

## Definition :-

### Terms Involved

- **Population** :- It Is A Representative Sample Of A Larger Group Of People (Or Even Things) With One Or More Characteristics In Common.
- **Sample** :- It Is An Analytic Subset Of A Larger Population
- **Parameter** :- It Is A Number Describing A Whole Population (E.G., Population Mean), While A **Statistic** Is A Number Describing A Sample (E.G., Sample Mean).
- **Statistics** :- Branch Of Mathematics Dealing With The Collection, Analysis, Interpretation, And Presentation Of Masses Of Numerical Data

### Measures of central tendency – 3M

<u>Property</u>	<u>Formula</u>	<u>Major Points</u>
Mean	$\bar{x} = \Sigma x / n$	Sensitive to Outliers. More useful when data is symmetrical.
Median	For n = odd; Median = $(n + 1) / 2$ For n = even; Median = Avg. of $n/2$ and $(n + 1) / 2$	Robust to outliers, useful when the data is skewed
Mode	Highest repeating value in dataset	Finds central value. Helpful in categorical features

### Measures of Variation

<u>Property</u>	<u>Formula</u>	<u>Major Points</u>
Sample Variance	$S^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1}$	
Sample Standard Deviation	$s = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1}}$	Square root of variance. Robust to outliers. Commonly used
IQR	$IQR = Q3 - Q1$	Less sensitive to outliers
Range	Highest Value - Lowest Value	Highly sensitive to unusual values. Not used often

This file is meant for personal use by anupriyambtech@gmail.com only.

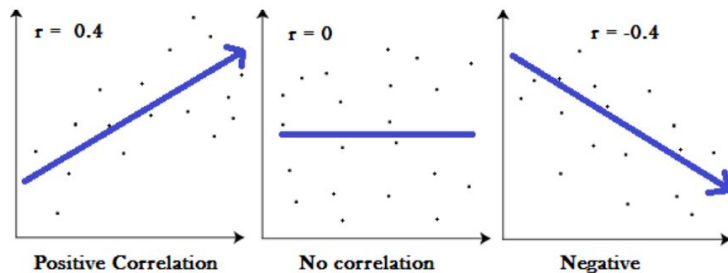
Proprietary content. ©Great Learning. All Rights Reserved. Use or full scale distribution prohibited.

## Measures of Relative Position

Property	Formula	Major Points
Percentile	Data is divided into 100 equal parts by increasing order	For applying normal distributions
Quartile	Data is divided into 4 equal parts. For ex. Q3(third is the value greater than $\frac{3}{4}$ of others.	Used to compute IQR
Z-Score	Data is divided into 4 equal parts. For ex. Q3(third is the value greater than $\frac{3}{4}$ of others.	Measures distance from mean in terms of standard deviation

**Correlation Coefficient - Correlation** coefficient formulas are used to find how strong a relationship is between data. The formulas return a value between -1 and 1, where:

- 1 indicates a strong positive relationship.
- -1 indicates a strong negative relationship.
- A result of zero indicates no relationship at all.



**Person Correlation –**

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

$r$  = correlation coefficient

$x_i$  = values of the x-variable in a sample

$\bar{x}$  = mean of the values of the x-variable

$y_i$  = values of the y-variable in a sample

$\bar{y}$  = mean of the values of the y-variable

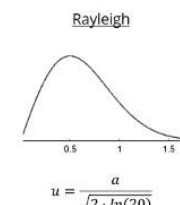
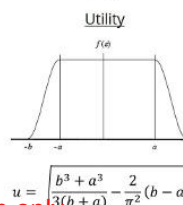
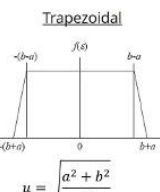
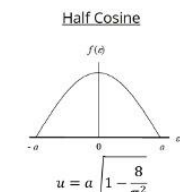
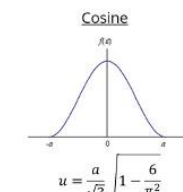
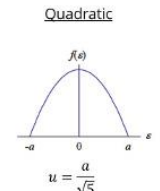
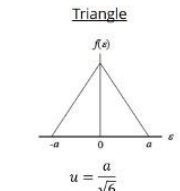
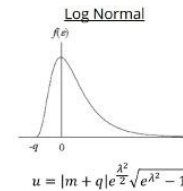
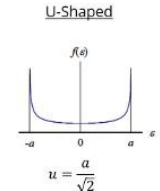
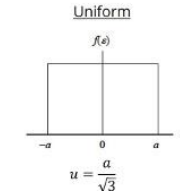
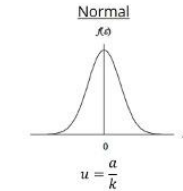
# Applied Statistics – Week 2 - Probability

- Population - It is a representative sample of a larger group of people (or even things) with one or more characteristics in common.
- Sample - It is an analytic subset of a larger population.
- Event - A probability event can be defined as a set of outcomes of an experiment.

