

Probability and Statistics

1. Gaussian Distribution and its PDF
2. CDF of Gaussian Normal distribution
3. Central Limit Theorem
4. QQ Plot
5. Correlation Coefficient and its type
6. Power Law
7. Pareto Distribution Function
8. Co-variance, Pearson coefficient, Spearman rank correlation coefficient
9. Confidence Interval, CI means for random variable
10. Hypothesis Testing, P-value

1. Random Variable:

- It means set of possible outcomes of an experiment.
- A random variable is a variable whose value is unknown or a function that assigns values to each of an experiment's outcomes.
- Types of Random Variable are -
 - Discrete Random Variable
 - Continuous Random Variable

1.1 Discrete Random Variable-

- Random Variable which has a finite set of outcomes or values is called a discrete random variable
- Eg - A random experiment of tossing a coin results in either heads or tails, Rolling of a dice

1.2 Continuous Random Variable-

- Random variable which can take any real value is called a continuous random variable.
- Eg- The probability of height of a group of students.

Outlier

An observation point that is distant from other observations is called an outlier.

Note:

Mean and variance gets corrupted by an outlier, hence we use median and median absolute deviation.

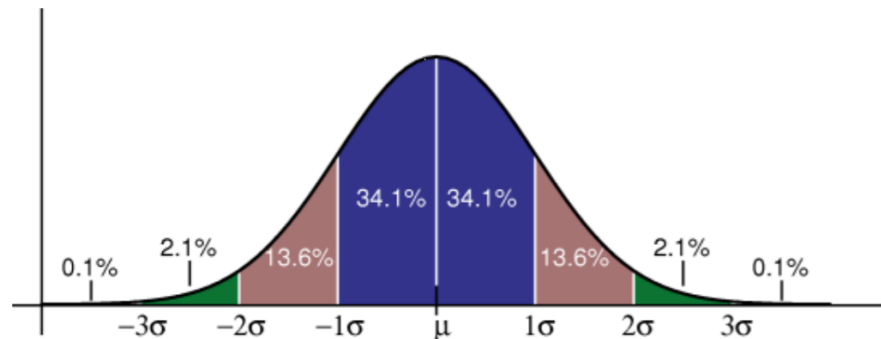
Median Absolute Deviation - MAD

It is the summation of all the absolute differences of the median from the individual points)

Gaussian/Normal Distribution -

If X is a continuous random variable that has a PDF like that of a bell shaped curve, then we say X has a distribution which is a **Gaussian Distribution**.

- Parameters of Gaussian Distribution-
 - Mean
 - Variance
- Variance** is the spread of the curve. So, if it is small, the curve is going to be steeper.
- Mean(μ)** is generally the peak of the distribution.

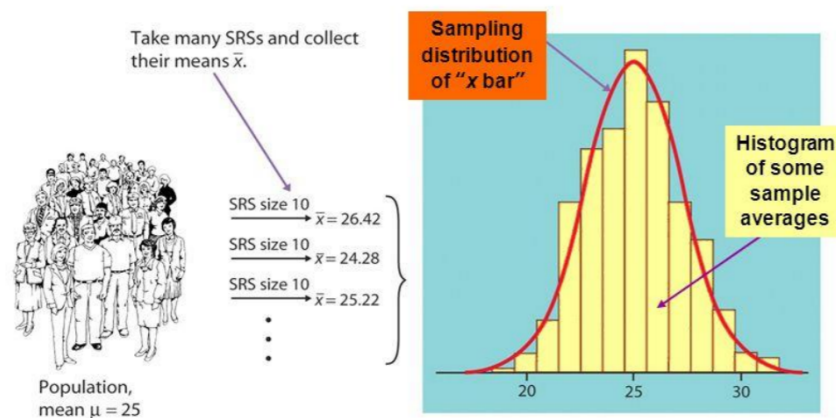


Standard Normal Distribution / Standard Normal Variate (Z):

A standard normal variate is a normal variate with mean $\mu=0$ and standard deviation $\sigma =1$.

