

# Data Science with R

## Day 3 - Statistics for Data Science – Basics and Advanced



# Today's Agenda

## ✓ Basics of Statistics

- Type of Random Variables -  
Based on Scale of Measurement
  - Nominal
  - Ordinal
  - Interval
  - Ratio
- Variance
- Standard Deviation

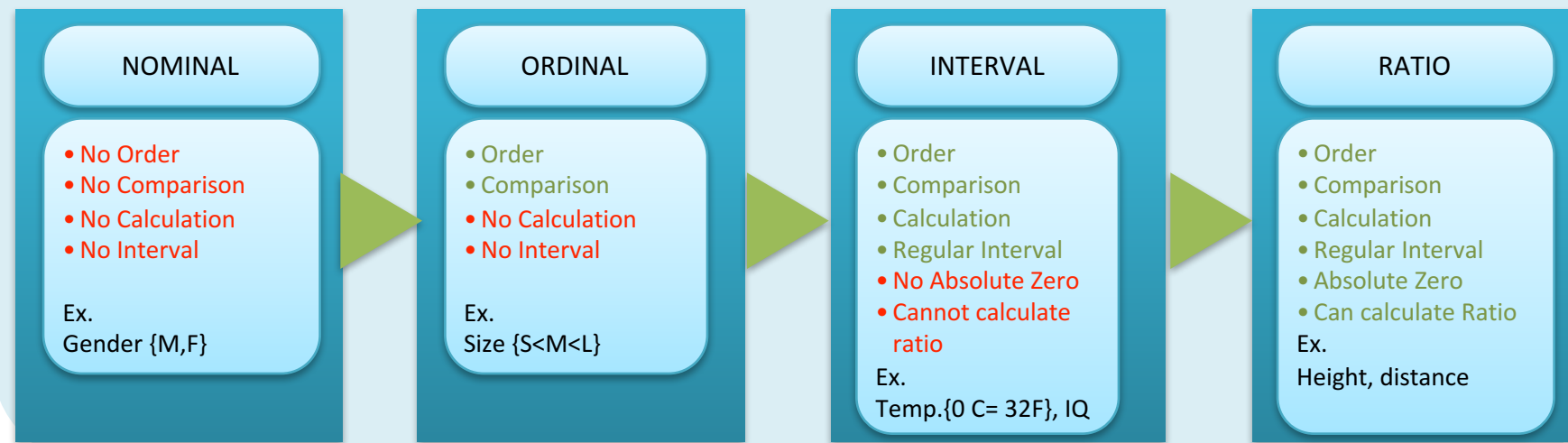
## ✓ Advanced Statistics P1

- Normal Distribution
- Standard Normal Distribution  
and Z- Score
- Binomial Distribution
- Poisson Distribution

# Basics of Statistics

## Type of Random Variable - Based on Scale of Measurement

Based on scale of measurement RV can belong to one of the following types -



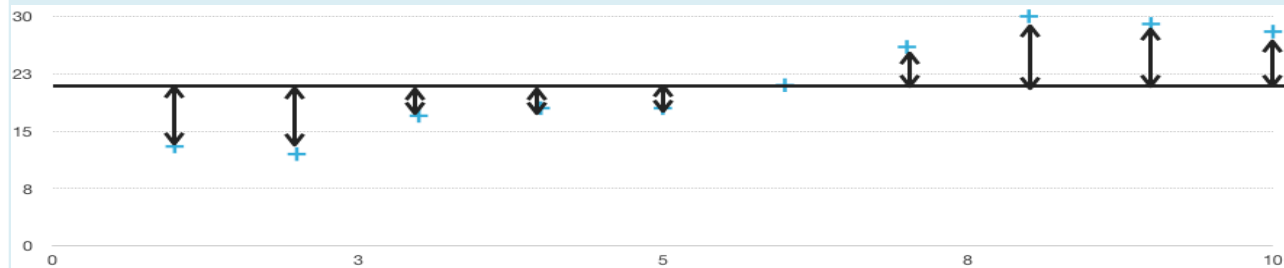
Hemant Rathore

# Basics of Statistics

## Variance( $\sigma^2$ ) and Standard Deviation( $\sigma$ ) –

**Variance ( $\sigma^2$ )** - Squared deviation of value from Mean :  $\text{Var}(X) = 1/n [\text{Sum } [X-\text{Mu}]^2]$

