

# Statistics

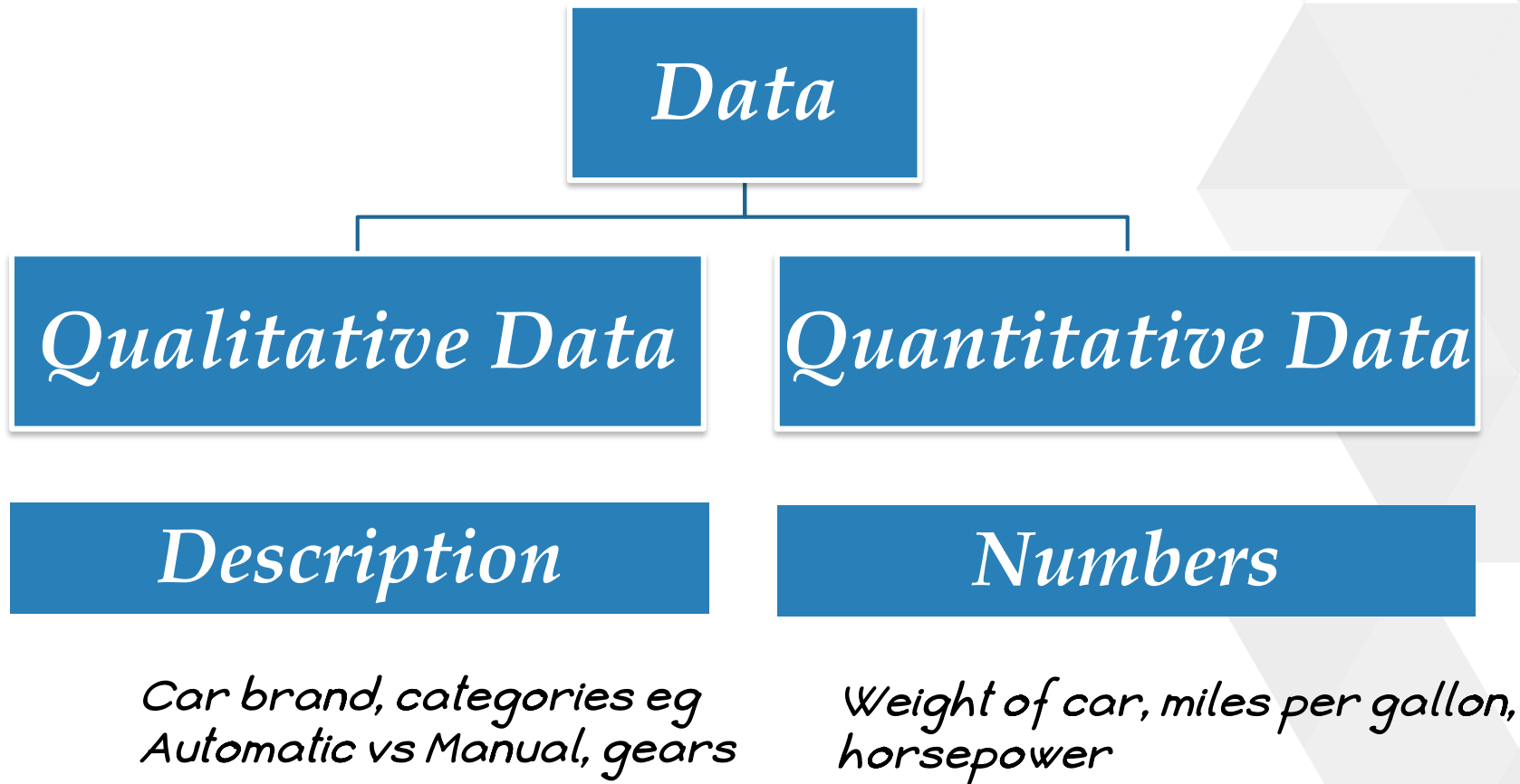
---

- ❖ **Statistics** is a branch of mathematics dealing with the collection, analysis, interpretation, presentation, and organization of data. (Wikipedia)
- ❖ Statistics converts data to information.

# Descriptive vs Inferential Statistics

- ❖ **Descriptive Statistics:** Summary of the data (population)
  - ❖ Dataset: State.X77 - Average population
  
- ❖ **Inferential Statistics:** Based on random samples from population make inference about the population.
  - ❖ Dataset: Chickwt - Effect of feed on chicken weight based on samples.

# Qualitative vs Quantitative



# Qualitative vs Quantitative

*Data*

```
graph TD; Data[Data] --> Continuous[Continuous Data]; Data --> Discrete[Discrete Data]; Continuous --> Infinite[Infinite categories]; Discrete --> Finite[Finite categories];
```

