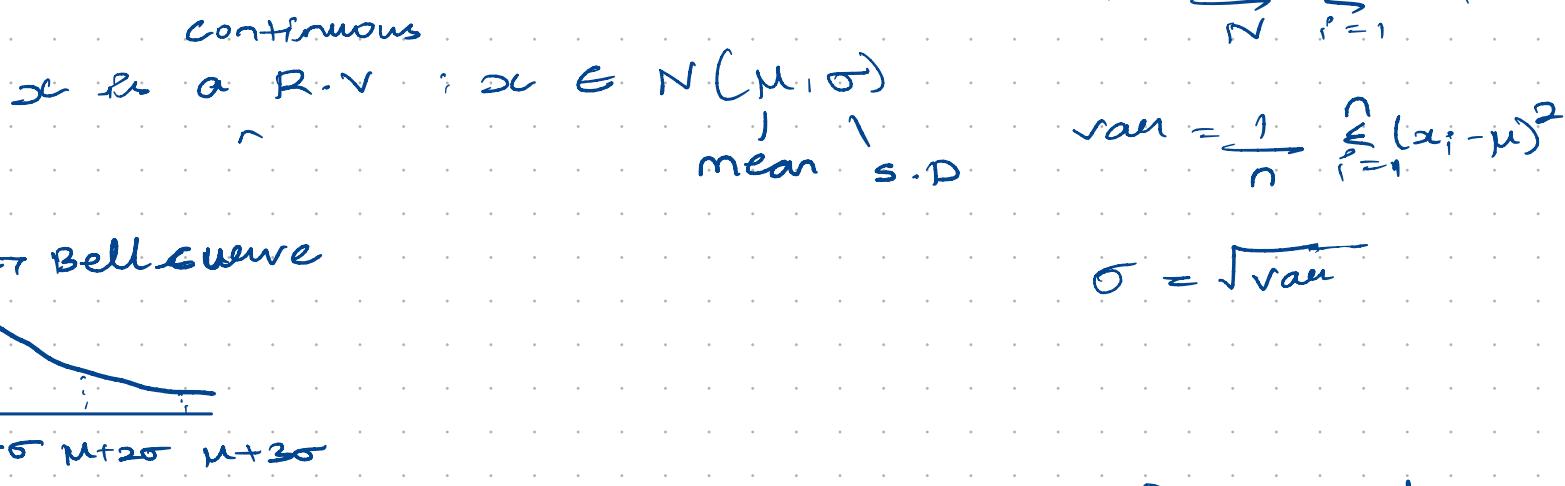


Gaussian Distribution (or) Normal distribution



$$P(\mu - \sigma \leq x \leq \mu + \sigma) = 68\%$$

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

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

Log Normal Distribution:

$$X \sim \text{Log Normal Distribution}$$

If $\ln(x)$ is normally distributed.

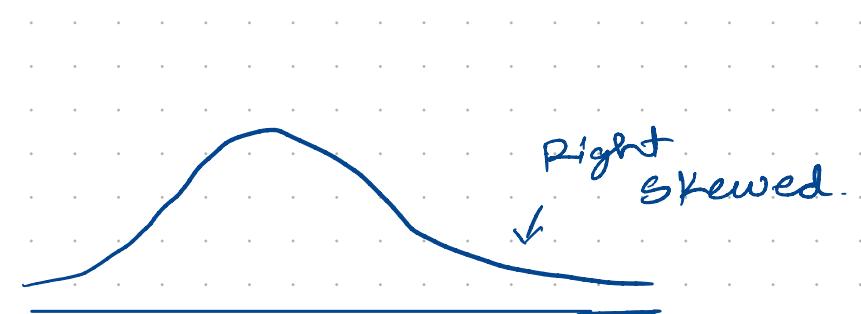
$$X = \{x_1, x_2, \dots, x_n\}$$

$$\{\ln(x_1), \ln(x_2), \dots, \ln(x_n)\}$$

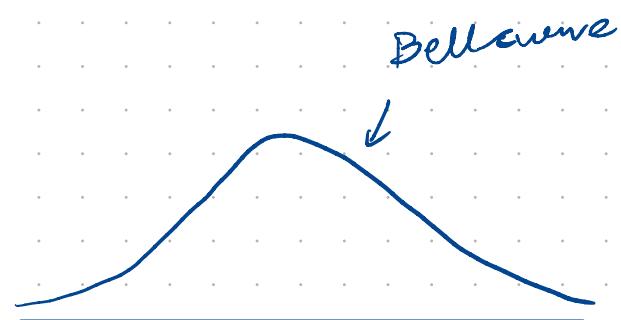
if this is normally distributed (formed a bell curve)

then X belongs to a Log Normal distribution.