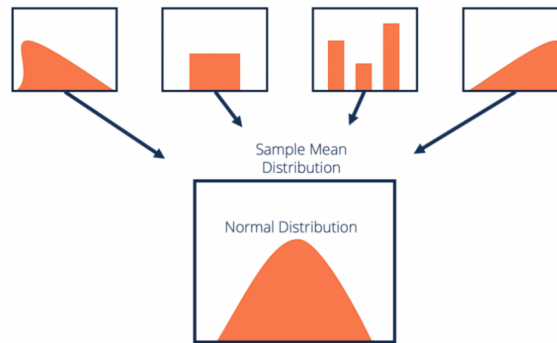
Central Limit Theorem CLT -

- Given a dataset with unknown distribution (it could be uniform, binomial or completely random), the sample means will approximate the normal distribution. Example -



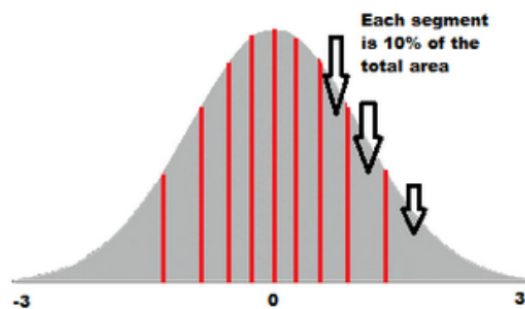
If we take m samples with n data points from a population distribution X which may not be Gaussian, and we calculate the mean of every samples, then, CLT states that-

- The distribution of the mean of the samples will be a Normal distribution with the population mean as it's Mean** (roughly).
- The standard deviation equal to (population variance/the number of data points in each sample) which can be written as, $N(m, s^2/2/n)$ as $n \rightarrow \infty$.
- If n is sufficiently large then sample means has Normal Distribution.



Q-Q Plot -

Quantiles are cut points dividing the range of a probability distribution into continuous intervals with equal probabilities or dividing the observations in a sample in the same way.



- 2 quantile is known as the Median
- 4 quantile is known as the Quartile
- 10 quantile is known as the Decile
- 100 quantile is known as the Percentile

Ques- How to check random variable(X) is Gaussian(Normal) Distributed ?

1. Sort the data of X in ascending order and then compute percentile of data.
2. Create a $Y \sim N(0,1)$ which is Standard Normal Distribution and sort them then compute Percentile.
3. Plot Q-Q plot of X and Y, if all points lie on same straight line then random variable X and Y have similar distribution. If X is Gaussian, Normal Distribution or Pereto then Y is also Gaussian, Normal Distribution or Pereto respectively.

Ques- How to check two distribution is statistically same, Here X and Y have different distribution ?

1. First, sort the given random variable in ascending order.
2. Then, generate a theoretical Quantile of the test distribution and sort it.
3. Then, graph the respective elements matching each of the given random variable and the theoretical distribution.
4. If the graph generates a straight line, then the 2 distributions are statistically same, else not.

Limitation- If no. of sample is small then its hard to interpret QQ-Plot.

Power Law Distribution-

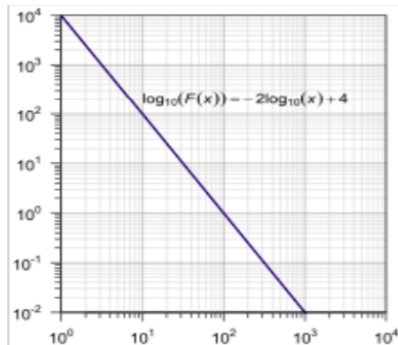
- The Power Law (Pereto Distribution/Scaling Law) states that a relative change in one quantity results in a proportional relative change in another.
- Eg- The simplest example of the law in action is a square; if you double the length of a side (say, from 2 to 4 inches) then the area will quadruple (from 4 to 16 inches squared).
- A power law distribution has the form $Y = k (X^\alpha)$, where: X and Y are variables of interest, α is the law's exponent, k is a constant.
- It roughly follows the **80-20 rule**, i.e., 80% of the points lie in the 20% of the region of the distribution.



Eg- Distribution of income, Magnitude of earthquakes, Size of corporations.