**Standard Deviation ( $\sigma$ )** - Square Root of Variance :  $\sqrt{1/n [\text{Sum } [X-\text{Mu}]^2]}$



$$\text{Variance} = 398/10 = 39.8$$

$$\text{SD} = 6.3$$

Data Point	Speed	X-Mu	(X-Mu) <sup>2</sup>
1	13	-8	64
2	12	-9	81
3	17	-4	16
4	18	-3	9
5	18	-3	9
6	21	0	0
7	26	5	25
8	30	9	81
9	29	8	64
10	28	7	49
Average	21		398

Hemant Rathore

# Basics of Statistics

## Mean( $\mu$ ), Variance( $\sigma^2$ ) and Standard Deviation( $\sigma$ ) -

Discrete random variable :

### Mean ( $\mu$ ) -

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

$$\mu = \sum_{i=1}^n p_i \cdot x_i.$$

### Variance ( $\sigma^2$ ) -

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

$$\text{Var}(X) = \sum_{i=1}^n p_i \cdot (x_i - \mu)^2,$$

### Standard Deviation ( $\sigma$ ) -

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

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

<https://en.wikipedia.org>

Hemant Rathore

# Basics of Statistics

## Variance( $\sigma^2$ ) and Standard Deviation( $\sigma$ ) -

Continues random variable -

Mean ( $\mu$ ) -

$$\mu = \int x f(x) dx$$

Variance ( $\sigma^2$ ) -

$$\text{Var}(X) = \sigma^2 = \int (x - \mu)^2 f(x) dx$$

Standard Deviation ( $\sigma$ ) -

$$\sigma = \sqrt{\int_{\mathbf{x}} (x - \mu)^2 p(x) dx},$$

[https://en.wikipedia.org/wiki/Mode\\_\(statistics\)](https://en.wikipedia.org/wiki/Mode_(statistics))

Hemant Rathore

# Advanced Statistics

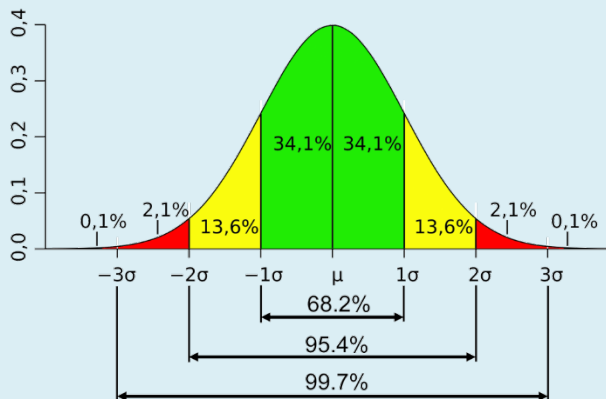
## Normal Distribution

Normal or Gaussian or bell shaped curve distribution is a very common continuous probability distribution. Normal Distribution has bell shaped curve, it's a symmetric single model distribution with highest density at and around the mean :

Ex. Age, Marks

### Some Important properties :

- Mean = Median = Mode
- Area within 1 Std. Dev around the mean ~ 68.3 %
- Area within 2 Std. Dev around the mean ~ 95.4 %
- Area within 3 Std. Dev around the mean ~ 99.7 %



$$f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

[https://en.wikipedia.org/wiki/Normal\\_distribution](https://en.wikipedia.org/wiki/Normal_distribution)

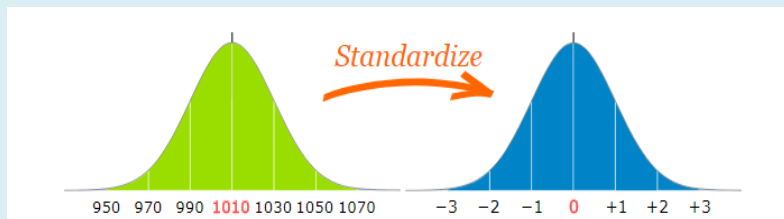
# Advanced Statistics

## Standard Normal Distribution and z Score/ z Statistic

Special case of Normal distribution with Mean = 0 and Variance = 1, Std. Deviation = 1. **it has total area under the curve = 1** which represents probability.

Any Normal distribution can be converted into Standard Normal Distribution by applying following transformation :

**$Z = (x - \mu) / \sigma$**  {This is called as **z Score**, it tells us how many SD far we are from mean in SND}



### Why do we need Standard Normal Distribution?

To easily make the decision about Probability distribution by using some known properties of Standard Normal Distribution and z table.