Log Normal Distribution



Gaussian distribution



E.g.: Income of the people.

Product comments (length of description)

If you have a data, which is not normally distributed & it's distribution looks like it is skewed, then apply log to the data & now your new data will follow normal distribution.

Covariance

Consider an example in which you have -

size of the house

1500 sqm

1800 sqm

2000 sqm

price of the house

100K \$

200K \$

300K \$

Both are directly proportional -

$$\text{cov}(\underset{\downarrow}{\text{size}}, \underset{\downarrow}{\text{price}}) = \frac{1}{n} \sum_{i=1}^n (\underset{x}{x_i} - \mu_x) * (\underset{y}{y_i} - \mu_y)$$

$$\text{var}(\text{size}) = \frac{1}{n} \sum_{i=1}^n (\underset{x}{x_i} - \mu_x) (\underset{x}{x_i} - \mu_x)$$

Relationship b/w variance & covariance is

$$\text{cov}(x, y) = \text{var}(x) \quad \text{by the formula.}$$

Covariance gives you an insight about

If $x \uparrow$ then $y \uparrow$ or \downarrow by what %?

$$x \uparrow \quad y \uparrow = +ve$$

$$x \uparrow \quad y \downarrow = -ve$$

Central Limit theorem

Consider a R.V 'x' which may/may not belong to gaussian

distribution

$$x \neq N(\mu, \sigma^2)$$

Take a sample from the R.V x , $n \geq 30$ Data is drawn randomly

$$\text{Sample } S_1 = x_1, x_2, \dots, x_{30} = \bar{x}_1$$

sample $s_2 = x_1, \dots, x_{30} = \bar{x}_2$

$s_{50} = \dots = \bar{x}_{50}$

Now take all the sample means & plot by using histogram (for ex).

$$\bar{x} = \text{G.D}(\mu, \frac{\sigma^2}{n})$$

proper defn of CLT:

If you have a population with μ & s.d σ & take sufficiently large random samples from the population with replacement, then the dist of sample means will be approx normally distributed.

Chebyshev's Inequality:

$X \sim \text{G.D}(\mu, \sigma)$, If a R.V follows a gaussian distribution, then you can make use of 68-95-99.7-100% rule.

If you have a R.V Y that does not follow a Gaussian distribution, then how can you find out?