## Rules for Probability Distribution

Rule	Formula
Complement Rule	$P(A) = 1 - P(A')$
Multiplication Rule (Dependant)	$P(A \cap B) = P(A) * P(B   A)$
Multiplication Rule (Independent)	$P(A \cap B) = P(A) * P(B)$
Addition Rule – Not Mutually Exclusive	$P(A \cup B) = P(A) + P(B) - P(A \cap B)$
Mutually Exclusive	$P(A \cup B) = P(A) + P(B)$
Conditional Probability	$P(A   B) = P(A \cap B) / P(B)$
Bayes Theorem	$P(A   B) = P(B   A) * P(A) / P(B)$

## Rules for Probability Distribution



This file is meant for personal use by anupriyambtech@gmail.com only.

## Poisson Distribution

The Poisson distribution is a probability distribution that represents the number of times an event occurs in a fixed time and/or space interval and is defined by parameter  $\lambda$  (lambda).

Examples of events that can be described by the Poisson distribution include the number of bikes crossing an intersection in a specific hour and the number of meteors seen in a minute of a meteor shower.

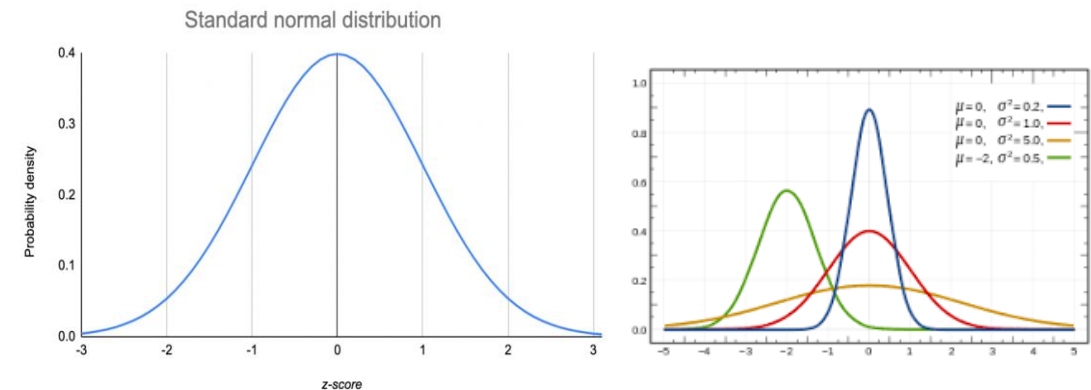
Mathematically:

$$X \sim \text{Binomial}(n, p), E(X) = n \times p$$

$$Y \sim \text{Poisson}(\lambda), E(Y) = \lambda$$

## Normal Distribution

It is a probability density function that looks like this where  $\sigma = 0$ ,  $\mu = 1$



# Applied Statistics – Week 3 – Hypothesis Testing

Hypothesis Term	Definition
Significance Level ( $\alpha$ )	Defines the strength of evidence in probabilistic terms. Specifically, alpha represents the probability that tests will produce statistically significant results when the null hypothesis is correct
Confidence Level (C)	The percentage of all possible samples that can be expected to include the true population parameter.
Critical Value( $z_c$ )	$z_c$ is the critical value of a standard normal distribution under $H_0$ . Critical values divide the rejection and non-rejection regions. Set using p-values or to a threshold value of 0.05 (5%) or 0.01 (1%), but always $\leq 0.10$ (10%)
Test Statistic ( $z_{data}$ )	$z_{data}$ is the test value of $z$ of a standard normal distribution under $H_0$ . If $z_{data}$ is inside the rejection region, demarked by $z_c$ , then we can reject the null hypothesis, $H_0$ .
p-value	Probability of obtaining a sample “more extreme” than the ones observed in your data, assuming $H_0$ is true.
Hypothesis	A premise or claim that we want to test.
Null Hypothesis: $H_0$	Currently accepted value for a parameter. (middle of the distribution) Is assumed true for the purpose of carrying out the hypothesis test. • Always contains an “=” {=, $\leq$ , $\geq$ } • The null value implies a specific sampling distribution for the test statistic • $H_0$ is the middle of the normal distribution curve at $z = 0$ . • Can be rejected, or not rejected, but NEVER supported
Alternative Hypotheses: $H_a$	Also called Research Hypothesis or $H_1$ . Is the opposite of $H_0$ and involves the claim to be tested. Is supported only by carrying out the test if the null hypothesis can be rejected. <ul style="list-style-type: none"> <li>• Always contains “&gt;” (right-tailed), “),” “&lt;” Left tailed or “!=” two tailed</li> <li>• Can be Supported (By rejecting the null), or not supported (by failing or rejecting the null) but never rejected</li> </ul>

