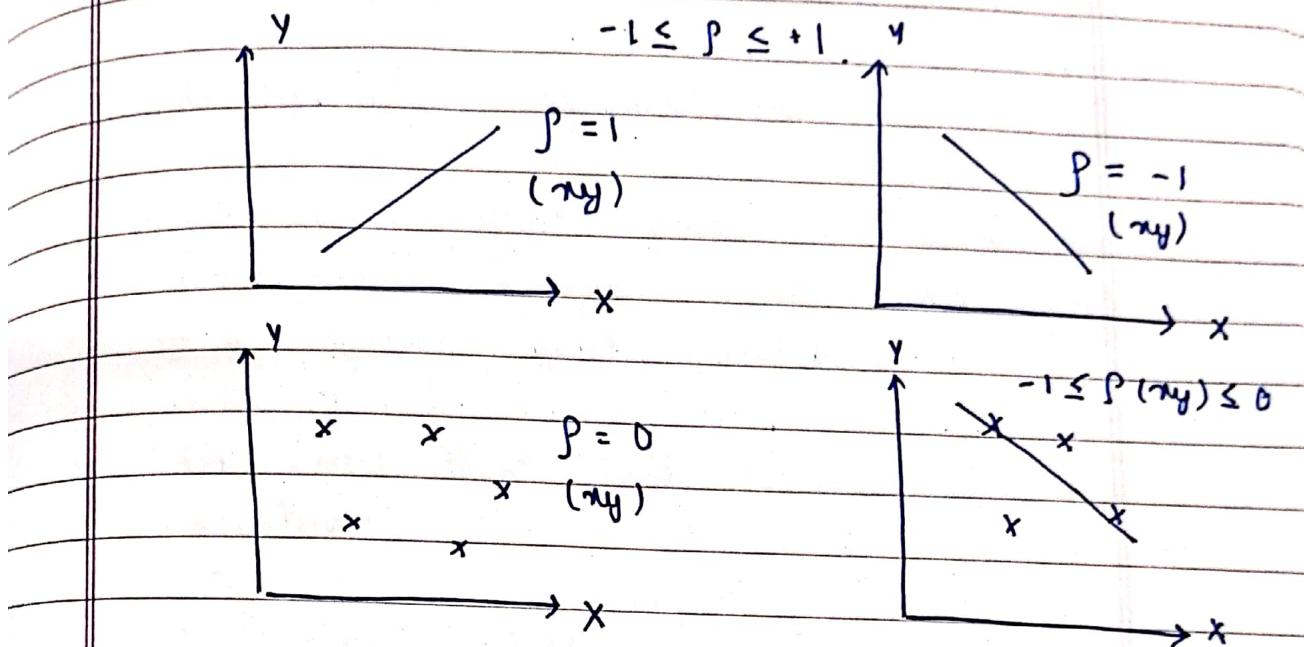
$$2 \quad \text{Pearson Coefficient} = \frac{\text{Covariance}(x, y)}{\sigma_x \cdot \sigma_y}$$

$$\text{Ansatz } (\mu, \sigma) \quad = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sigma_x \cdot \sigma_y}$$

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sigma_x \cdot \sigma_y}$$

If $x \uparrow y \uparrow$ or $x \downarrow y \downarrow$
 $+ve$ $-ve$

In covariance, we get to know the direction of relationship and do not know what is the strength of the relationship.



When two features have a very high correlation we can use either of them as a substitute for other.

Spearman's Rank Correlation

To find non-linear relationship between x & y

We try to find Pearson correlation of ranks of x and y .

Example:

100

106

100

86

101

99

103

97

113

112

110

Hours of TV / week

7

27

2

50

28

29

20

12

6

17

100 Hours / week (y_i) rank r_i ; rank y_i ; $d_i = d_i^2$

86 2 1 1 0 0

97 20 2 6 -4 16

99 28 3 8 -5 25

100 27 4 7 -3 9

101 50 5 9 100-3 25

103 29 6 9 100-3 25

106 7 7 7 100-7 16

110 17 8 5 +3 9

112 6 9 100-3 25

113 12 10 4 +6 36

ranked administration merit - on half of I

$$P = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$$

$$= 1 - \frac{6 \times 194}{10(100-1)} = -0.1175$$

with p-value = 0.62 using t-dist
 Correlation between 10 & hours per week watching TV is very low

Finding outliers in data with z-scores.