$$P(\mu - \sigma \leq x \leq \mu + \sigma) \geq 1 - \frac{1}{K^2}$$

Generalized formula is

$$P(\mu - K\sigma \leq x \leq \mu + K\sigma) \geq 1 - \frac{1}{K^2}$$

If $K = 3$

$$P(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \geq 1 - \frac{1}{9}$$

$$\geq \frac{8}{9} \approx 0.88 \approx 88\%$$

Also be written in the form

$$P(|x - \mu| = k\sigma) \leq \frac{1}{k^2}$$

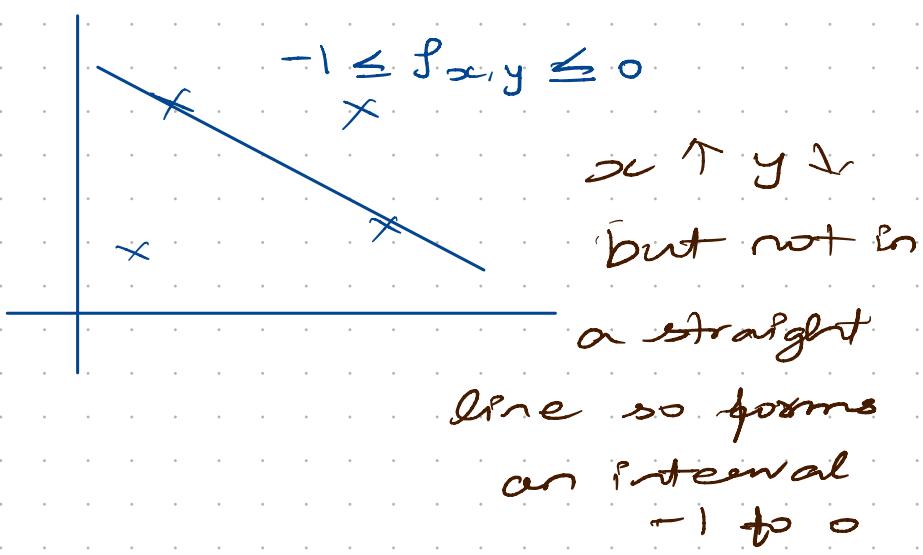
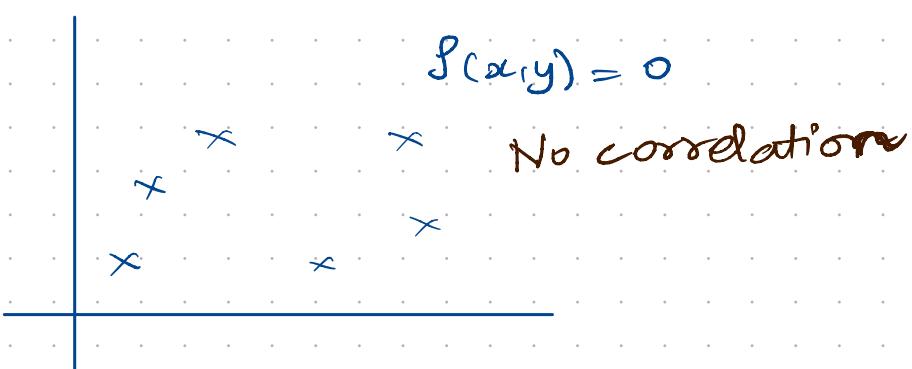
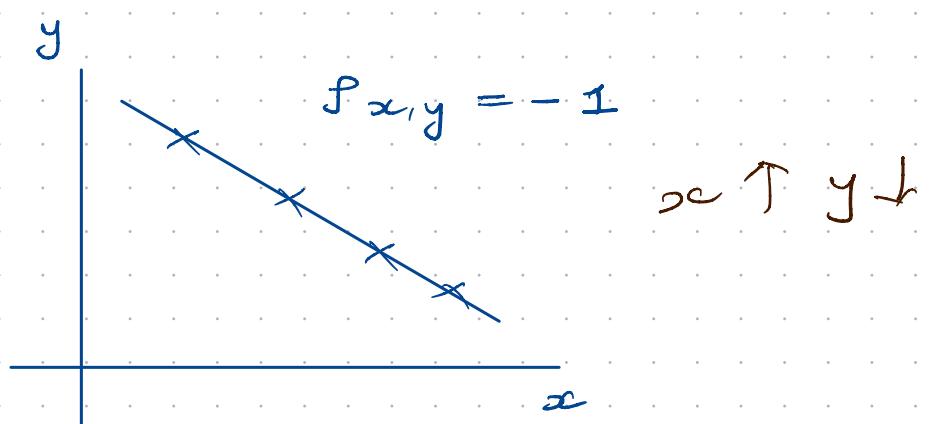
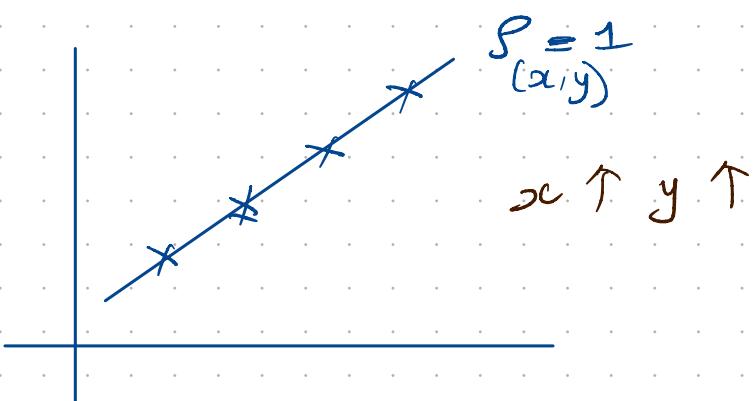
Pearson correlation Coefficient :

$$\rho_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

You know both in correlation.