# Applied Statistics – Week 3 – Hypothesis Testing

Hypothesis Testing	Steps
Hypothesis Testing	<ol style="list-style-type: none"> <li>1) Formula <math>H_0</math> And Has</li> <li>2) Graph: Sketch And Label Critical Value (Left-tailed, Right-tailed, Two-tailed)</li> <li>3) Decision Rule: Use Significance Level (<math>\alpha</math>), Confidence Level (C), Confidence Interval, Or Critical Value (<math>z_c</math>). E.G. We Will Reject <math>H_0</math> If <math>z_{data} &gt; 1.645</math></li> <li>4) Critical Value: Determine Critical Values (<math>z_c</math>) To Mark The Rejection Regions</li> <li>5) Test Statistic: Calculate The Test Statistic (<math>z_{data}</math>) From The Sample Data</li> <li>6) Conclusion: Reject The Null Hypothesis (Supporting The Alternative Hypothesis) Otherwise Fail To Reject The Null Hypothesis, Then State Claim</li> </ol>

## Hypothesis Formulation

If Claim Consist Of	Then Hypothesis Test Is	Represented By
“Is Equal To”, “Is Exactly”, “Is The Same As”, “Is Between” “Is At Least” “Is At Most”	Two-tailed = Left-tailed $\leq$ Right-tailed $\geq$	$H_0$
“Is Not Equal To”, “Is Different From”, “Has Changed From” “Is Less Than”, “Is Below”, “Is Lower Or Smaller Than” “Is Greater Than”, “Is Above”, “Is Longer Or Bigger Than” <i>Make sure <math>H_0 + H_a = \text{all possible outcome}</math>.</i>	Two-tailed $\neq$ Left-tailed $<$ Right-tailed $>$	$H_a$

## Decision Rule

P-value	Use Probability Value To Determine Z In Normal Distribution Table
Significance Level $\alpha$	Usually At A Threshold Value Of 5% Or 1% But Always $\leq 10\%$ $\alpha = 1 - C$
Confidence Level (C)	With A Confidence Of 0.95 (95%) Or 0.99 (99%), But Always $\geq 0.90$ (90%). $C = 1 - \alpha$
Examples:	E.G. We Will Reject $H_0$ If Significance Level Is Less Than 5% E.G. We Will Reject $H_0$ If Confidence Level Is Greater Than 95% E.G. We Will Reject $H_0$ If Confidence Interval Is Between 5% And 95% (e. g. $\pm 5\%$ ) E.G. We Will Reject $H_0$ If $z_{data} > z_c$ In A Right-tailed Test

## Determine Critical Value

Critical Values ( $z_c$ )	Determine $z_c$ by looking up $\alpha$ , C, or p-values in a standard normal distribution table. Two-tailed tests have two values for $z_c$ .
---------------------------	---

This file is meant for personal use by anupriyambtech@gmail.com only.

# Applied Statistics – Week 3 – Hypothesis Testing

Type	Description
Test For Independence	Tests For The Independence Of Two Categorical Variables
Homogeneity Of Variance	Test If More Than Two Subgroups Of A Population Share The Same Multivariate Distribution
Goodness Of Fit	Whether A Multinomial Model For The Population Distribution ( $P_1, \dots, P_m$ ) Fits Our Data
<b>Assumptions</b> <ol style="list-style-type: none"> <li>1. One Or Two Categorical Variables</li> <li>2. Independent Observations</li> <li>3. Outcomes Mutually Exclusive</li> <li>4. Large N And No More Than 20% Of Expected Counts <math>&lt; 5</math></li> </ol>	
Anova Analysis	Comparing The Means Of Two Or More Continuous Populations
One-way Layout	A Test That Allows One To Make Comparisons Between The Means Of Two Or More Groups Of Data.
Two-way Layout	A Test That Allows One To Make Comparisons Between The Means Of Two Or More Groups Of Data, Where Two Independent Variables Are Considered.
Assumptions About Data: <ol style="list-style-type: none"> <li>1. Each Data Y Is Normally Distributed</li> <li>2. The Variance Of Each Treatment Group Is Same</li> <li>3. All Observations Are Independent</li> </ol>	

This file is meant for personal use by anupriyambtech@gmail.com only.

Proprietary content. ©Great Learning. All Rights Reserved. Sharing or publishing the contents in part or full is liable for legal action. Distribution prohibited.