<https://www.mathsisfun.com/data/standard-normal-distribution.html>



# Advanced Statistics

## Binomial Distribution

**Bi → 2, Nomial → Nominal → Only 2 possible outcomes (Success or Failure)**

When we perform any given experiment multiple times and we are interested in knowing #successes, this type of experiments are known as Binomial experiments, Ex. Flipping the coins multiple times. Using Binomial Distribution we can answer probability related questions for any Binomial experiments.

Probability of getting 'x' #Successes out of 'n' trials using Binomial Distribution –

$$P(x) = {}^n C_x P^x (1-P)^{n-x} ; P = \text{Probability of Success in 1 trial}$$

**Some Important properties :**

- N Fixed Number of Trials
- Only 2 Possible Exclusive Outcomes
- Probability of success remain same during the experiment
- All the trials are independent

Hemant Rathore

# Advanced Statistics

## Poisson Distribution

When we analyze the probability of occurrence of any event during some specified interval of time or according to some other binding conditions.

Probability of 'x' occurrence using Poisson Distribution –

$$P(x) = (\lambda^x e^{-\lambda})/x!; \lambda = \text{Mean/Expected \#Occurrence}$$

### Some Important properties :

- All the occurrences are independent
- Expected #Occurrence doesn't change over the period of time

Hemant Rathore

“Qs & As”

# Data Science with R

## Day 4 - Statistics for Data Science – Basics and Advanced



# Today's Agenda

## ✓ Inferential Statistics

- Sampling
- Inferential Statistics
- Sampling Distribution
- Central Limit Theorem
- Central Limit Theorem Exercise

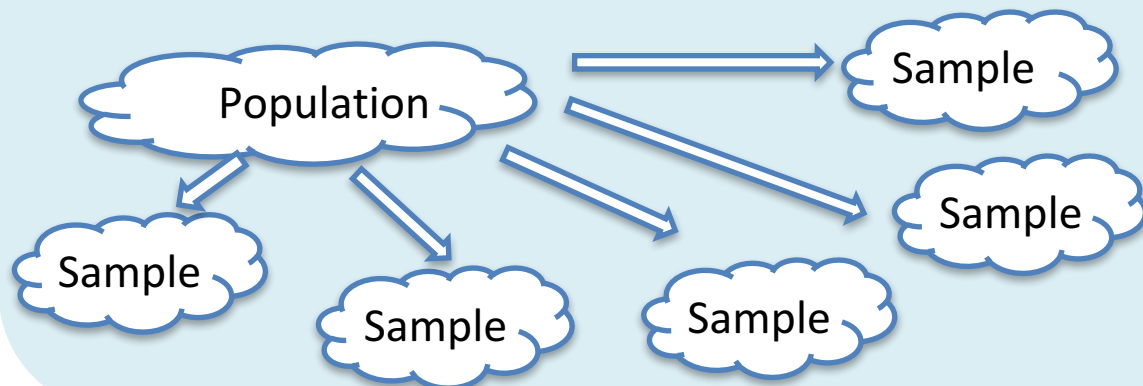
## ✓ Hypothesis Testing

- Hypothesis and hypothesis Testing
- One tail/Two tail test
- Type I and Type II Errors
- Hypothesis Testing using z test
- Hypothesis Testing t test

# Advanced Statistics - Inferential Statistics

## Sampling

Sampling is taking random samples from over all population, sampling is done in order to make some judgements about overall population because many a time it is not possible or practical to analyze the overall population and instead we can get approximately same results even using sampling with sufficient sample size



Hemant Rathore

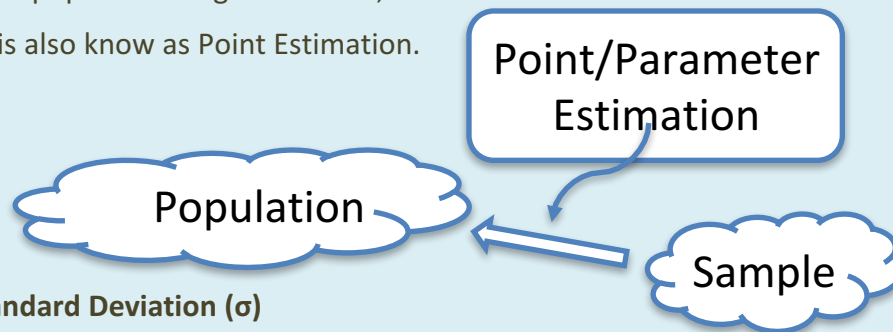
# Advanced Statistics - Inferential Statistics

## Inferential Statistics

With inferential statistics, we try to reach conclusions that extend beyond the immediate data alone. For instance, we use inferential statistics to try to infer from the sample data what the overall population might think. Or, we use inferential statistics to make judgments of the probability of the overall population. This is also known as Point Estimation.