- ① Strength (how much +ve or how much -ve).
- ② Direction of Relationship (+ve or -ve)

$$-1 \leq \rho \leq 1$$



Spearman's rank correlation coefficient:

$r_{g_x} \rightarrow$ Rank of x.

$$r_{g_x} = \frac{\text{cov}(r_{g_x}, r_{g_y})}{\sigma_{r_{g_x}} \sigma_{r_{g_y}}}$$

$$\rho_s = \text{Pearson's Correlation Coefficient}$$

Spearman's correlation coefficient (ρ_{S}) is simply calculating Pearson's correlation coefficient but for rank of x & y instead of just x & y .

Consider an example, where X is an IQ value and Y is number of hours spent in front of T.V.

X	Y
103	29
97	20
113	12
112	6
110	17

First find then Rank.

- Sort the first column. Create new column x_i^r & assign ranked values.
- Similarly do for Y .
- Create column d_i for diff b/w x_i^r & y_i^r
- Final column d_i^2

X	Y	rank x_i^r	rank y_i^r	d_i	d_i^2
97	20	1	4	-3	9
101	29	2	5	-3	9
110	17	3	3	0	0
112	6	4	1	3	9
113	12	5	2	3	9

$$P = 1 - \frac{6 \sum d_i^2}{n(n^2-1)} = 1 - \frac{6 \cdot 36}{20} = 1 - \frac{216}{5(24)} = 1 - \frac{216}{120}$$

$$= \frac{-96}{120} = -0.8$$

P = -0.8

$x \uparrow y \downarrow$

InterQuartile Range

4, 4, 6, 7, 10, 11, 12, 14, 15

Median = 10

IQR is finding the middle of the first & second half
diff b/w

(25 & 75th percentile).

4, 4, 6, 7 \Rightarrow Median = 5

11, 12, 14, 15 \Rightarrow Median = 13

IQR = 13 - 5 = 8

Find outliers in dataset using IQR:

1) Arrange data in ascending order.

2) calculate First & third quartile

3) Find IQR range ($Q_3 - Q_1$)

4) Find lower bound = $Q_1 - (1.5 * IQR)$

5) Find upper bound $Q_3 + (1.5 * IQR)$

Anything that lies outside of lower & upper is an outlier.

Standardization vs Normalization:

- * Normalization helps you to scale down the feature between 0 to 1.
- * Standardization helps you to scale down the feature based on standard normal distribution.

Min Max Normalization:

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

Min Max Scaler:

If you have two features whose magnitudes are huge away, so we should try to scale down between same scale.

Standardization (Z-score Normalization):

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

In standardization, all the features will be transformed in such a way that it will have the properties of a standard normal distb with $\mu=0$ & $\sigma=1$.

What is PDF, PMF, CDF?

Probability distribution functions.