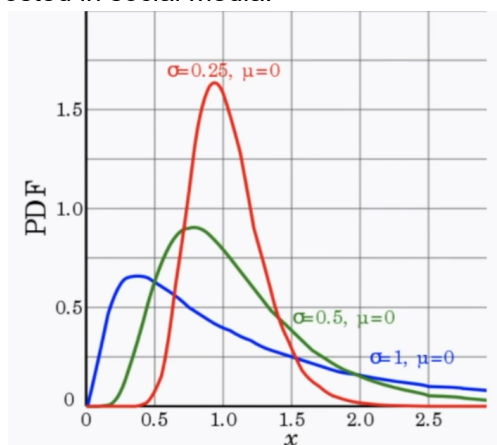
Ques- How to check Distribution is Power Law or not?

Ans- Take log of x and log of y and if they have a straight line from upper to lower in increasing order then it follows Power law.



Log Normal Distribution

- Random Variable X is considered to be Log-Normal if $\log(X)$ is Normally Distributed.
- It has a long tail like graph.
- Eg - Length of comments posted in social media.



Ques- How to test random variable $X(x_1, x_2, x_3 \dots x_n)$ have Log-Normal Distribution ??

1. Take $\log(x_1), \log(x_2) \dots \log(x_n)$.
2. Sort the data of X in ascending order and then compute percentile of data.
3. Create a $Y \sim N(0, 1)$ which is Standard Normal Distribution and sort them then compute Percentile.
4. Plot Q-Q plot of X and Y , if all points lie on same straight line then X have Log-Normal Distribution.

Power transform(Box-cos tranform)

- It is used to tranform Pareto distribution to Guassian distribution.
- By using Box-Cox method in X we get λ , if λ is equal to 0 then we get $\log(x)$ else we get in below eqn.

① $\text{box-cos}(X) = \sum_{i=1}^n \log(x_i) + \lambda \left(\frac{\sum_{i=1}^n x_i}{n} \right)$

② $y_i = \begin{cases} \frac{x_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \log(x_i) & \text{if } \lambda = 0 \end{cases}$

To find relationship betw two variable x, y belongs we use following methods are-

1. Co-variance
2. Pearson correlation coefficient(PCC)
3. Spearman rank corr coe.

1. Covariance -

- Covariance is the find of a relationship between two or more random variable.
- Covariance(X, Y) is +ve, if X and Y increases.
- Covariance(X, Y) is -ve, if X increases and Y decreases.

$$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n \{x_i - \mu_x\} * (y_i - \mu_y)$$

$$\text{Var}(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x) * (x_i - \mu_x)$$

✓ $\text{Cov}(X, X) = \text{Var}(X)$

$$\begin{cases} \text{Cov}(X, Y) = +ve \\ \text{Cov}(X, Y) = -ve \end{cases}$$

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

Drawback- By changing the units of measure, co-variance may differ.

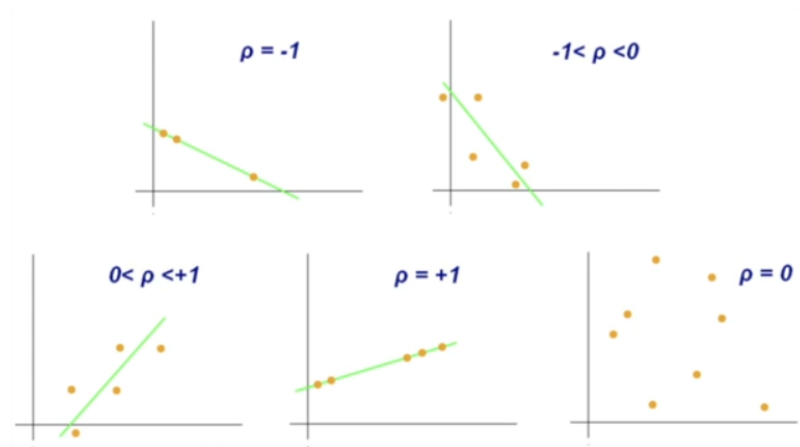
- It means for same dataset of X and Y(cm,kg) have different Covariance if we change the unit of X,Y(feet,pounds).

2. Pearson Correlation Coefficient PCC-

- PCC lie between $-1 \leq \text{PCC} < 1$, $\text{PCC} = 0$, when there is no relation between x and y, refer graph in each case.

$$\rho_{x,y} = \frac{\text{Cov}(X,Y)}{\sigma_x \sigma_y} \quad \sigma_x = \sqrt{\text{Var}(X)}$$

$\sigma(x)$ is Standard Deviation of x



- If $\rho = 0$, then there is no relation between X and Y
- The slope of the line doesn't matter
- **Drawback-** PCC cannot capture complex non-linear relationships. Like sinusoidal wave.