### Point Estimators

- Sampling Mean ( $\mu$  XBar)  $\rightarrow$  Population Mean ( $\mu$ )
- Sampling Standard Deviation ( $\sigma$  XBar)  $\rightarrow$  Population Standard Deviation ( $\sigma$ )



Hemant Rathore

# Advanced Statistics - Inferential Statistics

## Sampling Distribution

When we use the distribution of samples taken randomly from population to make judgment about the overall population. Different Samples taken from same population can show different characteristics this is known as sampling variability. Larger the sample size – less the variability.

### Expected Value $E(x)$ or Sampling Mean ( $\mu$ XBar)-

We take multiple samples from overall population and analyze the distribution of these samples to make the decision about overall population. The mean of these samples is known as expected value or sampling mean and it can be considered as Overall population mean ( $\mu$ ).

**Sample Size (n) -> Very large than Expected Value  $E(x) \rightarrow \mu$**

### Standard Error of the Mean ( $\sigma$ XBar)-

Standard deviation of sampling distribution is known as Standard Error of the Mean.

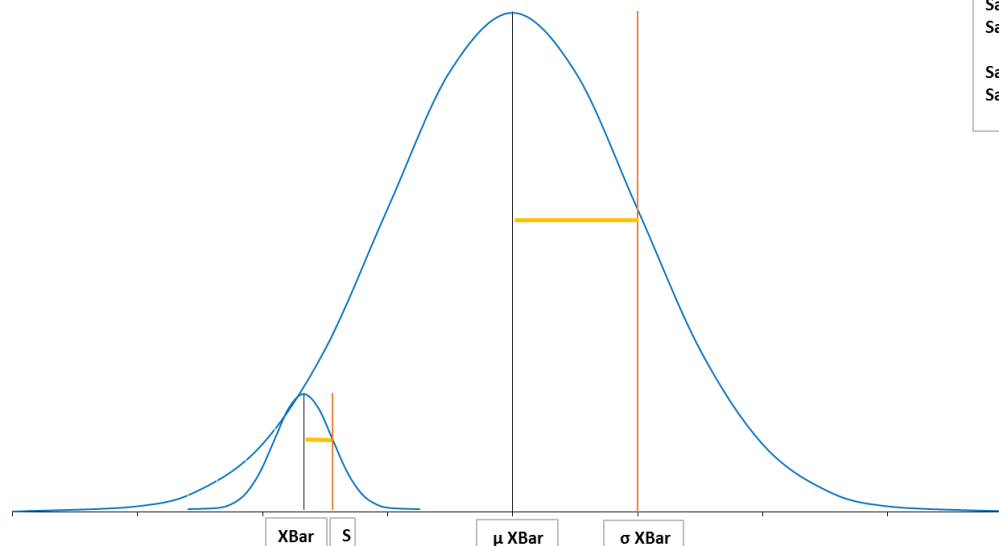
**Sample Size (n) -> Very large than Standard Error of the mean  $\rightarrow 0$**

Hemant Rathore



# Advanced Statistics - Inferential Statistics

## Sampling Distribution



Population Mean -  $\mu$   
Population SD -  $\sigma$

Sample Mean -  $\bar{X}$   
Sample SD -  $S$

Sampling Mean -  $\mu \bar{X}$   
Sampling SD -  $\sigma \bar{X}$

Hemant Rathore

# Advanced Statistics - Inferential Statistics

## Central Limit Theorem (CLT)

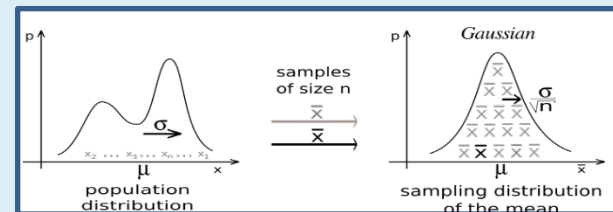
*"Sample Mean will be approximately normally distributed for larger sample size regardless of the original distribution from which we are taking samples."*

With Mean = Population Mean ( $\mu$ )

SD =  $\sigma / \sqrt{n}$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

{in case  $\sigma$  is not known then SD =  $s / \sqrt{n}$ ,  $s$  = Sample SD}



## Application of CLT

So we can use standard normal distribution concepts for any non normal population by taking samples because as per CLT Samples will be normally distributed for large sample size.