obtained by
measuring

1) PDF \Rightarrow Probability density function \Rightarrow Continuous R.V

2) PMF \Rightarrow Probability Mass function \Rightarrow Discrete R.V

3) CDF \Rightarrow Cumulative Distribution function

obtained by
counting

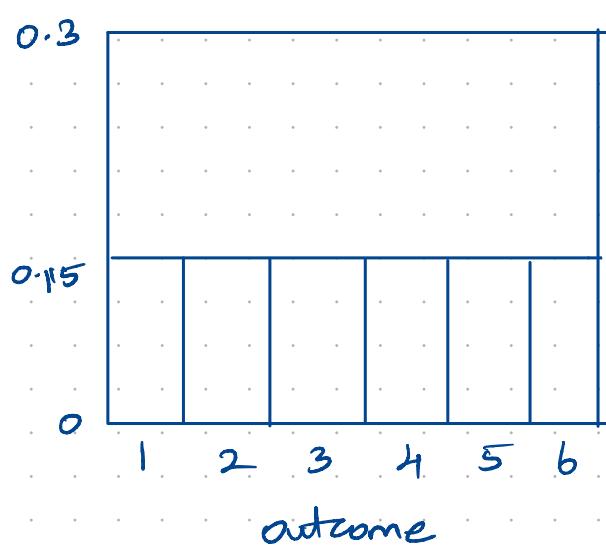
* PMF :-

Ex:- Dice Outcomes

All values in a dice has the probability of $\frac{1}{6}$. $= 0.167$

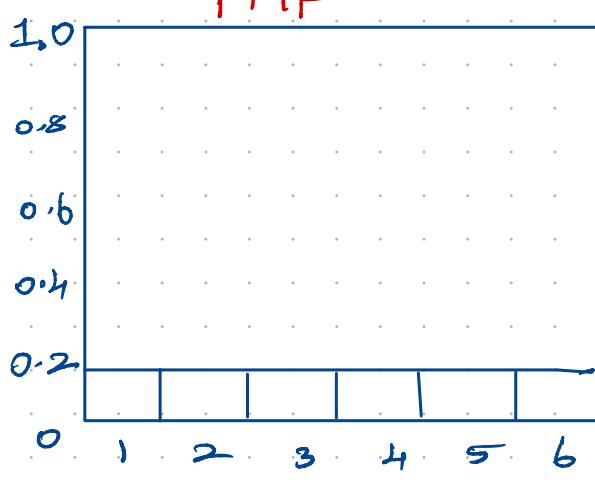
outcome

CDF for the same:

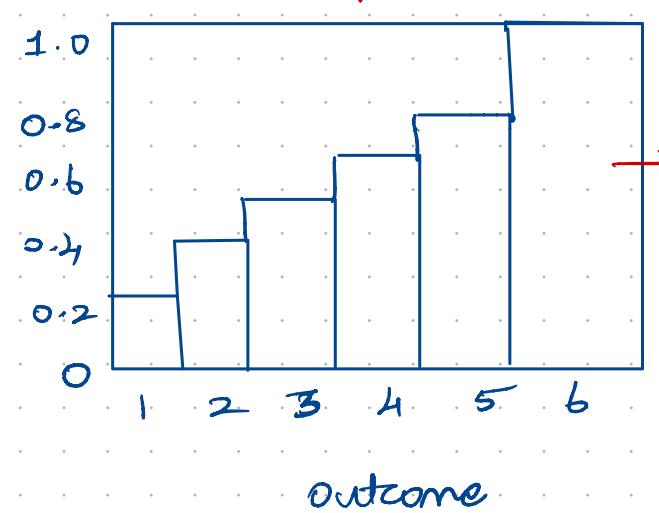


The diff b/w the above PMF & CDF diagram is its scale. PMF ranges from 0 to 0.3 & CDF from 0 to 1. Draw the PMF also in the scale of 0 to 1.

PMF



CDF



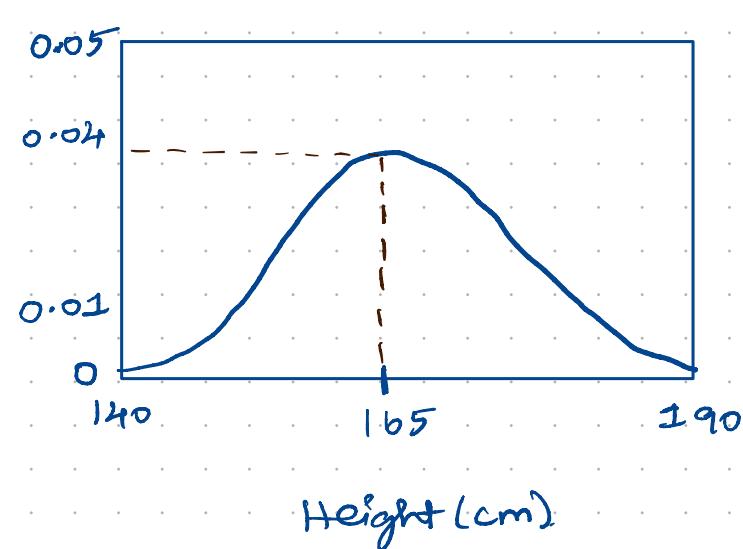
The final bar always needs to be 1 in CDF.

so in CDF, let's take $P(X \leq 3)$, then

$$\begin{aligned} P(X \leq 3) &= P(X = 1) + P(X = 2) + P(X = 3) \\ &= 0.167 + 0.167 + 0.167 \\ &\approx 0.5 \end{aligned}$$

* PDF :