*Continuous Data*

*Discrete Data*

*Infinite categories*

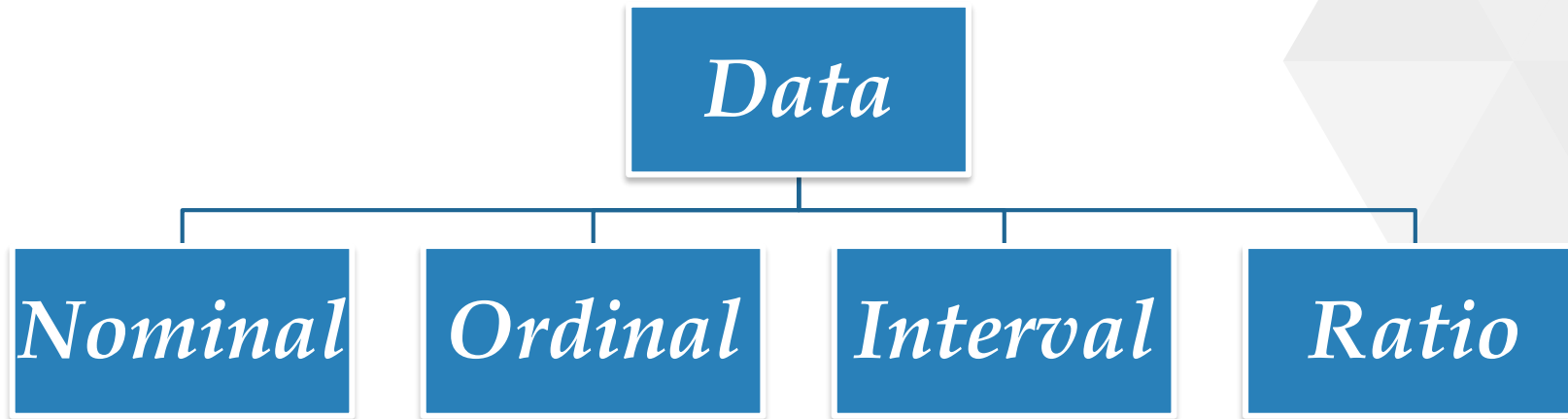
*Finite categories*

*Weight of car, miles per gallon,  
horsepower*

*Number of gears, Automatics vs Manual*

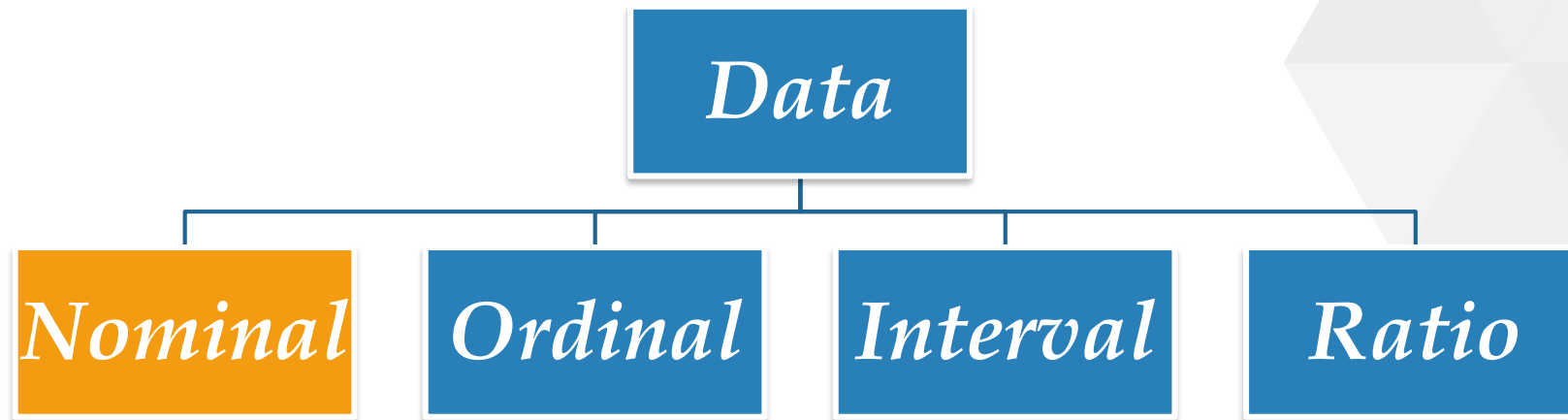
# Measurement Scales

---



# Measurement Scales

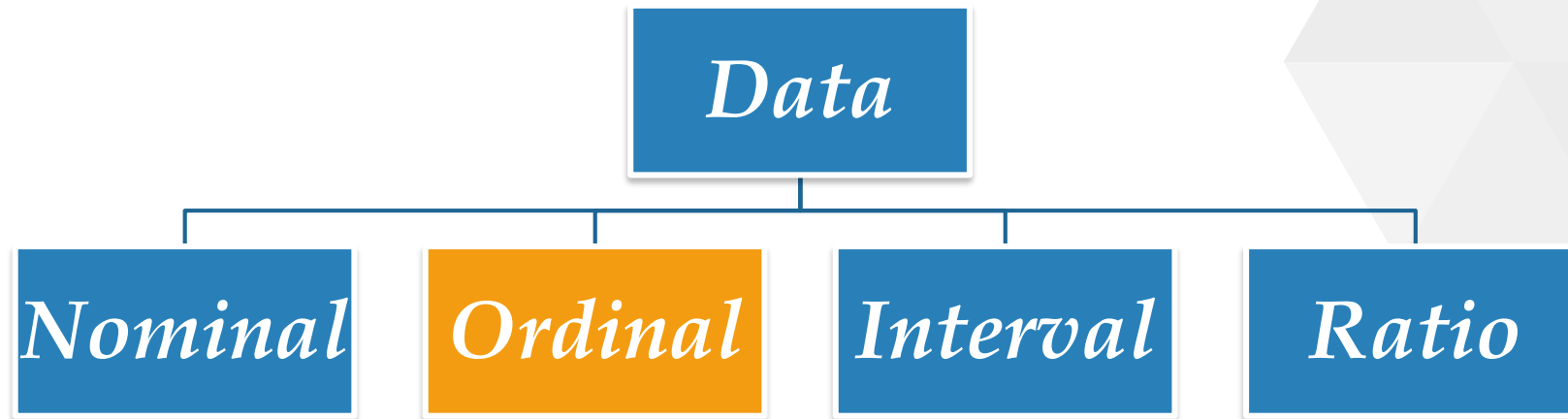
---



*Example:*  
*Color: Blue, Green, Red*

# Measurement Scales

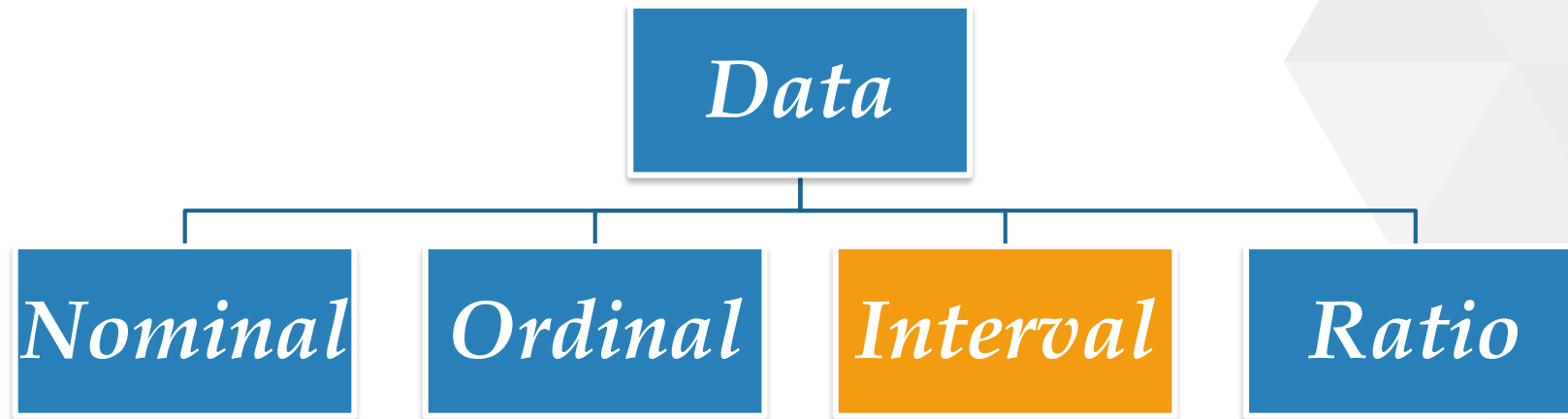
---



*Example:*  
*Pass/Fail*  
*Good, Bad, Worst*

# Measurement Scales

---

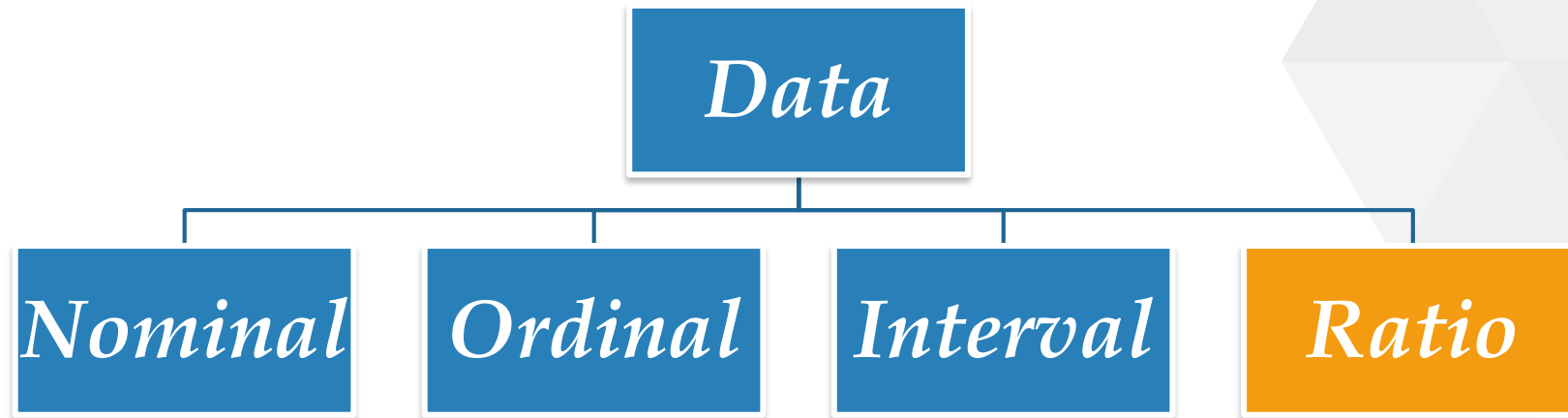


*Example:  
Temperature: Celsius*



# Measurement Scales

---



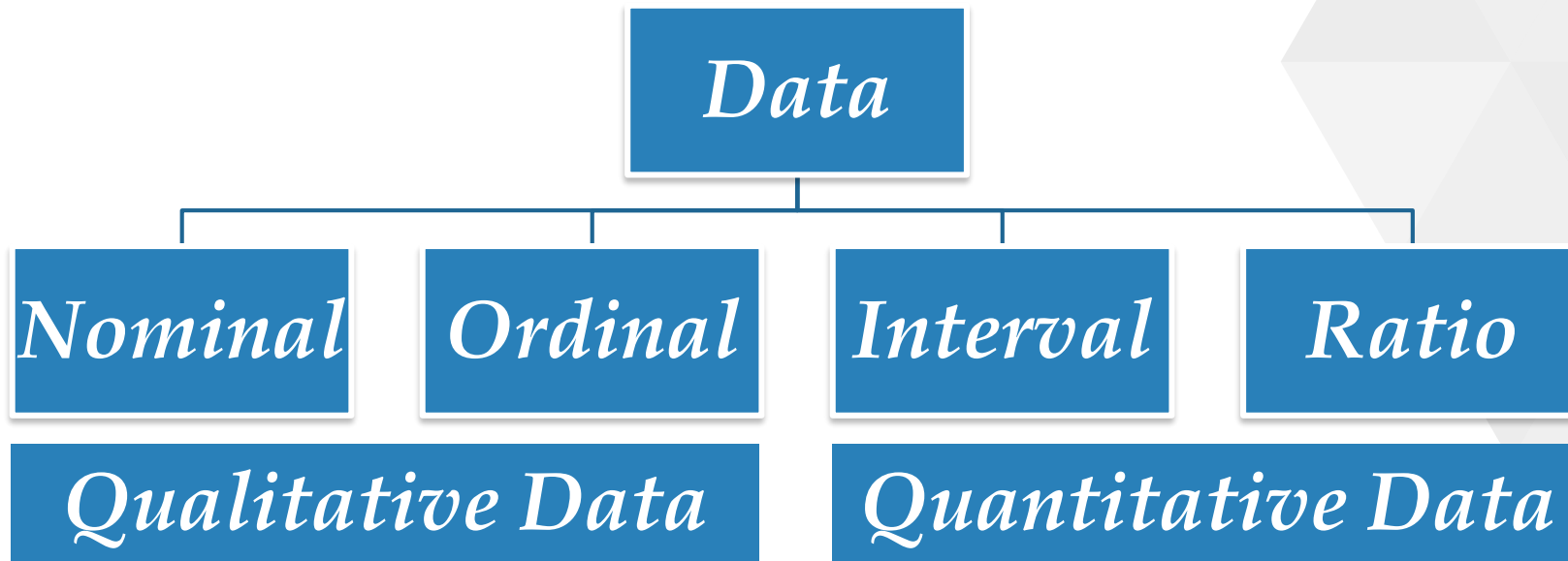
*Example:  
Height, mass, volume*

# Measurement Scales

---

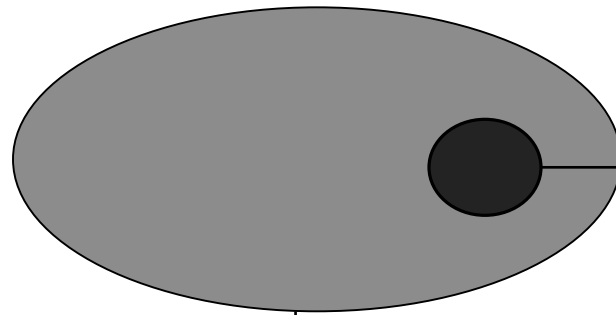
	Nominal	Ordinal	Interval	Ratio
Ordered	N	Y	Y	Y
Difference	N	N	Y	Y
Absolute Zero	N	N	N	Y
Example	Red, Blue	Good, Bad, Worst	Temperature : Degree C	Length, Weight
Central Tendency Measurement	Mode	Mode, Median	Mode, Median, Mean	Mode, Median, Mean

# Qualitative vs Quantitative



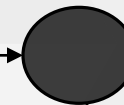
# Basic Statistical Terms

**Population:**  
Complete  
collection to  
be studied



Sampling  
Process

**Sample:** Part  
of population



Parameter

Characteristic of  
a population

$N$

number of members

$\mu$

mean

$\sigma$

standard deviation

Inference

Statistic

Characteristic  
of a sample

$n$

$\bar{x}$

$s$

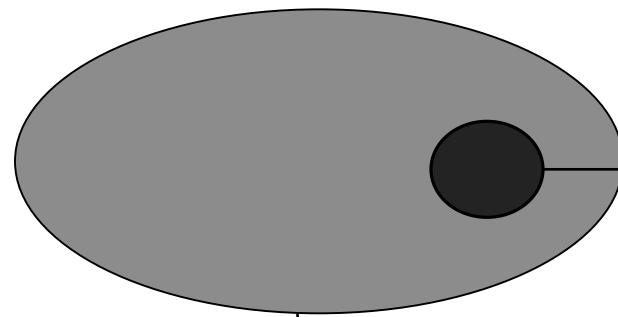
# Notations

---

	Population Parameters	Sample Statistics
Mean	$\mu$	$\bar{x}$
Standard Deviation	$\sigma$	$s$
Variance	$\sigma^2$	$s^2$
Proportion of population having an attribute	$P$	$p$
Proportion of population not having an attribute	$Q$ (=1-P)	$q$ (=1-p)
Correlation coefficient	$\rho$	$r$
Number of elements	$N$	$n$

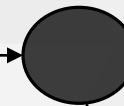
# Basic Statistical Terms

Population:  
Complete  
collection to  
be studied



Sampling  
Process

Sample: Part  
of population



Parameter

Inference

Statistic

Characteristic of  
a population

Characteristic  
of a sample

$N$  number of members  
 $\mu$  mean  
 $\sigma$  standard deviation

$n$   
 $\bar{x}$   
 $s$

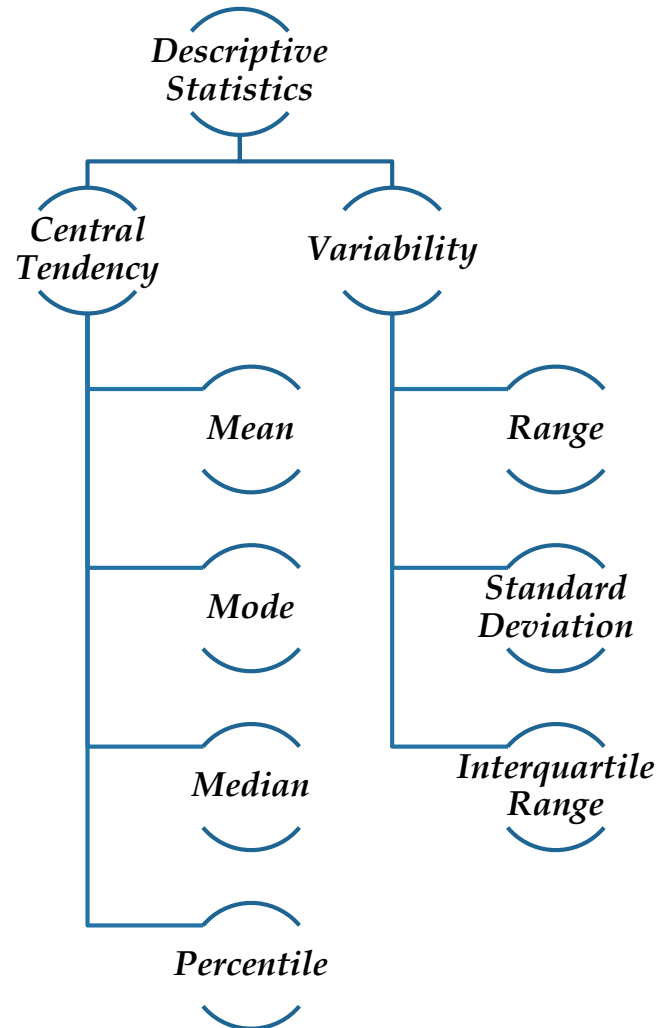
# Notations

---

	Population Parameters	Sample Statistics
Mean	$\mu$	$\bar{x}$
Standard Deviation	$\sigma$	$s$
Variance	$\sigma^2$	$s^2$
Proportion of population having an attribute	$P$	$p$
Proportion of population not having an attribute	$Q$ (=1-P)	$q$ (=1-p)
Correlation coefficient	$\rho$	$r$
Number of elements	$N$	$n$

# Descriptive Statistics

---

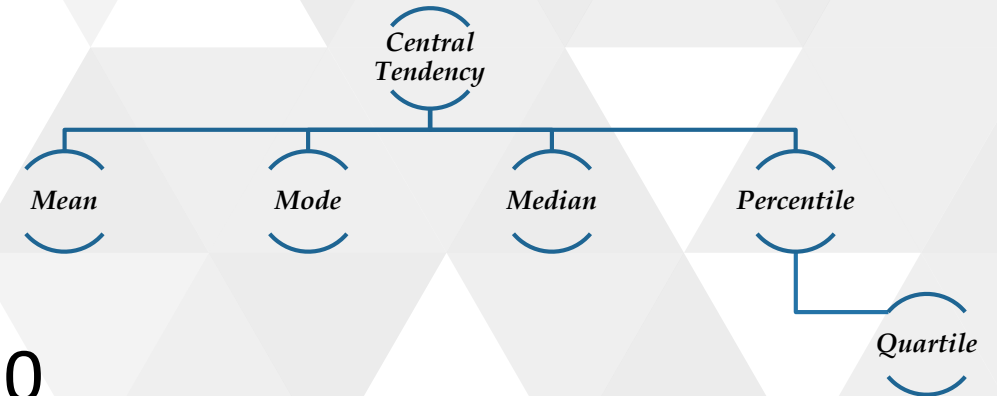




# Mean

---

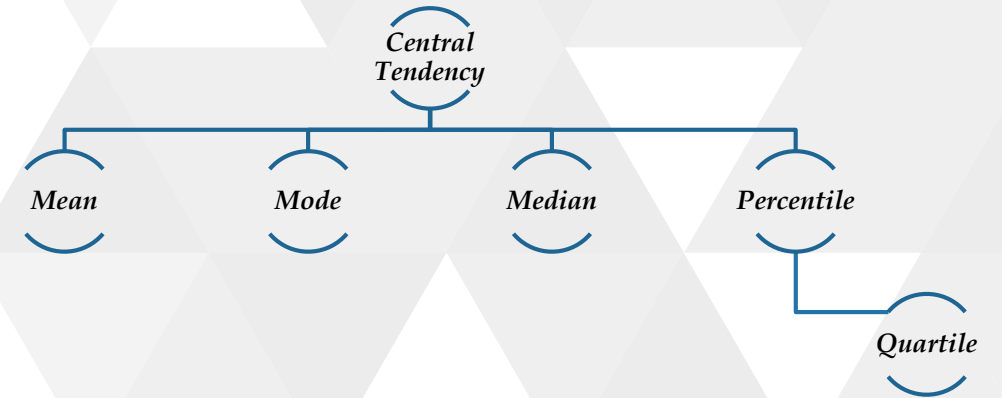
- ❖ Also known as Average
- ❖ Affected by extreme values
- ❖ Example: 10, 11, 14, 9, 6
- ❖  $\text{Mean} = (10+11+14+9+6)/5 = 50/5 = 10$



# Mode

---

- ❖ Most occurring item
- ❖ Example: 10, 11, 14, 9, 6, 10
- ❖ Mode = 10

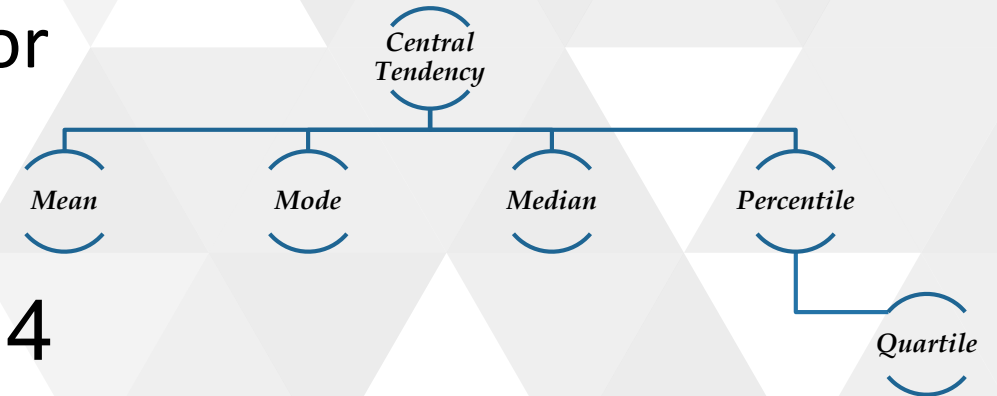


# Median

---

- ❖ Middle value when put in ascending or descending order.
- ❖ Example: 10, 11, 14, 9, 6
- ❖ In ascending order - 6,9,10,11,14
- ❖ Median = 10