From CLT we know, sampling SD  $\sigma_{\bar{x}} = \sigma / \sqrt{n}$

From Standard Normal distribution we know –  $Z = (x - \mu) / \sigma$

So for any sampling distribution we can say –  $Z = (X - \mu) / (\sigma / \sqrt{n})$ , so now we can calculate the probability using SND for any Non normal Population.

Hemant Rathore

# Advanced Statistics - Inferential Statistics

## Exercise - Central Limit Theorem

A large freight elevator can transport a maximum of 9800 pounds. Suppose a load of cargo containing 49 boxes must be transported via the elevator. Experience has shown that the weight of boxes of this type of cargo follows a distribution with mean  $\mu = 205$  pounds and standard deviation  $\sigma = 15$  pounds. Based on this information, what is the probability that all 49 boxes can be safely loaded onto the freight elevator and transported?

**Solution** – Given :  $\mu = 205$  ,  $\sigma = 15$  ,  $n=49$ ; Average total weight =  $49 \times 205 = 10045 > 9800$

We know nothing about the original probability distribution whether its normal or not but from CLT we know sample mean will be normally distributed,

We are interested in the weight of 49 boxes not 1 so let's calculate :  $\mu$  and  $\sigma$  for 49 boxes :

$\mu = 10045$  ,  $\sigma = 15 \times 49 = 735$ , now we need to know the probability that total weight would be  $\leq 9800$  so  $X=9800$

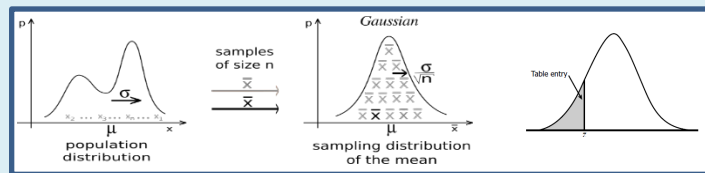
$$Z = (X - \mu) / \sigma / \sqrt{n}$$

$$(9800 - 10045) / (735 / 7)$$

$$= -245 / 105$$

$$Z = -2.33$$

let use z table to get the probability for  $z \leq -2.33 \rightarrow 0.0099$



# Advanced Statistics – Hypothesis Testing

## Hypothesis

A hypothesis (plural hypotheses) is a proposed explanation for a phenomenon. In Statistics Hypothesis can be any theory about the data that we want to validate (generally accept or reject) – we will be mainly working of two type of hypotheses :

1. **Null Hypothesis ( $H_0$ )** – Current Assumption or Theory which is currently assumed to be correct
2. **Alternative Hypothesis ( $H_1$ )** – Claim or theory that we want to prove

Ex.  $H_0$ : While flipping a coin the probability of getting head is 0.5;  $H_1$  : Probability of getting head is less than 0.5

Hemant Rathore

# Advanced Statistics – Hypothesis Testing

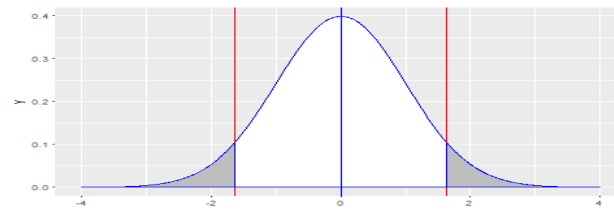
## Hypothesis Testing

Validating the null hypothesis ( $H_0$ ) against some Alternative Hypothesis ( $H_1$ ) based on some given **sample data**, Steps involved in Hypothesis Testing using P values –

1. Define your Null and Alternative hypothesis,  $H_0$  &  $H_1$
  2. Decide the type of test (One tail or two tail)
  3. Define level of significance  $\alpha$ , generally assumed to be 0.05 or 0.01
  4. Find the Test Statistics TS (t Test or z test)
  5. Find P Value
  6. Reject the null hypothesis or you may accept the alternative hypothesis if  $P < \alpha$
- **P value** – Probability of getting the given sample or even more extreme samples
  - **Significance level ( $\alpha$ )** – Minimum acceptable P Value/ border line