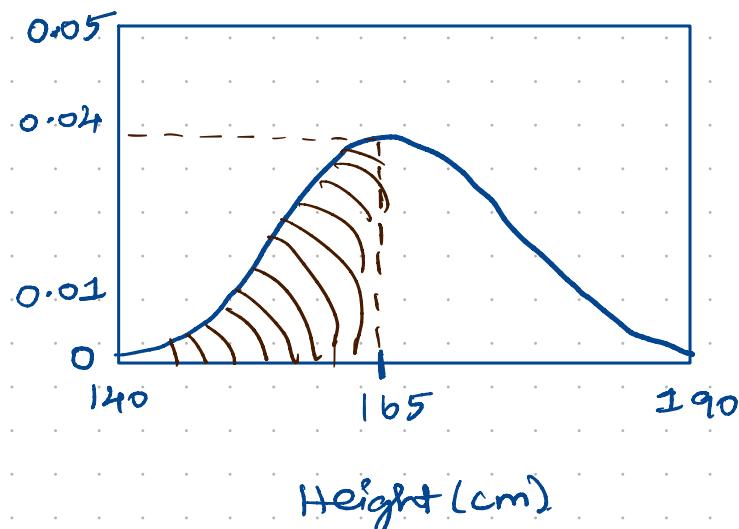
Ex:- Heights of Men.



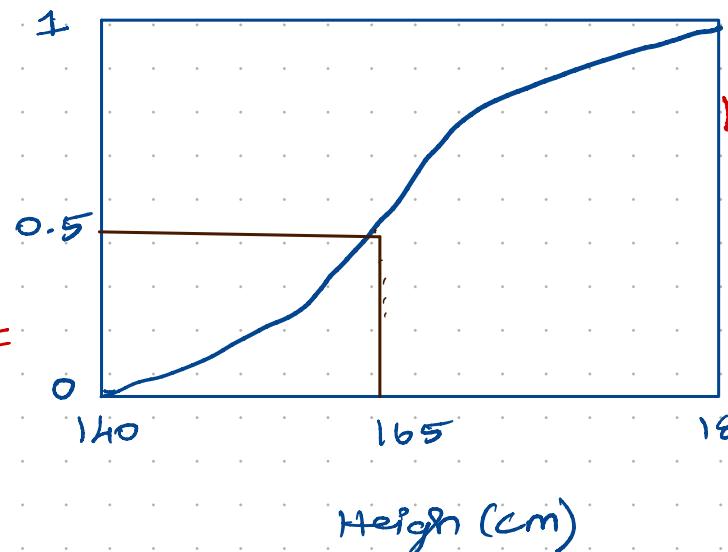
Does this plot means that $P(X = 165) = 0.03$?

NO

PDF



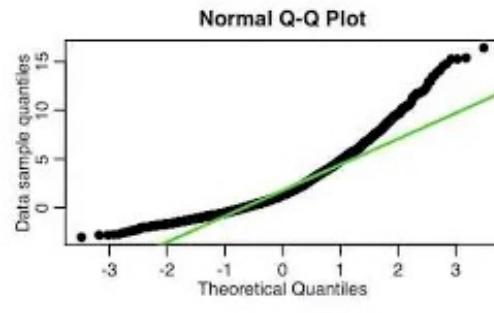
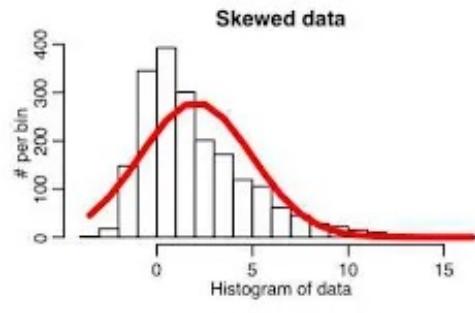
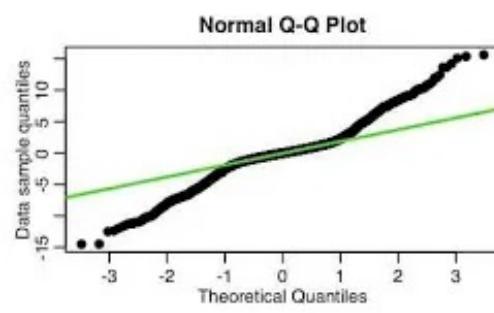
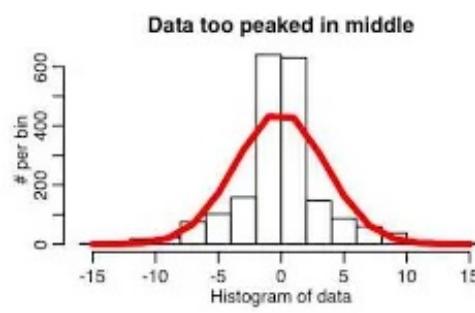
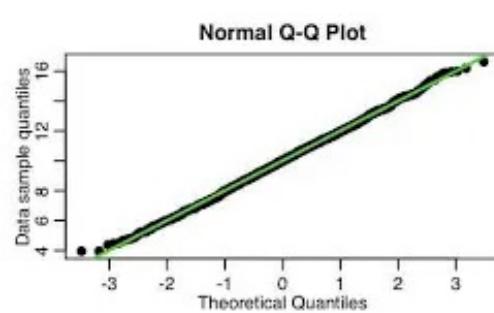
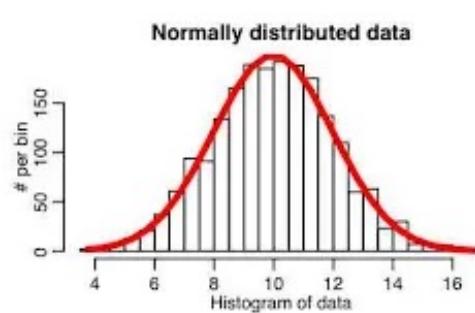
CDF