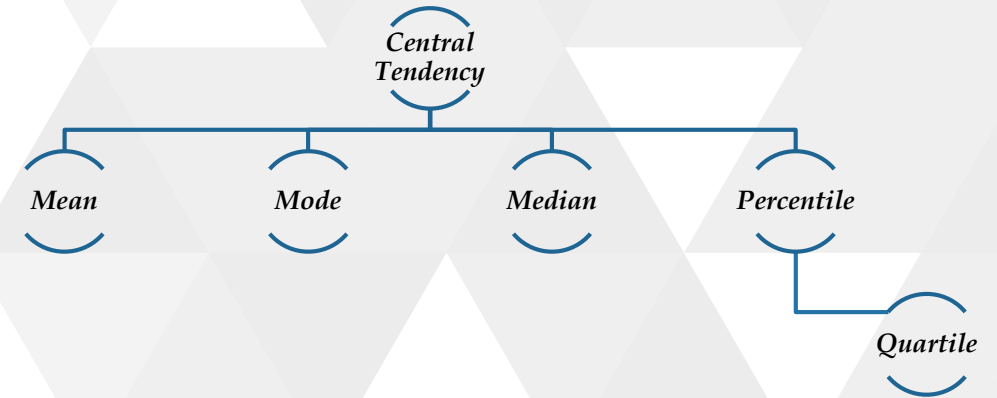
- ❖ Example: 10, 11, 14, 9, 6, 11
- ❖ In order - 6,9,10,11, 11,14
- ❖ Median = 10.5



# Percentile

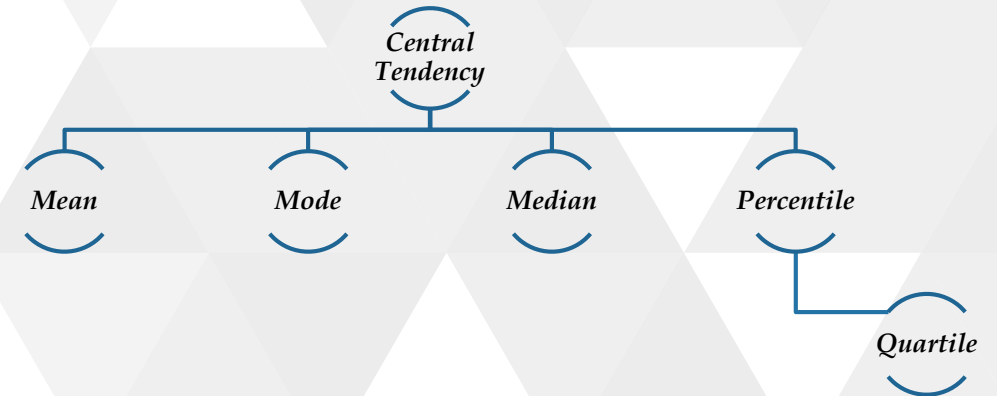
---

- ❖ Median divides the data in two equal parts when arranged in ascending or descending order
- ❖ Percentile divides data in 99 parts
- ❖ Quartile divides data in 4 parts
- ❖ Example: 6,9,10,11, 11,14
- ❖  $Q1=9$ ,  $Q2=10.5$ ,  $Q3=11$



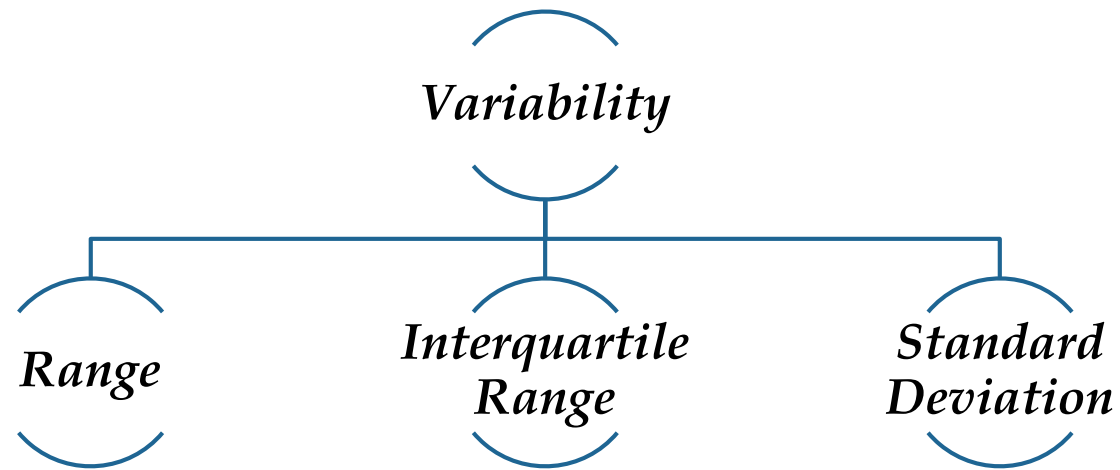
# Percentile/Quartile Steps

- ❖ Arrange in ascending or descending order
- ❖ Calculate location(i) =  $P.(n)/100$
- ❖ P=percentile, n=numbers in data set
- ❖ If i is whole number – Percentile is average of (i)th and (i+1)th location
- ❖ If i is “not” a whole number – Percentile is located at (i+1)th whole-num.
- ❖ Example: 6,9,10,11, 11,14
- ❖ Q1=9, Q2=10.5, Q3=11



# Descriptive Statistics

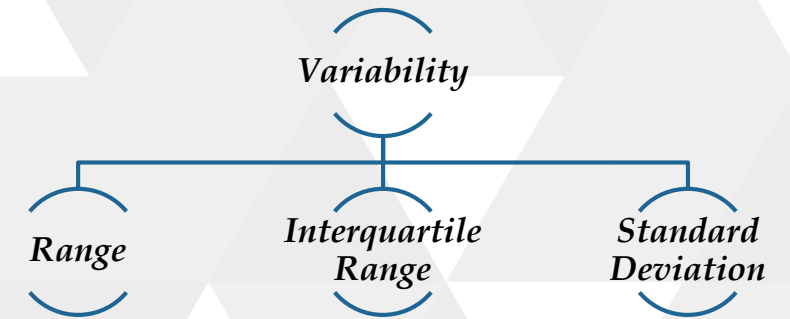
---



# Range

---

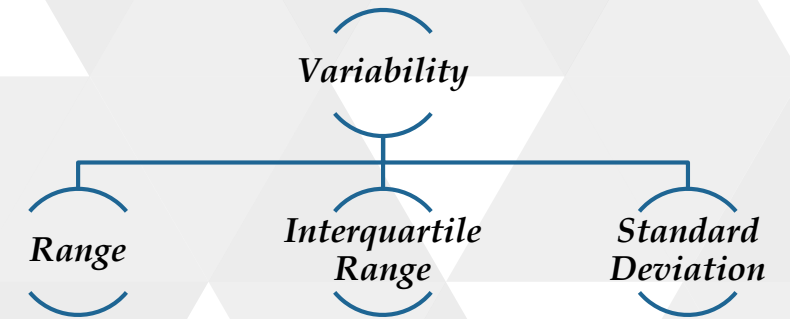
- ❖ Difference between lowest and the highest value.
- ❖ Example: 6,9,10,11, 11,14
- ❖  $\text{Range} = 14 - 6 = 8$



# Interquartile Range

---

- ❖ Range of middle 50% data
- ❖  $IQR = Q3 - Q1$
- ❖ Example: 6, 9, 10, 11, 11, 14
- ❖  $Q1 = 9$ ,  $Q2 = 10.5$ ,  $Q3 = 11$
- ❖  $IQR = 11 - 9 = 2$
- ❖ Box-and-Whisker Plot



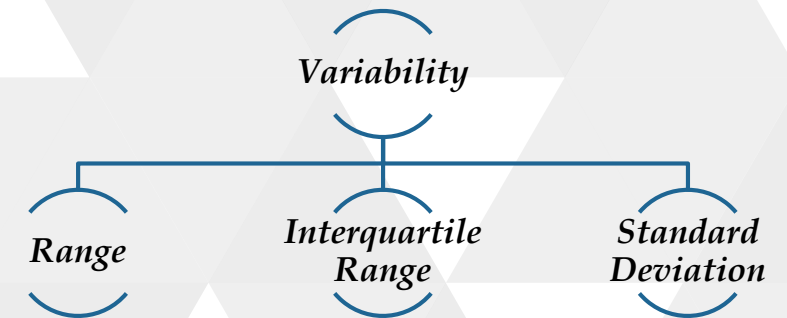


# Standard Deviation

- ❖ Variance = average of squared deviation about the arithmetic mean.
- ❖ Square root of variance is standard deviation

$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$$

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$$



# Standard Deviation

---

x	$x - \bar{x}$	$(x - \bar{x})^2$
100	0	0
101	1	1
99	-1	1
102	2	4
98	-2	4
100	0	0
$\bar{x} = 100$	$\sum(x - \bar{x}) = 0$	$\sum(x - \bar{x})^2 = 10$

$$s^2 = \frac{\sum(x - \bar{x})^2}{n-1}$$

$$S^2 = 10/5 = 2$$

$$S = \sqrt{2} = 1.414$$

# Skewness

---

## ❖ **Negative Skew or Left Skew**

- ❖ Left tail is longer

## ❖ **Positive Skew or Right Skew**

- ❖ Right tail is longer

## ❖ **Value of Skewness (Rule of Thumb)**

- ❖ Skew = 0 means perfect symmetric
- ❖ Skew between 0 and +/- 0.5 means approximately symmetric
- ❖ Skew between +/- 0.5 and 1.0 means moderately skewed
- ❖ Skew more than +1 or less than -1 means highly skewed

# Kurtosis

---

- ❖ Normal distribution has Kurtosis = 0
- ❖ Kurtosis < 0 means peak is short and broad, tails are shorter
- ❖ Kurtosis > 0 means peak is higher and thinner, tails are longer

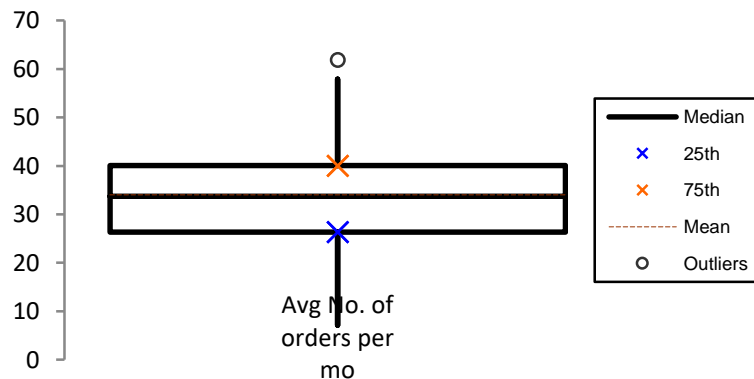
# Graphical Methods

---

- ❖ Box-and-whisker plots
- ❖ Scatter diagrams
- ❖ Histograms

# Box and Whisker Plots

- ❖ Also known as Box Plot
- ❖ Shows the median
- ❖ Shows Q1, Q3 and IQR



# Scatter Diagram

---

- ❖ One of seven basic quality tools
- ❖ To see relationship between two variables
- ❖ Relationship should make practical sense
- ❖ Temperature(X) vs Ice cream sale (Y)
- ❖ Some times relationship between two variables is because of a third variable. (ice cream sale vs heat stroke cases)
- ❖ Correlation/Regression is covered in the Analyze Phase

# Histogram

---

- ❖ Graphical representation of the distribution of numerical data
- ❖ Values are assigned “bins” and frequency for each bin is plotted.



# Histograms

---

# Scatter Diagram

---

- ❖ One of seven basic quality tools
- ❖ To see relationship between two variables
- ❖ Relationship should make practical sense
- ❖ Temperature(X) vs Ice cream sale (Y)
- ❖ Some times relationship between two variables is because of a third variable. (ice cream sale vs heat stroke cases)
- ❖ Correlation/Regression is covered in the Analyze Phase

# Probability

---

## ❖ Classic Model

Number of outcomes in which the event occurs

Total Number of possible outcomes of an experiment

# Probability

---

## ❖ Relative Frequency of Occurrence

Number of times an event occurred

---

Total number of opportunities for an event to occur

# Probability

---

- ❖ Experiment/Trial: Some thing done with an expectation of result.
- ❖ Event or Outcome: Result of experiment
- ❖ Sample Space: A sample space of an experiment is the set of all possible results of that random experiment.

$\{1, 2, 3, 4, 5, 6\}$

# Probability

---

- ❖ Union: Probability that events A or B occur:  $P(A \cup B)$
- ❖ Intersection: Probability that events A and B occur:  $P(A \cap B)$

# Probability

---

- ❖ Mutually Exclusive Events: When two events cannot occur at the same time
- ❖ Independent Events: The occurrence of Event A does not change the probability of Event B
- ❖ Complementary Events: The probability that Event A will NOT occur is denoted by  $P(A')$ .

# Probability

---

## ❖ Rule of Addition

The probability that Event A or Event B occurs

=

Probability that Event A occurs

+

Probability that Event B occurs

-

Probability that both Events A and B occur

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



# Probability

---

## ❖ Rule of Multiplication:

The probability that Events A and B both occur

=

Probability that Event A occurs

x

Probability that Event B occurs, given that A has occurred

$$P(A \cap B) = P(A) P(B | A)$$

# Probability

---

## ❖ Independent Events

# Probability

---

## ❖ Dependent Events

# Factorial

---

- ❖ **Factorial** of a non-negative integer  $n$ , denoted by  $n!$ , is the product of all positive integers less than or equal to  $n$

# Permutation/ Combination

- ❖ Permutation: A set of objects in which position (or order) is important.
  - ❖ e.g. Lock combination: 3376
  
- ❖ Combination: A set of objects in which position (or order) is NOT important.
  - ❖ e.g. Selecting 2 students out of 5

# Central Limit Theorem

---

# Central Limit Theorem

---

- ❖ For almost all populations, the sampling distribution of the mean can be approximated closely by a normal distribution, provided the sample size is sufficiently large.

# Central Limit Theorem

---



# Central Limit Theorem

---

# Normal Probability Distribution

- ❖ Symmetrically distributed
- ❖ Long Tails / Bell Shaped
- ❖ Mean/ Mode and Median are same

# Normal Probability Distribution

- ❖ Two factors define the shape of the curve:
  - ❖ Mean
  - ❖ Standard Deviation

# Normal Probability Distribution

- ❖ About 68% of the area under the curve falls within **1 standard deviation** of the mean.
- ❖ About 95% of the area under the curve falls within **2 standard deviations** of the mean.
- ❖ About 99.7% of the area under the curve falls within **3 standard deviations** of the mean.