3. Spearman's Rank-Correlation Coefficient-

- Here we find $\text{rank}(x)$ and $\text{rank}(y)$, and calculate **PCC on $\text{rank}(x,y)$**
- To find rank just range in increasing order of x and y.
- As X and Y increases for linear or non-linear distri. $\text{rank} = 1$
- As X increases and Y decreases linear or non-linear distri. $\text{rank} = -1$.
- It works fine in monotonous increasing curve.

Correlation VS Causation

If two variable are co-related with each other that does not means one is the cause of other.

Hypothesis Testing

Hypothesis Testing a concept by which a statement is validated through a Proof by Contradiction.

1. We first start by choosing a Test-Statistics($\mu_2 - \mu_1$)
2. We provide a Null hypothesis rejecting the test statistics. We choose an Alternate hypothesis which is the complement of Null hypothesis.
3. Set the Significance Level or Critical Value.
4. We check the p-value or the Probability of the Test Statistics, and assume the Null hypothesis is True.
 - If, p-value is closer to 1, we accept the Null Hypothesis.
 - Else if p-value is closer to 0, we accept the Alternate Hypothesis.

Critical Value-

- * Normally Critical Value is 5%.
- * If our test score lies in the critical zone, we reject the Null Hypothesis and accept the Alternate Hypothesis

p-Value-

- p-Value is the probability of observation that we already made given Null-Hypothesis is True.
 $p\text{-value} = P(\text{Obs}/H_0)$

Example -

- Given a coin, determine if the coin is biased towards Head or Not.
 Basic probability we know-
 Biased towards Head- $P(H) > 0.5$
 Not-Biased towards Head - $P(H) = 0.5$

Experiment - Flip a coin 5 times and count the no. of heads. To check Coin is Biased or not.
 Count no. of Head(X) is our Test-Statistic.
 After performing experiment we get no. of Head, $X = 5$
 Our Null-Hypothesis(H_0) is Coin is not biased towards Heads.
 Our Alternate Hypothesis(H_1) is Coin is biased towards Heads.<

$$P(X=5|\text{coin is not Biased}) = P(\text{Obs}|H_0)$$

Probability of coin is not biased = $1/2$

Probability of getting 5 Heads in 5 toss assuming to the coin is not biased is $1/2^5$

$$P(\text{Obs}|H_0) = 1/2^5$$

$$= 1/32 = 0.03 = 3\%$$

So here p-value = 3%

- It means that there is 3% chance of getting 5 Heads in 5 flips if the coin is not biased towards Head.
- If $P(\text{Obs}|H_0) < 5\%$, then H_0 may be incorrect.
 It means our assumption or H_0 is not True \Rightarrow Reject H_0 or Accept $H_1 \Rightarrow$ Coin is Biased Towards Head.

Permutation Test(Resampling) -

- It is used to find p-value from sample.
- In Real World there is no exact p-value exist, its depend on where we are using p-value, Eg in medicine we take p-value is 1% it means here are very-very sure about certain things.

Example-

We are here to find which class students have heights greater than other class.

- Let there are two class A and class B having 50-50 students in each class with μ_1 and μ_2 as mean of Class A,B respectively. Mean Difference is $\mu_1 - \mu_2 = \Delta$
- We jumble both class students and then randomly separate with 50-50 students in let say class X and Y, with μ_1 and μ_2 as mean value.
- Difference of μ_1 and μ_2 is δ_1 .
- Repeat step 2 and 3 to 10k times and we get $\delta_2, \delta_3 \dots \delta_{10k}$.
- After getting all $\delta_1, \delta_2, \delta_3 \dots \delta_{10k}$, Sort them to get - $\delta_1', \delta_2', \delta_3' \dots \delta_{10k}'$.
- Then we fit our Actual Δ to above sorted list of sorted means.
- After placing our Δ in above sorted list of sorted means and Lets assume there are x% points greater than sorted means then we can say our p-value is x%. If x is 5 then p-value = 5%