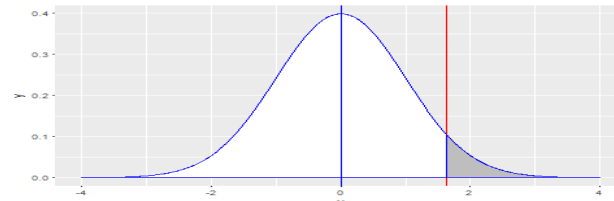
Two tail test –

$H_0 : \mu = 200$   
 $H_1 : \mu \neq 200$



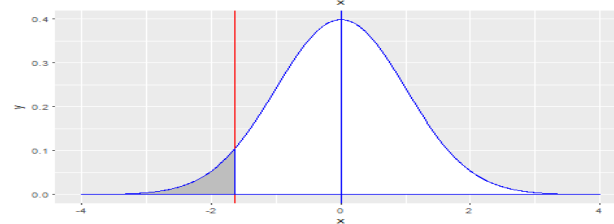
One tail Test –

$H_0 : \mu \leq 200$   
 $H_1 : \mu > 200$



One tail Test –

$H_0 : \mu \geq 200$   
 $H_1 : \mu < 200$



Hemant Rathore

# Advanced Statistics – Hypothesis Testing

## Type I Error or False Positive

Getting Positive result when it should be Negative in reality.

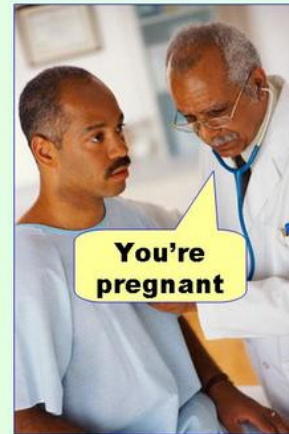
Rejecting null hypothesis ( $H_0$ ) while  $H_0$  is correct and should not be rejected, Probability of Type I error is known as Alpha Risk.

## Type II Error or False Negative

Getting Negative result when it should be Positive in reality.

Failing to Reject null hypothesis ( $H_0$ ) while  $H_0$  is not correct and should be rejected, Probability of Type II error is known as Beta Risk.

**Type I error**  
(false positive)



**Type II error**  
(false negative)



Hemant Rathore

# Advanced Statistics – Hypothesis Testing

## Z-Test for Hypothesis Testing

Z-Test is used to perform hypothesis testing when we know population Standard Deviation ( $\sigma$ ) or the sample size  $n > 30$ , Steps for 1 sample z test –

1. Define your Null and Alternative hypothesis,  $H_0$  &  $H_1$
2. Define level of significance  $\alpha$ , generally assumed to be 0.05 or 0.01
3. Find the Test Statistics using  $TS = (X - \mu) / (\sigma / \sqrt{n})$  or  $TS = (X - \mu) / (s / \sqrt{n})$  {when  $\sigma$  is not known but  $n \geq 30$ }
4. Find P Value using Z Table and TS, if it's a two sided test then double P value
5. Reject the null hypothesis or you may accept the alternative hypothesis if  $P < \alpha$

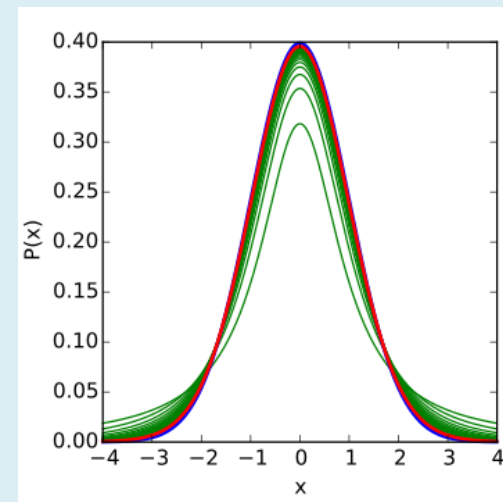
Hemant Rathore

# Advanced Statistics – Hypothesis Testing

## T-Test for Hypothesis Testing

T-Test is used to perform hypothesis testing when we don't know population Standard Deviation ( $\sigma$ ) and sample size  $n < 30$ , Steps for 1 sample t test –

1. Define your Null and Alternative hypothesis,  $H_0$  &  $H_1$
2. Define level of significance  $\alpha$ , generally assumed to be 0.05 or 0.01
3. Calculate Degree of Freedom  $DF = n-1$
4. Find the Test Statistics using  $TS = (X - \mu) / (s / \sqrt{n})$
5. Find P Value or P Range using T Table for calculated TS and DF
6. Reject the null hypothesis or you may accept the alternative hypothesis if  $P < \alpha$



Hemant Rathore



**“Qs & As”**



For more information or to set up an appointment, please contact us today.

**[jointact@tactlearn.com](mailto:jointact@tactlearn.com)**