# Normal Probability Distribution

- ❖ The total area under the normal curve = 1.
- ❖ The probability of any particular value is 0.
- ❖ The probability that  $X$  is greater than or less than a value = area under the normal curve in that direction

# Normal Probability Distribution

- ❖ The value of the random variable Y is:

$$Y = \left\{ \frac{1}{\sigma \cdot \sqrt{2\pi}} \right\} \cdot e^{-\frac{(x - \mu)^2}{2\sigma^2}}$$

- ❖ where X is a normal random variable,
- ❖  $\mu$  = mean,
- ❖  $\sigma$  = standard deviation,
- ❖  $\pi$  is approximately 3.14159,
- ❖ e is approximately 2.71828.

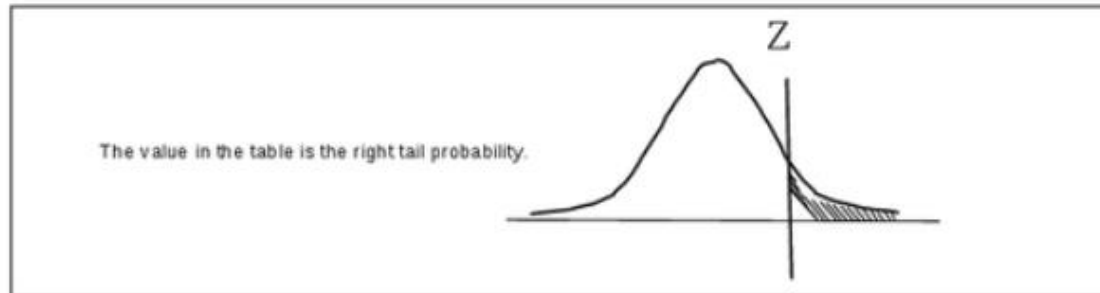


# Normal Probability Distribution

- ❖ Z Value / Standard Score
- ❖ How many standard deviations an element is from the mean.
- ❖  $z = (X - \mu) / \sigma$
- ❖  $z$  is the z-score,
- ❖  $X$  is the value of the element,
- ❖  $\mu$  is the population mean,
- ❖  $\sigma$  is the standard deviation.



# Z Table



Hundredth place for Z-value

Z-Value	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.50000	0.49601	0.49202	0.48803	0.48405	0.48006	0.47608	0.47210	0.46812	0.46414
0.1	0.46017	0.45620	0.45224	0.44828	0.44433	0.44038	0.43644	0.43251	0.42858	0.42465
0.2	0.42074	0.41683	0.41294	0.40905	0.40517	0.40129	0.39743	0.39358	0.38974	0.38591
0.3	0.38209	0.37828	0.37448	0.37070	0.36693	0.36317	0.35942	0.35569	0.35197	0.34827
0.4	0.34458	0.34090	0.33724	0.33360	0.32997	0.32636	0.32276	0.31918	0.31561	0.31207
0.5	0.30854	0.30503	0.30153	0.29806	0.29460	0.29116	0.28774	0.28434	0.28096	0.27760
0.6	0.27425	0.27093	0.26763	0.26435	0.26109	0.25785	0.25463	0.25143	0.24825	0.24510
0.7	0.24196	0.23885	0.23576	0.23270	0.22965	0.22663	0.22363	0.22065	0.21770	0.21476
0.8	0.21186	0.20897	0.20611	0.20327	0.20045	0.19766	0.19489	0.19215	0.18943	0.18673
0.9	0.18406	0.18141	0.17879	0.17619	0.17361	0.17106	0.16853	0.16602	0.16354	0.16109
1.0	0.15866	0.15625	0.15386	0.15151	0.14917	0.14686	0.14457	0.14231	0.14007	0.13786
1.1	0.13567	0.13350	0.13136	0.12924	0.12714	0.12507	0.12302	0.12100	0.11900	0.11702
1.2	0.11507	0.11314	0.11123	0.10935	0.10749	0.10565	0.10383	0.10204	0.10027	0.09853
1.3	0.09680	0.09510	0.09342	0.09176	0.09012	0.08851	0.08691	0.08534	0.08379	0.08226
1.4	0.08076	0.07927	0.07780	0.07636	0.07493	0.07353	0.07215	0.07078	0.06944	0.06811
1.5	0.06681	0.06552	0.06426	0.06301	0.06178	0.06057	0.05938	0.05821	0.05705	0.05592
1.6	0.05480	0.05370	0.05262	0.05155	0.05050	0.04947	0.04846	0.04746	0.04648	0.04551



# Continuous Probability Distributions

- ❖ **Normal probability distribution**
- ❖ Student's t distribution
- ❖ Chi-square distribution
- ❖ F distribution

# Continuous vs Discrete Variable

- ❖ If a variable can take on any value between two specified values, it is called a **continuous variable**; otherwise, it is called a **discrete variable**.

# Discrete Probability Distributions

- ❖ **Binomial Probability Distribution**
- ❖ Bernoulli Distribution
- ❖ Hypergeometric Probability Distribution
- ❖ Geometric Distribution
- ❖ Negative Geometric Distribution
- ❖ **Poisson Probability Distribution**

# Binomial Probability Distribution

- ❖ **A binomial experiment** has the following properties:
  - ❖ The experiment consists of  $n$  repeated trials.
  - ❖ Each trial can result in just two possible outcomes. We call one of these outcomes a success and the other, a failure.
  - ❖ The probability of success, denoted by  $p$ , is the same on every trial.
  - ❖ The trials are independent; that is, the outcome on one trial does not affect the outcome on other trials.

# Binomial Probability Distribution

- ❖ **A binomial experiment** has the following properties:
  - ❖ The experiment consists of  $n$  repeated trials.
  - ❖ Each trial can result in just two possible outcomes. We call one of these outcomes a success and the other, a failure.
  - ❖ The probability of success, denoted by  $p$ , is the same on every trial.
  - ❖ The trials are independent; that is, the outcome on one trial does not affect the outcome on other trials.

# Binomial Probability Distribution

❖  **$x$** : The number of successes that result from the binomial experiment.

$$P(x) = {}_nC_x \cdot p^x \cdot (1 - p)^{n - x}$$

❖  **$n$** : The number of trials in the binomial experiment.

❖  **$p$** : The probability of success on an individual trial.

❖  **$q$** : The probability of failure on an individual trial. (This is equal to  $1 - p$ .)

❖  **$n!$** : The factorial of  $n$  (also known as  $n$  factorial).

❖  **$P(x)$** : Binomial probability - the probability that an  $n$ -trial binomial experiment results in exactly  $x$  successes, when the probability of success on an individual trial is  $p$ .

❖  **${}_nC_x$** : The number of combinations of  $n$  things, taken  $x$  at a time.

# Binomial Probability Distribution

- ❖ The **binomial probability** refers to the probability that a binomial experiment results in exactly  $x$  successes.
- ❖ Suppose a binomial experiment consists of  $n$  trials and results in  $x$  successes. If the probability of success on an individual trial is  $p$ , then the binomial probability is:
- ❖  $P(x) = {}_n C_x \cdot p^x \cdot (1 - p)^{n - x}$   
or  
 $P(x) = \{ n! / [ x! (n - x)! ] \} \cdot p^x \cdot (1 - p)^{n - x}$

# Binomial Probability Distribution

❖ The mean of the distribution ( $\mu_x$ ) is  
 **$n \cdot p$**

*$n$ : The number of trials in the binomial experiment.*

❖ The variance ( $\sigma^2_x$ ) is  
 **$n \cdot p \cdot (1 - p)$**

*$p$ : The probability of success on an individual trial.*

❖ The standard deviation ( $\sigma_x$ ) is  
 **$\text{sqrt}[n \cdot p \cdot (1 - p)]$**



# Five Conditions - Binomial

---

- ❖ 1. There is a fixed number,  $n$ , of identical trials.
- ❖ 2. For each trial, there are only two possible outcomes (success/failure).
- ❖ 3. The probability of success,  $p$ , remains the same for each trial.
- ❖ 4. The trials are independent of each other.
- ❖ 5.  $x$  = the number of successes observed for the  $n$  trials.

$$P(x) = {}_nC_x \cdot p^x \cdot (1 - p)^{n - x}$$

# Bernoulli Distribution

---

- ❖ Distribution of successes on a single trial.
  - ❖ What is the probability of getting head in tossing of a coin once?

# Hypergeometric Distribution

$$P(x) = {}_n C_x \cdot p^x \cdot (1 - p)^{n - x}$$

- ❖ There is a fixed number,  $n$ , of identical trials.
- ❖ For each trial, there are only two possible outcomes (success/failure).
- ~~❖ The probability of success,  $p$ , remains the same for each trial.~~
- ~~❖ The trials are independent of each other.~~
- ❖ Finite and known population without replacement.
- ❖ Number of successes in population are known
- ❖  $x$  = the number of successes observed for the  $n$  trials.

# Hypergeometric Distribution

---

- ❖ **N**: size of population
- ❖ **A**: number of successes in population
- ❖ **x**: The number of successes that result from the experiment.
- ❖ **n**: The number of trials without replacement.
- ~~❖ **p**: The probability of success on an individual trial.~~
- ~~❖ **q**: The probability of failure on an individual trial. (This is equal to  $1 - p$ .)~~
- ❖ **P(x)** : The probability that an  $n$ -trial experiment results in exactly  $x$  successes
- ❖  ${}_n\mathbf{C}_x$ : The number of combinations of  $n$  things, taken  $x$  at a time.

$$P(x) = {}_A\mathbf{C}_x \cdot {}_{N-A}\mathbf{C}_{n-x} / {}_N\mathbf{C}_n$$

# Hypergeometric Distribution

- ❖ Out of 10 people (6M, 4F), 3 people are selected without replacement. What is the probability that two of them are females?

$$P(x) = {}_A C_x \cdot {}_{N-A} C_{n-x} / {}_N C_n$$

- ❖  $P(2) = {}_4 C_2 \cdot {}_{10-4} C_{3-2} / {}_{10} C_3$

- ❖  $= {}_4 C_2 \cdot {}_6 C_1 / {}_{10} C_3$

- ❖ When sample size is less than 5% population then can use Binomial.

# Geometric Distribution

---

- ❖ Number of trials needed to get the first success.
- ❖ What is the probability that if the coin is tossed repeatedly the first head appears on 5<sup>th</sup> trial?

# Negative Binomial Distribution

- ❖ Generalization of the Geometric distribution
- ❖ Number of trials needed to get the first number of successes.
  - ❖ What is the probability that if the coin is tossed repeatedly the first third time head appears on 5<sup>th</sup> trial?
- ❖ In Binomial distribution trials are fixed, in Negative Binomial number of successes are fixed.

# Poisson Distribution

---

- ❖ A **Poisson experiment** has the following properties:
- ❖ The experiment results in outcomes that can be classified as successes or failures.
- ❖ The average number of successes ( $\mu$ ) that occurs in a specified region is known.
- ❖ Outcomes are random. Occurrence of one outcome does not influence the chance of another outcome of interest.
- ❖ The outcomes of interest are rare relative to the possible outcomes.



# Poisson Distribution

---

- ❖  $e$ : A constant equal to approximately 2.71828. (Actually,  $e$  is the base of the natural logarithm system)
- ❖  $\mu$ : The mean number of successes that occur in a specified region.
- ❖  $x$ : The actual number of successes that occur in a specified region.
- ❖  $P(x; \mu)$ : The **Poisson probability** that exactly  $x$  successes occur in a Poisson experiment, when the mean number of successes is  $\mu$ .

$$P(x; \mu) = (e^{-\mu}) (\mu^x) / x!$$

# Poisson Distribution

---

- ❖ **Poisson Formula.** Suppose we conduct a Poisson experiment, in which the average number of successes within a given region is  $\mu$ . Then, the Poisson probability is:
- ❖  $P(x; \mu) = (e^{-\mu}) (\mu^x) / x!$
- ❖ where  $x$  is the actual number of successes that result from the experiment, and  $e$  is approximately equal to 2.71828.

# Poisson Distribution

---

- ❖ The Poisson distribution has the following properties:
- ❖ The mean of the distribution is equal to  $\mu$ .
- ❖ The variance is also equal to  $\mu$ .

# Poisson Distribution

---

- ❖ On a booking counter on the average 3.6 people come every 10 minute on weekends. What is the probability of getting 7 people in 10 minutes?
- ❖  $\mu = 3.6, x=7$
- ❖  $P(x; \mu) = (e^{-\mu}) (\mu^x) / x! = (e^{-3.6}) (3.6^7) / 7!$
- ❖  $= 0.02732 \times 7836.41 / 5040 = 0.0424$

# Errors of Statistical Tests

		True State of Nature	
		$H_0$ Is true	$H_a$ Is true
Conclusion	Support $H_0$ / Reject $H_a$	Correct Conclusion	Type II Error
	Support $H_a$ / Reject $H_0$	Type I Error	Correct Conclusion (Power)

# Errors of Statistical Tests

	Type I error (alpha)	Type II error (beta)
Name	Producer's risk/ Significance level	Consumer's risk
1 minus error is called	Confidence level	Power of the test
Example of Fire Alarm	False fire alarm leading to inconvenience	Missed fire leading to disaster
Effects on process	Unnecessary cost increase due to frequent changes	Defects may be produced
Control method	Usually fixed at a pre-determined level, 1%, 5% or 10%	Usually controlled to < 10% by appropriate sample size
Simple definition	Innocent declared as guilty	Guilty declared as innocent

# Significance Level

---

Level of Confidence / Confidence Interval:

$C = 0.90, 0.95, 0.99$  (90%, 95%, 99%)

Level of Significance:

$\alpha = 1 - C$  (0.10, 0.05, 0.01)

# Power

---

- ❖ Power =  $1 - \beta$  (or 1 - type II error)
- ❖ Type II Error: Failing to reject null hypothesis when null hypothesis is false.
- ❖ Power: Likelihood of rejecting null hypothesis when null hypothesis is false.
- ❖ Or: Power is the ability of a test to correctly reject the null hypothesis.



# Alpha vs Beta

---

- ❖ Researcher can not commit both Type I and II error. Only one can be committed.
- ❖ As the value of  $\alpha$  increases (say 0.01 to 0.05)  $\beta$  goes down and the Power of test increases.
- ❖ To reduce both Type I and II errors increase sample size.