CDF is supposed to be in shape of "S". But in this diagram it is not. Bear with it 😊

Q - Q Plots (Quantile - Quantile Plots) - How to test if a R.V is normally distributed or not.

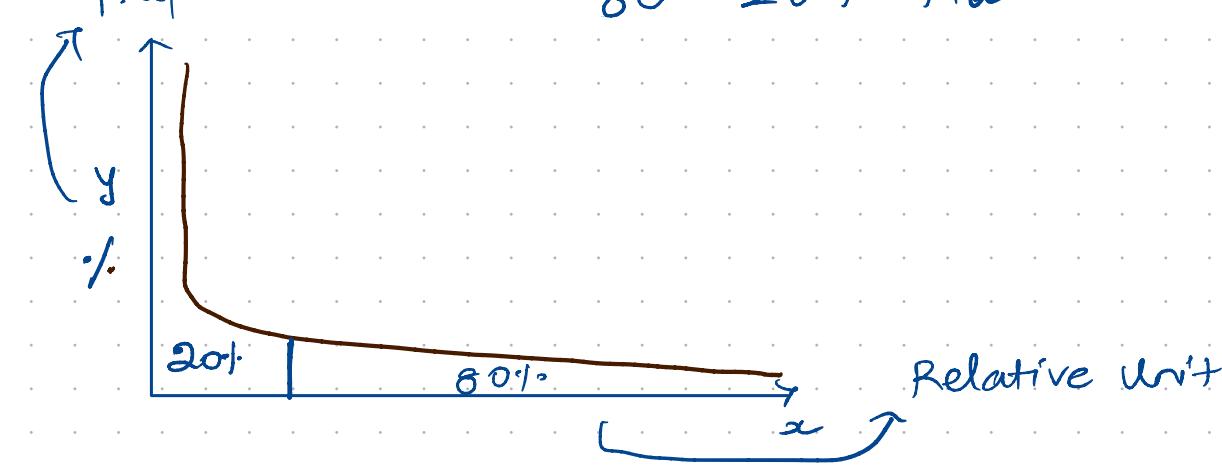
- * Helps in graphically analyzing, and comparing two probability distributions by plotting their quantiles against each other.
- * As we are comparing two distributions, if they are equal, then the Q-Q plot will lie in the straight line ($y = x$)



Power law distribution -

Proportional Unit

80 - 20% Rule



Ex:-

80% of sales are coming from 20% of overall products.

It describes a particular pattern in which few extreme values occurs much more frequently than most other values. It follows a long tail distribution, where small number of events have high frequency & vast majority have a low frequency.

Power-law:

~~XX~~ Relationship between two quantities, wherein a relative change in one quantity results in a proportional change in other quantity.

Confidence Intervals

Point estimate:

In many cases, knowing population mean is difficult bcoz of the large input size, so we calculate sample mean (\bar{x}). With the help of \bar{x} , we will estimate the parameter (μ) of population mean.

* Confidence Interval is the mean of the estimate plus & minus the variation in that estimate.

Ex:-

Avg size of sharks in the sea.