An outlier is a data point in a dataset that is distant from all other observations. A data point lies outside the overall distribution of the dataset.

Identification of outliers:

1 Data point that falls 1.5 times outside interquartile range above third quartile and below first quartile.

2 Data point that falls outside of 3 standard deviations

3 If z-score falls outside of 2 standard deviation

$$z = \frac{x - \mu}{\sigma}$$

Why should outliers exist in data?

1 Variability in data is high

2 An experimental measurement error

$\text{error} = \mu + \sigma Z$

What are the impacts of having outliers?

• Biases the mean

1 Causes various problems to statistical analysis

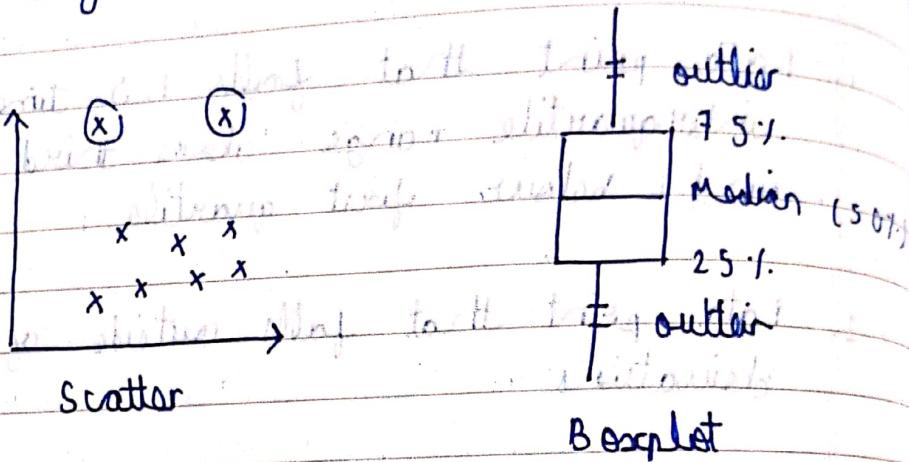
2 It impacts mean and standard deviation

IQR : 25% to 75% in $\frac{1}{4}^{\text{th}}$ to $\frac{3}{4}^{\text{th}}$

$$IQR = 7.5\% - 2.5\% = 5\%$$

Finding outliers in dataset

1. Using scatter plots
2. Box plots
3. Using Z-score
4. Using interquartile range.



$$\text{Outlier} = \frac{(x - \bar{x})}{\sigma}$$

```
def detect_outliers(data):
```

```
threshold = 3.
```

```
mean = np.mean(data)
```

```
std = np.std(data)
```

```
for i in data:
```

```
z-score = (i - mean) / std
```

```
if np.abs(z-score) > threshold:
```

```
outliers.append(i)
```

```
return outliers
```

Interquartile Range

1. Arrange the data in increasing order

2. Calculate first quartile and third quartile

3. Find interquartile range = $q_3 - q_1$

4. Find lower bound $q_1 + 1.5 \text{ IQR}$

5. Find upper bound $q_3 + 1.5 \text{ IQR}$

6. All data points between these bounds are in range.

Anything that lies outside the interquartile lower & upper bound, we remove them.

7. White area outside : Outliers

8. Feature scaling is the process of

normalizing features by dividing by their mean

Features has units and magnitude

Normalization helps to scale down the features between 0 to 1

Standardization is used to scale down the features based on standard normal dist

$$x_{\text{norm}} = \frac{x - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

$$z_{\text{std}} = \frac{x - \mu}{\sigma}$$

($\mu = 0, \sigma = 1$) after standardization

From now onwards all features are standardised

Mean is removed for intercept

standard deviation of each feature

Max minus min / 2