# Hypothesis Testing

---

1. State the Alternate Hypothesis.
2. State the Null Hypothesis.
3. Select a probability of error level (alpha level). Generally 0.05
4. Select and compute the test statistic (e.g t or z score)
5. Critical test statistic
6. Interpret the results.

# Hypothesis Testing

---

## ❖ Lower Tail Tests

❖  $H_0: \mu \geq 150\text{cc}$

❖  $H_a: \mu < 150\text{cc}$

## ❖ Upper Tail Tests

❖  $H_0: \mu \leq 150\text{cc}$

❖  $H_a: \mu > 150\text{cc}$

# Hypothesis Testing

---

- ❖ Two Tail Tests
  - ❖  $H_0: \mu = 150\text{cc}$
  - ❖  $H_a: \mu \neq 150\text{cc}$

# Calculate Test Statistic

---

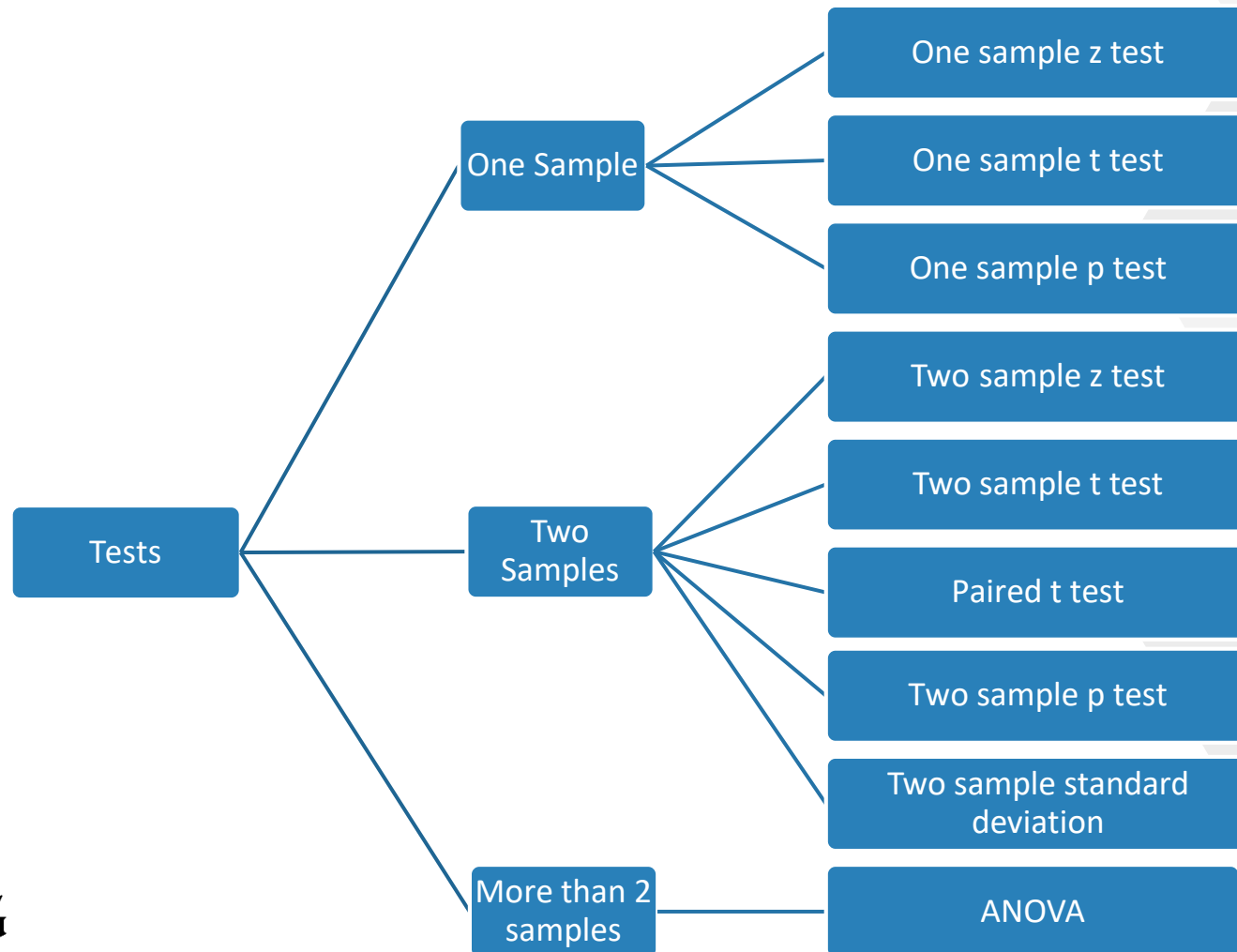
❖ Single sample

❖  $z = (x - \mu) / \sigma$

❖ Mean of Multiple samples

❖  $z = (\bar{x} - \mu) / (\sigma / \sqrt{n})$

# Tests for Mean, Variance & Proportion



# One Sample z Test

---

- ❖ Calculated value
- ❖  $z = [\bar{x} - \mu] / [\sigma / \text{sqrt}(n)]$
- ❖ Example: Perfume bottle producing 150cc with sd of 2 cc, 100 bottles are randomly picked and the average volume was found to be 152cc. Has mean volume changed? (95% confidence)
- ❖  $z_{\text{calculated}} = (152-150)/[2 / \text{sqrt}(100)] = 2/0.2 = 10$
- ❖  $z_{\text{critical}} = ?$

# One Sample z Test

$$z_{critical} = 1.96$$

z	0	1	2	3	4	5	6	7	8	9
0.0	.5000	.4960	.4920	.4880	.4840	.4801	.4761	.4721	.4681	.4641
0.1	.4602	.4562	.4522	.4483	.4443	.4404	.4364	.4325	.4286	.4247
0.2	.4207	.4168	.4129	.4090	.4052	.4013	.3974	.3936	.3897	.3859
0.3	.3821	.3783	.3745	.3707	.3669	.3632	.3594	.3557	.3520	.3483
0.4	.3446	.3409	.3372	.3336	.3300	.3264	.3228	.3192	.3156	.3121
0.5	.3085	.3050	.3015	.2981	.2946	.2912	.2877	.2843	.2810	.2776
0.6	.2743	.2709	.2676	.2643	.2611	.2578	.2546	.2514	.2483	.2451
0.7	.2420	.2389	.2358	.2327	.2296	.2266	.2236	.2206	.2177	.2148
0.8	.2119	.2090	.2061	.2033	.2005	.1977	.1949	.1922	.1894	.1867
0.9	.1841	.1814	.1788	.1762	.1736	.1711	.1685	.1660	.1635	.1611
1.0	.1587	.1562	.1539	.1515	.1492	.1469	.1446	.1423	.1401	.1379
1.1	.1357	.1335	.1314	.1292	.1271	.1251	.1230	.1210	.1190	.1170
1.2	.1151	.1131	.1112	.1093	.1075	.1056	.1038	.1020	.1003	.0985
1.3	.0968	.0951	.0934	.0918	.0901	.0885	.0869	.0853	.0838	.0823
1.4	.0808	.0793	.0778	.0764	.0749	.0735	.0721	.0708	.0694	.0681
1.5	.0668	.0655	.0643	.0630	.0618	.0606	.0594	.0582	.0571	.0559
1.6	.0548	.0537	.0526	.0516	.0505	.0495	.0485	.0475	.0465	.0455
1.7	.0446	.0436	.0427	.0418	.0409	.0401	.0392	.0384	.0375	.0367
1.8	.0359	.0351	.0344	.0336	.0329	.0322	.0314	.0307	.0301	.0294
1.9	.0287	.0281	.0274	.0268	.0262	.0256	.0250	.0244	.0239	.0233
2.0	.0228	.0222	.0217	.0212	.0207	.0202	.0197	.0192	.0188	.0183
2.1	.0179	.0174	.0170	.0166	.0162	.0158	.0154	.0150	.0146	.0143
2.2	.0139	.0136	.0132	.0129	.0125	.0122	.0119	.0116	.0113	.0110
2.3	.0107	.0104	.0102	.0099	.0096	.0094	.0091	.0089	.0087	.0084
2.4	.0082	.0080	.0078	.0075	.0073	.0071	.0069	.0068	.0066	.0064
2.5	.0062	.0060	.0059	.0057	.0055	.0054	.0052	.0051	.0049	.0048
2.6	.0047	.0045	.0044	.0043	.0041	.0040	.0039	.0038	.0037	.0036
2.7	.0035	.0034	.0033	.0032	.0031	.0030	.0029	.0028	.0027	.0026
2.8	.0026	.0025	.0024	.0023	.0023	.0022	.0021	.0021	.0020	.0019
2.9	.0019	.0018	.0018	.0017	.0016	.0016	.0015	.0015	.0014	.0014
3.0	.0013	.0013	.0013	.0012	.0012	.0011	.0011	.0011	.0010	.0010
3.1	.0010	.0009	.0009	.0009	.0008	.0008	.0008	.0008	.0007	.0007
3.2	.0007	.0007	.0006	.0006	.0006	.0006	.0006	.0005	.0005	.0005
3.3	.0005	.0005	.0005	.0004	.0004	.0004	.0004	.0004	.0004	.0003
3.4	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0002
3.5	.0002	.0002	.0002	.0002	.0002	.0002	.0002	.0002	.0002	.0002
3.6	.0002	.0002	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0001
3.7	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0001
3.8	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0001
3.9	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000



# One Sample z Test

---

- ❖ Calculated value
- ❖  $z = [\bar{x} - \mu] / [\sigma / \text{sqrt}(n)]$
- ❖ Example: Perfume bottle producing 150cc with sd of 2 cc, 100 bottles are randomly picked and the average volume was found to be 152cc. Has mean volume changed? (95% confidence)
- ❖  $z_{\text{calculated}} = (152-150)/[2 / \text{sqrt}(100)] = 2/0.2 = 10$
- ❖  $z_{\text{critical}} = 1.96 > \text{Reject } H_0$

# p Value

---

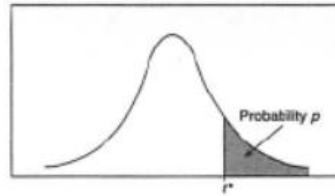
- ❖ p value is the lowest value of alpha for which the null hypothesis can be rejected. (Probability that the null hypothesis is correct)
- ❖ If  $p = 0.01$  you can reject the null hypothesis at  $\alpha = 0.05$
- ❖ p is low the null must go / p is high the null fly.

# One Sample t Test

---

- ❖ Calculated value
- ❖  $t = [\bar{x} - \mu] / [s / \text{sqrt}(n)]$
- ❖ Example: Perfume bottle producing 150cc, 4 bottles are randomly picked and the average volume was found to be 151cc and sd of sample was 2 cc. Has mean volume changed? (95% confidence)
- ❖  $t_{\text{cal}} = (151-150)/[2 / \text{sqrt}(4)] = 1/1 = 1$
- ❖  $t_{\text{critical}} = ?$

# One Sample t Test



$$t_{critical} = 3.182$$

TABLE B t Distribution Critical Values

df	TAIL PROBABILITY P											
	.25	.20	.15	.10	.05	.025	.02	.01	.005	.0025	.001	.0005
1	1.000	1.376	1.963	3.078	6.314	12.71	15.89	31.82	63.66	127.3	318.3	636.6
2	.816	1.061	1.386	1.886	2.920	4.303	4.849	6.965	9.925	14.09	22.33	31.60
3	.765	.978	1.250	1.638	2.353	3.182	3.482	4.541	5.841	7.453	10.21	12.92
4	.741	.941	1.190	1.533	2.132	2.776	2.999	3.747	4.604	5.598	7.173	8.610
5	.727	.920	1.156	1.476	2.015	2.571	2.757	3.365	4.032	4.773	5.893	6.869
6	.718	.906	1.134	1.440	1.943	2.447	2.612	3.143	3.707	4.317	5.208	5.959
7	.711	.896	1.119	1.415	1.895	2.365	2.517	2.998	3.499	4.029	4.785	5.408
8	.706	.889	1.108	1.397	1.860	2.306	2.449	2.896	3.355	3.833	4.501	5.041
9	.703	.883	1.100	1.383	1.833	2.262	2.398	2.821	3.250	3.690	4.297	4.781
10	.700	.879	1.093	1.372	1.812	2.228	2.359	2.764	3.169	3.581	4.144	4.587
11	.697	.876	1.088	1.363	1.796	2.201	2.328	2.718	3.106	3.497	4.025	4.437
12	.695	.873	1.083	1.356	1.782	2.179	2.303	2.681	3.055	3.428	3.930	4.318
13	.694	.870	1.079	1.350	1.771	2.160	2.282	2.650	3.012	3.372	3.852	4.221
14	.692	.868	1.076	1.345	1.761	2.145	2.264	2.624	2.977	3.326	3.787	4.140
15	.691	.866	1.074	1.341	1.753	2.131	2.249	2.602	2.947	3.286	3.733	4.073
16	.690	.865	1.071	1.337	1.746	2.120	2.235	2.583	2.921	3.252	3.686	4.015
17	.689	.863	1.069	1.333	1.740	2.110	2.224	2.567	2.898	3.222	3.646	3.965
18	.688	.862	1.067	1.330	1.734	2.101	2.214	2.552	2.878	3.197	3.611	3.922
19	.688	.861	1.066	1.328	1.729	2.093	2.205	2.539	2.861	3.174	3.579	3.883
20	.687	.860	1.064	1.325	1.725	2.086	2.197	2.528	2.845	3.153	3.552	3.850
21	.686	.859	1.063	1.323	1.721	2.080	2.189	2.518	2.831	3.135	3.527	3.819
22	.686	.858	1.061	1.321	1.717	2.074	2.183	2.508	2.819	3.119	3.505	3.792
23	.685	.858	1.060	1.319	1.714	2.069	2.177	2.500	2.807	3.104	3.485	3.768
24	.685	.857	1.059	1.318	1.711	2.064	2.172	2.492	2.797	3.091	3.467	3.745
25	.684	.856	1.058	1.316	1.708	2.060	2.167	2.485	2.787	3.078	3.450	3.725
26	.684	.856	1.058	1.315	1.706	2.056	2.162	2.479	2.779	3.067	3.435	3.707
27	.684	.855	1.057	1.314	1.703	2.052	2.158	2.473	2.771	3.057	3.421	3.690
28	.683	.855	1.056	1.313	1.701	2.048	2.154	2.467	2.763	3.047	3.408	3.674
29	.683	.854	1.055	1.311	1.699	2.045	2.150	2.462	2.756	3.038	3.396	3.659
30	.683	.854	1.055	1.310	1.697	2.042	2.147	2.457	2.750	3.030	3.385	3.646

# One Sample t Test

---

- ❖ Calculated value
- ❖  $t = [\bar{x} - \mu] / [s / \text{sqrt}(n)]$
- ❖ Example: Perfume bottle producing 150cc, 4 bottles are randomly picked and the average volume was found to be 151cc and sd of sample was 2 cc. Has mean volume changed? (95% confidence)
- ❖  $t_{\text{cal}} = (151-150)/[2 / \text{sqrt}(4)] = 1/1 = 1$
- ❖  $t_{\text{critical}} = 3.182 > \text{Fail to reject } H_0$

# One Sample p Test

---

- ❖  $H_0: p = p_0$
- ❖ Calculated value

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1 - p_0)}{n}}}$$

- ❖ Example: Smoking rate in a town in past was 21%, 100 samples were picked and found 14 smokers. Has smoking habit changed?