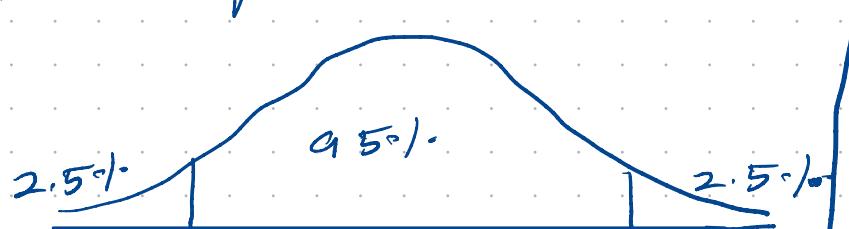
$$\text{population S.D} = 100$$

$$\text{sample } n = 30$$

$$\bar{x} = 500$$

Confidence Interval = Point Estimate \pm Margin Error.

Take 95% confidence Interval, then



$$= \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$Z = \text{significance value}$

As I have taken 95% CI, then 5% is the significance so 0.05

$$C.I = 500 \pm Z_{0.05/2} \frac{100}{\sqrt{n}} \quad Z = 1 - 0.025$$

$$= 0.975$$

$$= 500 \pm 2.025 \frac{100}{\sqrt{n}}$$

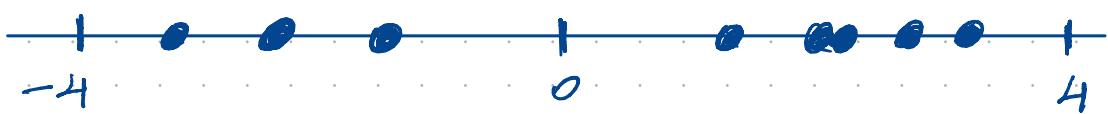
(Lookup in z-table
for $Z_{0.975}$)

$$C.I = 500 \pm 1.96 \left(\frac{100}{\sqrt{n}} \right)$$

$$500 - 1.96 * \frac{100}{\sqrt{n}} < C.I < 500 + 1.96 * \frac{100}{\sqrt{n}}$$

Confidence Interval Using Bootstrapping

Bootstrapping:



Create a new number line and add 8 values randomly sampled with replacement.

One such example will be:



As this is diff from original dataset, we will get diff mean.

Bootstrapping has 4 steps:

- 1) Make a Bootstrapped Dataset
- 2) calculate mean (any statistics)
- 3) Keep track of mean (any statistics)
- 4) Repeat 1-3 for bunch of time

Bootstrap for 1000 times & plot it in the histogram.

we can find 95% confidence Interval using the Bootstrapping.

Hypothesis Testing & Null Hypothesis

The hypothesis that there is no difference b/w things is called the Null Hypothesis.

Hypothesis testing is used to assess the credibility of a hypothesis by using Sample Data.

P-value:

P-values are numbers b/w 0 and 1. It tells you how likely is that the data could have occurred under null hypothesis. (Measure of evidence against null hypothesis)

The commonly used threshold in p-value is 0.05. It means we did the exact same experiment a bunch of times, then only 5% of the experiments results in wrong decision.

Ex:-

Test whether a new drug is effective in reducing blood pressure.

Null hypothesis (H_0):- Drug has no effect

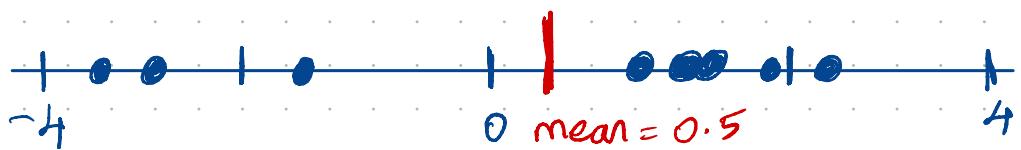
Alternate hypothesis :- Drug is effective.

We calculated a p-value of 0.02, but our significance level (threshold) is 0.05, then we would reject the null hypothesis

Smaller the p-value, stronger the evidence that you should reject null hypothesis.

Using Bootstrapping to calculate the p-value :-

Your original data is



Create a new line & shift the data points in such a way that mean will be 0.

* Then create a bootstrapped datasets from a collection of measurements with mean = 0, so the resulting histogram gives us a sense of what would happen if the Null hypothesis was true.

<https://youtu.be/N4ZQQqyIf6k> ⇒ Refer for the example with diagrams.

As the mean on the original data is 0.5, on the bootstrapped data check the value for

$$p\text{value} = P(\text{bootstrapped mean} \geq 0.5) + P(\text{bootstrapped mean} \leq -0.5)$$

By the ex: p-value = 0.63

As this is greater than 0.05, we fail to reject the null hypothesis.