# One Sample p Test

❖ Example: Smoking rate in a town in past was 21%, 100 samples were picked and found 14 smokers. Has smoking habit changed at 95% confidence? (two tail)

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1 - p_0)}{n}}}$$

❖  $p_0 = 0.21$ ,  $p = 0.14$

❖  $np_0 = 0.21 \times 100 = 21$  and  $n(1 - p_0) = 0.79 \times 100 = 79$

❖  $> 5$  means sample size is sufficient.

❖  $z = (0.14 - 0.21) / \sqrt{0.21 \times 0.79 / 100}$

❖  $z = -0.07 / 0.0407 = -1.719$

❖  $z_{\text{critical}} = 1.96$



# One Sample p Test

❖ Example: Smoking rate in a town in past was 21%, 100 samples were picked and found 14 smokers. Has smoking habit reduced at 95% confidence? (one tail)

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1 - p_0)}{n}}}$$

❖  $H_0: p < p_0$

❖  $p_0 = 0.21, p = 0.14$

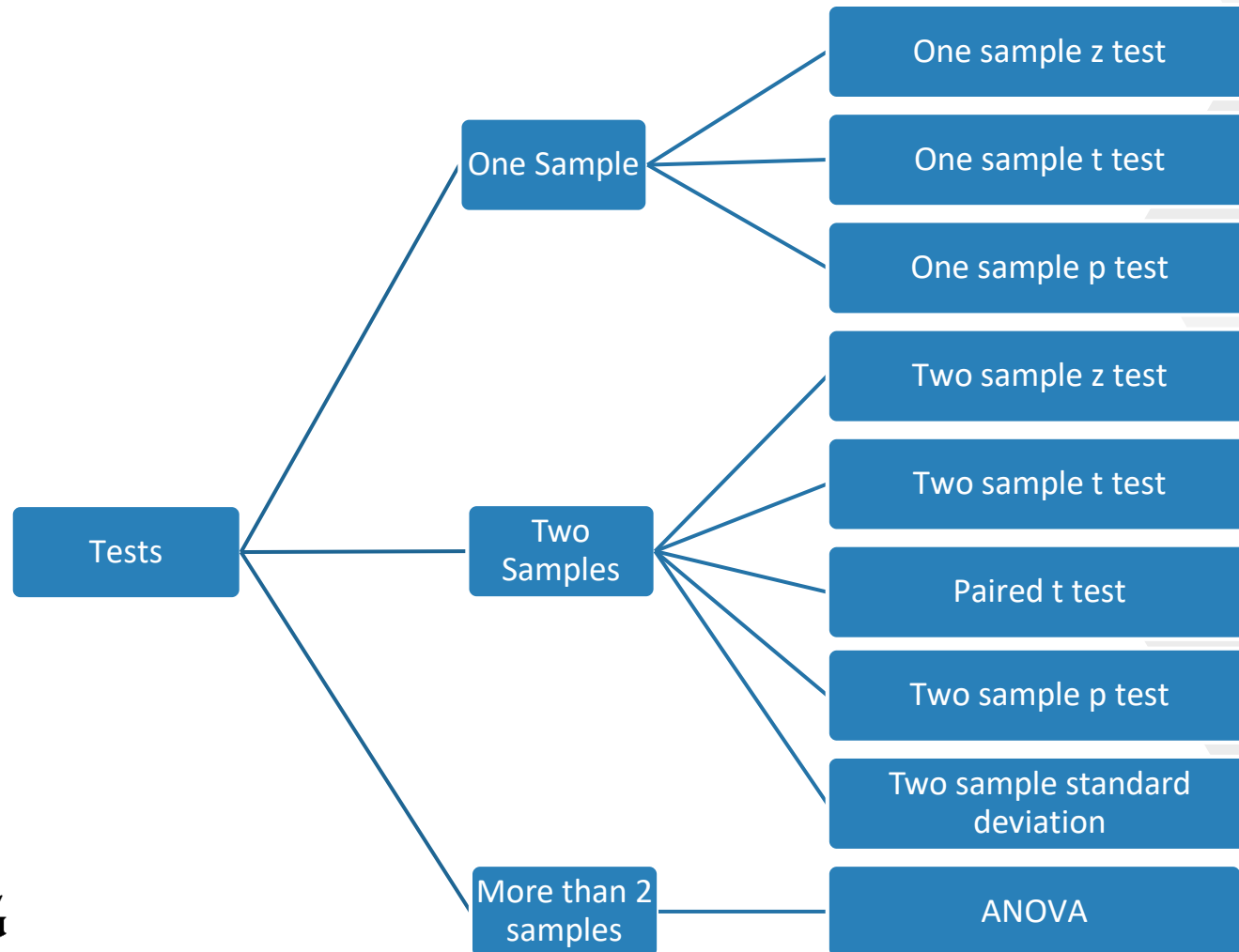
❖  $z = (0.14 - 0.21) / \sqrt{0.21 \times 0.79 / 100}$

❖  $z = -0.07 / 0.0407 = -1.719$

❖  $z_{\text{critical}} = 1.645$



# Tests for Mean, Variance & Proportion



# Two Sample z Test

---

- ❖ **Null hypothesis:**  $H_0: \mu_1 = \mu_2$
- ❖ or  $H_0: \mu_1 - \mu_2 = 0$
- ❖ **Alternative hypothesis:**  $H_a: \mu_1 \neq \mu_2$

$$z = \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

# Two Sample z Test

---

- ❖ Example: From two machines 100 samples each were drawn.
- ❖ Machine 1: Mean = 151.2 / sd = 2.1
- ❖ Machine 2: Mean = 151.9 / sd = 2.2
- ❖ Is there difference in these two machines.  
Check at 95% confidence level.

$$z = \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

# Two Sample z Test

- ❖ Example: From two machines 100 samples each were drawn.
  - ❖ Machine 1: Mean = 151.2 / sd = 2.1
  - ❖ Machine 2: Mean = 151.9 / sd = 2.2
  - ❖ Is there difference in these two machines.  
Check at 95% confidence level.
- ❖  $Z_{cal} = -0.7 / 0.304 = -2.30$
- ❖  $Z_{critical} = 1.96$
- ❖ Reject Null.
- ❖ There is a difference.

$$z = \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

# Two Sample z Test

- ❖ Example: From two machines 100 samples each were drawn.
  - ❖ Machine 1: Mean = 151.9 / sd = 2.1
  - ❖ Machine 2: Mean = 151.2 / sd = 2.2
  - ❖ Is there **difference of more than 0.2 cc** in these two machines. Check at 95% confidence level.
- ❖  $H_0: \mu_1 - \mu_2 = 0.2$

$$z = \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

# Two Sample z Test

- ❖ Example: From two machines 100 samples each were drawn.
  - ❖ Machine 1: Mean = 151.9 / sd = 2.1
  - ❖ Machine 2: Mean = 151.2 / sd = 2.2
  - ❖ Is there **difference of more than 0.2 cc** in these two machines. Check at 95% confidence level.
- ❖  $Z_{\text{cal}} = 0.5/0.304 = 2.30$
- ❖  $Z_{\text{critical}} = 1.64$
- ❖ Reject Null Hypothesis.

$$z = \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

# Two Sample t Test

---

- ❖ If two set of data are independent or dependent.
  - ❖ If the values in one sample reveal no information about those of the other sample, then the samples are independent.
    - ❖ Example: Blood pressure of male/female
- ❖ If the values in one sample affect the values in the other sample, then the samples are dependent.
  - ❖ Example: Blood pressure before and after a specific medicine

# Two Sample t Test

---

- ❖ If two set of data are independent or dependent.
  - ❖ If the values in one sample reveal no information about those of the other sample, then the samples are independent.
    - ❖ Example: Blood pressure of male/female
- ❖ If the values in one sample affect the values in the other sample, then the samples are dependent.
  - ❖ Example: Blood pressure before and after a specific medicine

*Two sample t test*

*Paired t test*



# Two Sample t Test

---

- ❖ Is variance for two samples equal?
- ❖ If yes: Pooled variance calculate  $S_p$  for finding out t

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

$$t = \frac{(\bar{x}_1 - \bar{x}_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

# Paired t Test

---

- ❖ Where you have two samples in which observations in one sample can be paired with observations in the other sample.
- ❖ Or
- ❖ If the values in one sample affect the values in the other sample, (the samples are dependent.)
  - ❖ Example: Blood pressure before and after a specific medicine

# Paired t Test

---

- ❖ Find the difference between two set of readings as  $d_1, d_2 \dots d_n$ .
- ❖ Find the mean and standard deviation of these differences.

$$t = \frac{\bar{d}}{\sqrt{s^2/n}}$$

# Paired t Test

---

- ❖ Example: Before and after medicine BP was measured. Is there a difference at 95% confidence level?

$$t = \frac{\bar{d}}{\sqrt{s^2/n}}$$

Patient	Before	After
1	120	122
2	122	120
3	143	141
4	100	109
5	109	109

# Paired t Test

- ❖ Example: Before and after medicine BP was measured. Is there a difference at 95% confidence level?

$$t = \frac{\bar{d}}{\sqrt{s^2/n}}$$

Patient	Before	After	difference
1	120	122	2
2	122	120	-2
3	143	141	-2
4	100	109	9
5	109	109	0

- ❖  $\bar{d} = 1.4$  ,  $s = 4.56$  ,  $n=5$

- ❖  $t_{\text{cal.}} = 1.4/1.99 = 0.70$

# Paired t Test

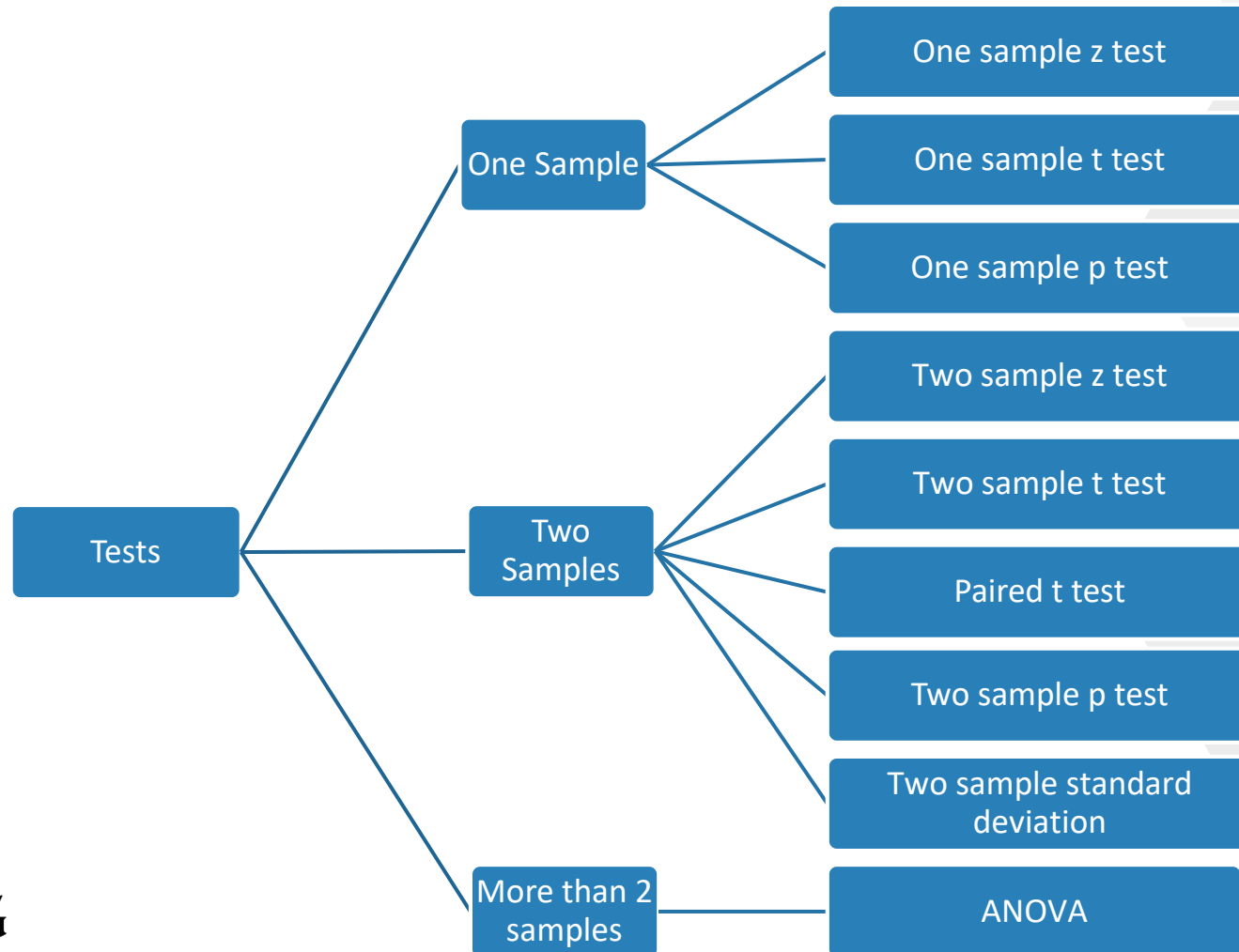
- ❖ Example: Before and after medicine BP was measured. Is there a difference at 95% confidence level?

$$t = \frac{\bar{d}}{\sqrt{s^2/n}}$$

Patient	Before	After	difference
1	120	122	2
2	122	120	-2
3	143	141	-2
4	100	109	9
5	109	109	0

- ❖  $t_{\text{cal.}} = 1.4/1.99 = 0.70$
- ❖  $t_{0.025, 4} = 2.766$
- ❖ Fail to reject null hypothesis

# Tests for Mean, Variance & Proportion



# Two Sample p Test

---

$$Z = \frac{\hat{p} - p_o}{\sqrt{\frac{p_o(1 - p_o)}{n}}}$$

- ❖ **Null hypothesis:**  $H_0: p_1 = p_2$
- ❖ or  $H_0: p_1 - p_2 = 0$
- ❖ **Alternative hypothesis:**  $H_a: p_1 \neq p_2$
- ❖ **Normal approximation – Pooled**
- ❖ **Normal approximation – Un-pooled**

$$\frac{(\bar{p}_1 - \bar{p}_2)}{\sqrt{\bar{p}(1 - \bar{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

$$Z = \frac{(p_1 - p_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}}$$



# Tests for Variance

---

## ❖ F-test

- ❖ for testing equality of *two* variances from different populations
- ❖ for testing equality of several means with technique of ANOVA.

## ❖ Chi-square test

- ❖ For testing the population variance against a specified value
- ❖ testing goodness of fit of some probability distribution

# Two Sample Variance – F Test

## ❖ F-test

- ❖ for testing equality of *two* variances from different populations

- ❖  $H_0: \sigma^2_1 = \sigma^2_2$

- ❖  $F_{\text{calculated}}$  
$$F_{\text{calculated}} = \frac{S_1^2}{S_2^2}$$

- ❖ Keep higher value at the top for right tail test.
- ❖ Remember: Variance is square of standard deviation

# Two Sample Variance – F Test

❖  $F_{critical}$

❖ Use table with appropriate degrees of freedom

❖ For two tail test use the table for  $\alpha/2$

F - Distribution ( $\alpha = 0.05$  in the Right Tail)

df <sub>1</sub>	1	2	3	4
1	161.45	199.50	215.71	224.58
2	18.513	19.000	19.164	19.247
3	10.128	9.5521	9.2766	9.1172
4	7.7086	6.9443	6.5914	6.3882
5	6.6079	5.7861	5.4095	5.1622
6	5.9874	5.1433	4.7571	4.5337
7	5.5914	4.7374	4.3468	4.1203
8	5.3177	4.4599	4.0662	3.8379
9	5.1174	4.2565	3.8625	3.6331
10	4.9646	4.1028	3.7083	3.4780
11	4.8443	3.9823	3.5874	3.3567
12	4.7472	3.8853	3.4903	3.2592
13	4.6672	3.8056	3.4105	3.1791
14	4.6001	3.7389	3.3439	3.1122
15	4.5431	3.6823	3.2874	3.0556
16	4.4940	3.6337	3.2389	3.0069
17	4.4513	3.5915	3.1968	2.9647
18	4.4139	3.5546	3.1599	2.9277
19	4.3807	3.5219	3.1274	2.8951
20	4.3512	3.4928	3.0984	2.8661
21	4.3248	3.4668	3.0725	2.8401
22	4.3009	3.4434	3.0491	2.8167
23	4.2793	3.4221	3.0280	2.7955
24	4.2597	3.4028	3.0085	2.7763
25	4.2417	3.3852	2.9912	2.7587
26	4.2252	3.3690	2.9752	2.7426
27	4.2100	3.3541	2.9604	2.7278
28	4.1960	3.3404	2.9467	2.7141
29	4.1830	3.3277	2.9340	2.7014
30	4.1709	3.3158	2.9223	2.6896
40	4.0847	3.2317	2.8387	2.6060
60	4.0012	3.1504	2.7581	2.5252
120	3.9201	3.0718	2.6802	2.4472
=	3.8415	2.9957	2.6049	2.3719

F - Distribution ( $\alpha = 0.01$  in the Right Tail)

df <sub>1</sub>	1	2	3	4	5	6	7	8	9
1	4052.2	4999.5	5403.4	5624.6	5763.4	5859.0	5928.4	5981.1	6022.5
2	98.503	99.000	99.166	99.249	99.299	99.333	99.356	99.374	99.388
3	34.116	30.817	29.457	28.710	28.237	27.911	27.672	27.489	27.345
4	21.198	18.000	16.694	15.977	15.522	15.207	14.976	14.799	14.659
5	16.258	13.274	12.008	11.392	10.967	10.672	10.456	10.289	10.158
6	13.745	10.925	9.7795	9.1483	8.7499	8.4661	8.2600	8.1017	7.9761
7	12.246	9.5466	8.4513	7.8466	7.4604	7.1914	6.9928	6.8400	6.7188
8	11.259	8.6491	7.5910	7.0061	6.6318	6.3707	6.1776	6.0289	5.9106
9	10.561	8.0125	6.9919	6.4221	6.0569	5.8018	5.6129	5.4671	5.3511
10	10.044	7.5594	6.5523	5.9943	5.6363	5.3858	5.2001	5.0567	4.9424
11	9.6460	7.2077	6.2167	5.6683	5.3160	5.0692	4.8861	4.7445	4.6315
12	9.3302	6.9266	5.9525	5.4120	5.0643	4.8206	4.6395	4.4994	4.3875
13	9.0718	6.7010	5.7394	5.2053	4.8616	4.6204	4.4410	4.3021	4.1911
14	8.8616	6.5149	5.5639	5.0354	4.6950	4.4558	4.2779	4.1399	4.0297
15	8.6831	6.3589	5.4170	4.8932	4.5556	4.3183	4.1415	4.0045	3.8948
16	8.5310	6.2262	5.2922	4.7726	4.4374	4.2016	4.0259	3.8896	3.7804
17	8.3997	6.1121	5.1850	4.6690	4.3359	4.1015	3.9267	3.7910	3.6822
18	8.2854	6.0129	5.0919	4.5790	4.2479	4.0146	3.8406	3.7054	3.5971
19	8.1849	5.9259	5.0103	4.5003	4.1708	3.9386	3.7653	3.6305	3.5225
20	8.0960	5.8489	4.9382	4.4207	4.1027	3.8714	3.6987	3.5644	3.4567
21	8.0166	5.7804	4.8740	4.3608	4.0421	3.8117	3.6396	3.5056	3.3981
22	7.9434	5.7190	4.8166	4.3134	3.9980	3.7683	3.5967	3.4630	3.3558
23	7.8811	5.6637	4.7649	4.2636	3.9502	3.7212	3.5490	3.4157	3.3086
24	7.8229	5.6136	4.7181	4.2184	3.9051	3.6767	3.4999	3.3679	3.2610
25	7.7698	5.5680	4.6755	4.1774	3.8650	3.6372	3.4608	3.3293	3.2227
26	7.7213	5.5263	4.6366	4.1400	3.8313	3.6041	3.4281	3.2970	3.1910
27	7.6767	5.4881	4.6009	4.1056	3.7948	3.5680	3.3925	3.2619	3.1564
28	7.6356	5.4529	4.5681	4.0740	3.7599	3.5336	3.3585	3.2285	3.1235
29	7.5977	5.4204	4.5378	4.0449	3.7254	3.4995	3.3249	3.1954	3.0909
30	7.5625	5.3903	4.5097	4.0179	3.6990	3.4735	3.3005	3.1716	3.0675
40	7.3141	5.1705	4.3126	3.8263	3.5138	3.2910	3.1238	2.9950	2.8876
60	7.0771	4.9774	4.1259	3.6400	3.3389	3.1187	2.9530	2.8233	2.7165
120	6.8909	4.7865	3.9491	3.4795	3.1735	2.9559	2.7916	2.6629	2.5566
=	6.6349	4.6052	3.7816	3.3192	3.0173	2.8020	2.6393	2.5113	2.4073

# Two Sample Variance – F Test

❖ Example: We took 8 samples from machine A and the standard deviation was 1.1. For machine B we took 5 samples and the variance was 11. Is there a difference in variance at 90% confidence level?

❖  $n_1 = 8, s_1 = 1.1, s_1^2 = 1.21, df = 7$  (denominator)

❖  $n_2 = 5, s_2^2 = 11, df = 4$  (numerator)

❖  $F_{\text{calculated}} = 11/1.21 = 9.09$  (higher value at top)

F - Distribution ( $\alpha = 0.05$  in the Right Tail)

df <sub>1</sub>	1	2	3	4
1	161.45	199.50	215.71	224.58
2	18.513	19.000	19.164	19.247
3	10.128	9.5521	9.2796	9.1172
4	7.7086	6.9443	6.5914	6.3882
5	6.6079	5.7861	5.4095	5.1622
6	5.9874	5.1433	4.7571	4.5337
7	5.5914	4.7374	4.3488	4.1203
8	5.3177	4.4590	4.0662	3.8379
9	5.1174	4.2565	3.8625	3.6331
10	4.9646	4.1028	3.7083	3.4780
11	4.8443	3.9823	3.5874	3.3567
12	4.7472	3.8853	3.4903	3.2592
13	4.6672	3.8056	3.4105	3.1791
14	4.6001	3.7389	3.3439	3.1122
15	4.5431	3.6823	3.2874	3.0556
16	4.4940	3.6337	3.2389	3.0069
17	4.4513	3.5915	3.1968	2.9647
18	4.4139	3.5546	3.1599	2.9277
19	4.3807	3.5219	3.1274	2.8951
20	4.3512	3.4928	3.0984	2.8661
21	4.3248	3.4668	3.0725	2.8401
22	4.3009	3.4434	3.0491	2.8167
23	4.2793	3.4221	3.0280	2.7955
24	4.2597	3.4028	3.0085	2.7763
25	4.2417	3.3852	2.9912	2.7587
26	4.2252	3.3690	2.9752	2.7426
27	4.2100	3.3541	2.9604	2.7278
28	4.1960	3.3404	2.9467	2.7141
29	4.1830	3.3277	2.9340	2.7014
30	4.1709	3.3158	2.9223	2.6896
40	4.0847	3.2317	2.8387	2.6060
60	4.0012	3.1504	2.7581	2.5252
120	3.9201	3.0718	2.6802	2.4472
=	3.8415	2.9957	2.6049	2.3719

F - Distribution ( $\alpha = 0.01$  in the Right Tail)

df <sub>1</sub>	1	2	3	4	5	6	7	8	9
1	4052.2	4999.5	5403.4	5624.6	5763.6	5859.0	5928.4	5981.1	6022.5
2	98.503	99.000	99.166	99.249	99.299	99.333	99.356	99.374	99.388
3	38.116	38.817	39.457	39.870	40.170	40.371	40.502	40.589	40.645
4	21.199	18.000	16.694	15.977	15.522	15.207	14.976	14.799	14.659
5	16.258	13.274	12.000	11.392	10.967	10.672	10.456	10.289	10.158
6	13.745	10.925	9.7795	9.1483	8.7499	8.4661	8.2600	8.1017	7.9761
7	12.246	9.5466	8.4513	7.8466	7.4604	7.1914	6.9928	6.8400	6.7188
8	11.259	8.6491	7.5910	7.0061	6.6318	6.3707	6.1776	6.0289	5.9106
9	10.561	8.0215	6.9919	6.4221	6.0569	5.8018	5.6129	5.4671	5.3511
10	10.044	7.5594	6.5523	5.9943	5.6363	5.3858	5.2001	5.0567	4.9424
11	9.6460	7.2077	6.2167	5.6683	5.3160	5.0692	4.8861	4.7445	4.6315
12	9.3302	6.9266	5.9525	5.4120	5.0643	4.8206	4.6395	4.4994	4.3875
13	9.0718	6.7010	5.7394	5.2053	4.8616	4.6204	4.4410	4.3021	4.1911
14	8.8616	6.5149	5.5639	5.0354	4.6950	4.4558	4.2779	4.1399	4.0297
15	8.6831	6.3589	5.4170	4.8932	4.5556	4.3183	4.1415	4.0045	3.8948
16	8.5310	6.2262	5.2922	4.7726	4.4374	4.2016	4.0259	3.8896	3.7804
17	8.3997	6.1121	5.1850	4.6690	4.3359	4.1015	3.9267	3.7910	3.6822
18	8.2854	6.0129	5.0919	4.5790	4.2479	4.0146	3.8406	3.7054	3.5971
19	8.1849	5.9259	5.0103	4.5003	4.1708	3.9386	3.7653	3.6305	3.5225
20	8.0960	5.8489	4.9382	4.4207	4.1027	3.8714	3.6987	3.5644	3.4567
21	8.0166	5.7804	4.8740	4.3608	4.0461	3.8177	3.6456	3.5116	3.4041
22	7.9434	5.7190	4.8166	4.3134	3.9980	3.7783	3.5867	3.4530	3.3458
23	7.8811	5.6637	4.7649	4.2656	3.9502	3.7302	3.5390	3.4057	3.2986
24	7.8229	5.6136	4.7181	4.2184	3.8951	3.6767	3.4959	3.3629	3.2560
25	7.7698	5.5680	4.6755	4.1774	3.8550	3.6372	3.4568	3.3239	3.2172
26	7.7213	5.5263	4.6366	4.1400	3.8183	3.5911	3.4110	3.2784	3.1721
27	7.6767	5.4881	4.6009	4.1056	3.7848	3.5580	3.3822	3.2500	3.1444
28	7.6356	5.4529	4.5681	4.0740	3.7539	3.5276	3.3511	3.2199	3.1145
29	7.5977	5.4204	4.5378	4.0449	3.7254	3.4995	3.3303	3.1982	3.0920
30	7.5625	5.3903	4.5097	4.0179	3.6990	3.4735	3.3045	3.1726	3.0665
40	7.3141	5.1705	4.3126	3.8253	3.5138	3.2910	3.1238	2.9930	2.8876
60	7.0771	4.9774	4.1259	3.6400	3.3389	3.1187	2.9530	2.8233	2.7185
120	6.8909	4.7865	3.9491	3.4795	3.1735	2.9559	2.7916	2.6629	2.5586
=	6.6349	4.6052	3.7816	3.3192	3.0173	2.8020	2.6393	2.5113	2.4073

# Two Sample Variance – F Test

F - Distribution ( $\alpha = 0.01$  in the Right Tail)

df <sub>2</sub> \ df <sub>1</sub>		Numerator Degrees of Freedom								
		1	2	3	4	5	6	7	8	9
1	4052.2									
2	98.50									
3	34.1									
4	21.1									
5	16.2									
6	13.7									
7	12.2									
8	11.2									
9	10.5									
10	10.0									
11	9.6									
12	9.3									
13	9.0									
14	8.8									
15	8.6									
16	8.5									
17	8.3									
18	8.2									
19	8.1									
20	8.0									
21	8.0									
22	7.9									
23	7.8									
24	7.8									
25	7.7									
26	7.7									
27	7.6									
28	7.6									
29	7.5									
30	7.5									
40	7.3									
60	7.0									
120	6.8									
∞	6.6									

F - Distribution ( $\alpha = 0.05$  in the Right Tail)

df <sub>2</sub> \ df <sub>1</sub>		Numerator Degrees of Freedom								
		1	2	3	4	5	6	7	8	9
1	161.45									
2	18.513									
3	10.128									
4	7.7086									
5	6.6079									
6	5.9874									
7	5.5914									
8	5.3177									
9	5.1174									
10	4.9646									
11	4.8443									
12	4.7472									
13	4.6672									
14	4.6001									
15	4.5431									
16	4.4940									
17	4.4513									
18	4.4139									
19	4.3807									
20	4.3512									
21	4.3248									
22	4.3009									
23	4.2793									
24	4.2597									
25	4.2417									
26	4.2252									
27	4.2100									
28	4.1960									
29	4.1830									
30	4.1709									
40	4.0847									
60	4.0012									
120	3.9201									
∞	3.8415									

$$F_{\text{calculated}} = \frac{S_1^2}{S_2^2}$$

$$F_{\text{critical}} = 4.1203$$

# Two Sample Variance – F Test

❖ Example: We took 8 samples from machine A and the standard deviation was 1.1. For machine B we took 5 samples and the variance was 11. Is there a difference in variance at 90% confidence level?

❖  $n_1 = 8, s_1 = 1.1, s_1^2 = 1.21, df = 7$  (denominator)

❖  $n_2 = 5, s_2^2 = 11, df = 4$  (numerator)

❖  $F_{\text{calculated}} = 11/1.21 = 9.09$  (higher value at top)

❖ Reject  $H_0$

F - Distribution ( $\alpha = 0.05$  in the Right Tail)

df <sub>1</sub>	Numerator Degrees of Freedom			
	1	2	3	4
1	161.45	199.50	215.71	224.58
2	18.513	19.000	19.164	19.247
3	10.128	9.5521	9.2766	9.1172
4	7.7086	6.9443	6.5914	6.3882
5	6.5910	5.7861	5.4095	5.1622
6	5.9648	5.1433	4.7571	4.5337
7	5.5014	4.7374	4.3468	4.1203
8	5.1174	4.4500	4.0622	3.8379
9	4.7864	4.1628	3.7683	3.5450
10	4.4943	3.8853	3.4903	3.2592
11	4.2545	3.6374	3.2424	3.0113
12	4.0443	3.4133	3.0173	2.7863
13	3.8543	3.2033	2.8073	2.5763
14	3.6843	2.9933	2.5973	2.3663
15	3.5343	2.7933	2.3973	2.1663
16	3.3943	2.6033	2.2073	1.9763
17	3.2643	2.4233	2.0273	1.7863
18	3.1443	2.2533	1.8473	1.6063
19	3.0343	2.0933	1.6773	1.4363
20	2.9343	1.9433	1.5173	1.2763
25	2.6543	1.6633	1.2373	1.0063
30	2.4543	1.4633	1.0373	0.8063
40	2.2543	1.2633	0.8373	0.6063
60	2.0543	1.0633	0.6373	0.4063
120	1.8543	0.8633	0.4373	0.2063
=	3.8415	2.9957	2.6049	2.3719

F - Distribution ( $\alpha = 0.01$  in the Right Tail)

df <sub>1</sub>	Numerator Degrees of Freedom			
	1	2	3	4
1	4052.2	4999.5	5403.4	5624.6
2	98.503	99.000	99.164	99.247
3	34.116	30.817	29.457	28.710
4	21.199	18.000	16.694	15.972
5	16.258	13.274	12.000	11.392
6	13.745	10.925	9.7795	9.1483
7	12.246	9.5466	8.4513	7.8466
8	11.259	8.4491	7.5910	7.0041
9	10.561	7.6215	6.9919	6.4221
10	10.044	7.0594	6.5223	5.9943
11	9.6460	6.6207	6.1047	5.6083
12	9.3302	6.2866	5.7925	5.3160
13	9.0718	5.9916	5.5053	5.0204
14	8.8616	5.7419	5.2559	4.7708
15	8.6811	5.5219	5.0359	4.5558
16	8.5310	5.3262	4.8322	4.3726
17	8.3997	5.1421	4.6440	4.2005
18	8.2854	4.9679	4.4709	4.0386
19	8.1849	4.8029	4.3113	3.8861
20	8.0960	4.6469	4.1648	3.7424
21	8.0166	4.4984	4.0291	3.6070
22	7.9434	4.3569	3.8939	3.4794
23	7.8751	4.2219	3.7588	3.3582
24	7.8129	4.0929	3.6239	3.2429
25	7.7558	3.9689	3.4891	3.1329
26	7.7031	3.8499	3.3543	3.0279
27	7.6541	3.7359	3.2195	2.9279
28	7.6086	3.6269	3.0847	2.8329
29	7.5677	3.5219	2.9500	2.7429
30	7.5305	3.4209	2.8152	2.6579
40	7.3141	3.1705	2.5126	2.3519
60	7.0771	2.9774	2.2599	2.1000
120	6.8009	2.7865	2.0491	1.8795
=	6.6349	2.6052	1.8716	1.7023

$$F_{\text{critical}} = 4.1203$$

# Tests for Variance

---

## ❖ F-test

- ❖ for testing equality of *two* variances from different populations
- ❖ for testing equality of several means with technique of ANOVA.

## ❖ Chi-square test

- ❖ For testing the population variance against a specified value
- ❖ testing goodness of fit of some probability distribution

# One Sample Chi Square

---

- ❖ For testing the population variance against a specified value  $\sigma$

$$\chi^2 = \frac{(n - 1)S^2}{\sigma^2}$$



# One Sample Chi Square

---

- ❖ Example: A sample of 25 bottles was selected. The variance of these 25 bottles as 5 cc. Has it **increased** from established 4 cc? 95% confidence level.
- ❖  $\chi^2 = 24 \times 5 / 4 = 30$
- ❖ What is critical value of Chi Square for 24 degrees of freedom?

$$\chi^2 = \frac{(n - 1)S^2}{\sigma^2}$$

# One Sample Chi Square

---

- ❖ Example: A sample of 25 bottles was selected. The variance of these 25 bottles as 5 cc. Has it **increased** from established 4 cc? 95% confidence level.
- ❖  $\chi^2 = 24 \times 5 / 4 = 30$
- ❖ What is critical value of Chi Square for 24 degrees of freedom?

$$\chi^2 = \frac{(n - 1)S^2}{\sigma^2}$$

# One Sample Chi Square

Percentage Points of the Chi-Square Distribution

Degrees of Freedom	Probability of a larger value of $\chi^2$								
	0.99	0.95	0.90	0.75	0.50	0.25	0.10	0.05	0.01
1	0.000	0.004	0.016	0.102	0.455	1.32	2.71	3.84	6.63
2	0.020	0.103	0.211	0.575	1.386	2.77	4.61	5.99	9.21
3	0.115	0.352	0.584	1.212	2.366	4.11	6.25	7.81	11.34
4	0.297	0.711	1.064	1.923	3.357	5.39	7.78	9.49	13.28
5	0.554	1.145	1.610	2.675	4.351	6.63	9.24	11.07	15.09
6	0.872	1.635	2.204	3.455	5.348	7.84	10.64	12.59	16.81
7	1.239	2.167	2.833	4.255	6.346	9.04	12.02	14.07	18.48
8	1.647	2.733	3.490	5.071	7.344	10.22	13.36	15.51	20.09
9	2.088	3.325	4.168	5.899	8.343	11.39	14.68	16.92	21.67
10	2.558	3.940	4.865	6.737	9.342	12.55	15.99	18.31	23.21
11	3.053	4.575	5.578	7.584	10.341	13.70	17.28	19.68	24.72
12	3.571	5.226	6.304	8.438	11.340	14.85	18.55	21.03	26.22
13	4.107	5.892	7.042	9.299	12.340	15.98	19.81	22.36	27.69
14	4.660	6.571	7.790	10.165	13.339	17.12	21.06	23.68	29.14
15	5.229	7.261	8.547	11.037	14.339	18.25	22.31	25.00	30.58
16	5.812	7.962	9.312	11.912	15.338	19.37	23.54	26.30	32.00
17	6.408	8.672	10.085	12.792	16.338	20.49	24.77	27.59	33.41
18	7.015	9.390	10.865	13.675	17.338	21.60	25.99	28.87	34.80
19	7.633	10.117	11.651	14.562	18.338	22.72	27.20	30.14	36.19
20	8.260	10.851	12.443	15.452	19.337	23.83	28.41	31.41	37.57
22	9.542	12.338	14.041	17.240	21.337	26.04	30.81	33.92	40.29
24	10.856	13.848	15.659	19.037	23.337	28.24	33.20	36.42	42.98
26	12.198	15.379	17.292	20.843	25.336	30.43	35.56	38.89	45.64
28	13.565	16.928	18.939	22.657	27.336	32.62	37.92	41.34	48.28
30	14.953	18.493	20.599	24.478	29.336	34.80	40.26	43.77	50.89
40	22.164	26.509	29.051	33.660	39.335	45.62	51.80	55.76	63.69
50	27.707	34.764	37.689	42.942	49.335	56.33	63.17	67.50	76.15
60	37.485	43.188	46.459	52.294	59.335	66.98	74.40	79.08	88.38

$$\chi^2 = \frac{(n - 1)S^2}{\sigma^2}$$

# One Sample Chi Square

- ❖ Example: A sample of 25 bottles was selected. The variance of these 25 bottles as 5 cc. Has it **increased** from established 4 cc? 95% confidence level.
- ❖  $\chi^2 = 24 \times 5 / 4 = 30$
- ❖ Critical value of Chi Square for 24 degrees of freedom = 36.42
- ❖ Fail to reject  $H_0$

$$\chi^2 = \frac{(n - 1)S^2}{\sigma^2}$$

# ANOVA

---

## ❖ F-test

- ❖ for testing equality of *two* variances from different populations
- ❖ for testing equality of several means with technique of ANOVA.

## ❖ Chi-square test

- ❖ For testing the population variance against a specified value
- ❖ testing goodness of fit of some probability distribution
- ❖ testing for independence of two attributes

# ANOVA

---

## ❖ F-test

- ❖ for testing equality of *two* variances from different populations

- ❖  $H_0: \sigma^2_1 = \sigma^2_2$

- ❖  $F_{\text{calculated}}$

$$F_{\text{calculated}} = \frac{S_1^2}{S_2^2}$$

- ❖ Keep higher value at the top for right tail test.
- ❖ Remember: Variance is square of standard deviation

# ANOVA

---

## ❖ Why ANOVA?

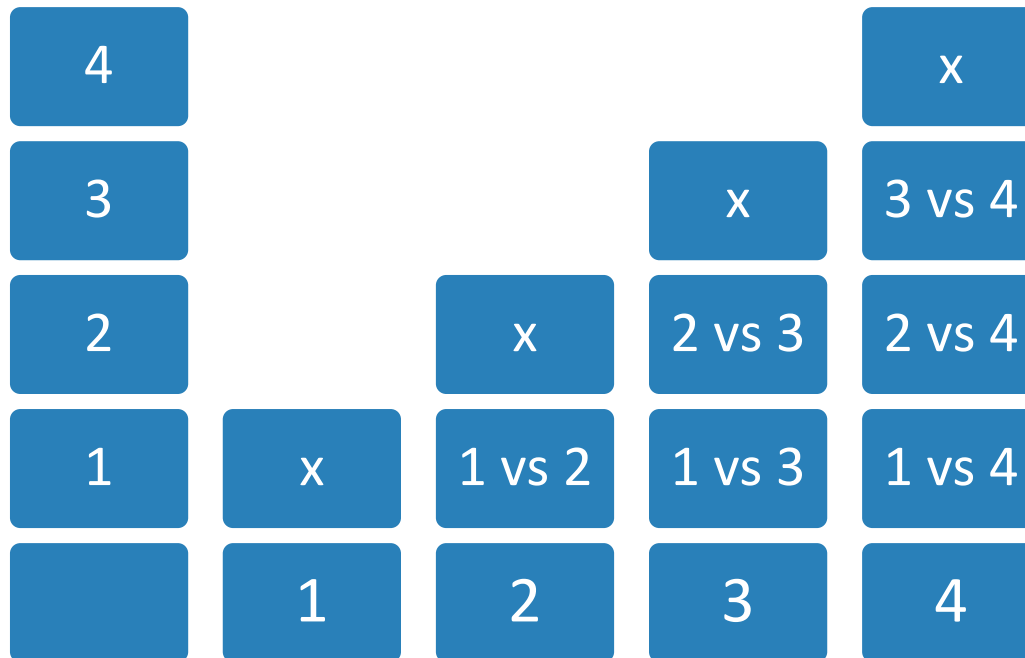
- ❖ We used t test to compare the means of two populations.
- ❖ What if we need to compare more than two populations? With ANOVA we can find out if one or more populations have different mean or comes from a different population.
- ❖ We could have conducted multiple t Test.
- ❖ How many t Test we need to conduct if have to compare 4 samples? ... 6

# ANOVA

---

## ❖ Why ANOVA?

- ❖ How many t Test we need to conduct if have to compare 4 samples? ... 6





# ANOVA

---

## ❖ Why ANOVA?

- ❖ How many t Test we need to conduct if have to compare 4 samples? ... 6
- ❖ Each test is done with  $\alpha = 0.05$  or 95% confidence.
- ❖ 6 tests will result in confidence level of  $0.95 \times 0.95 \times 0.95 \times 0.95 \times 0.95 \times 0.95 = 0.735$

# ANOVA

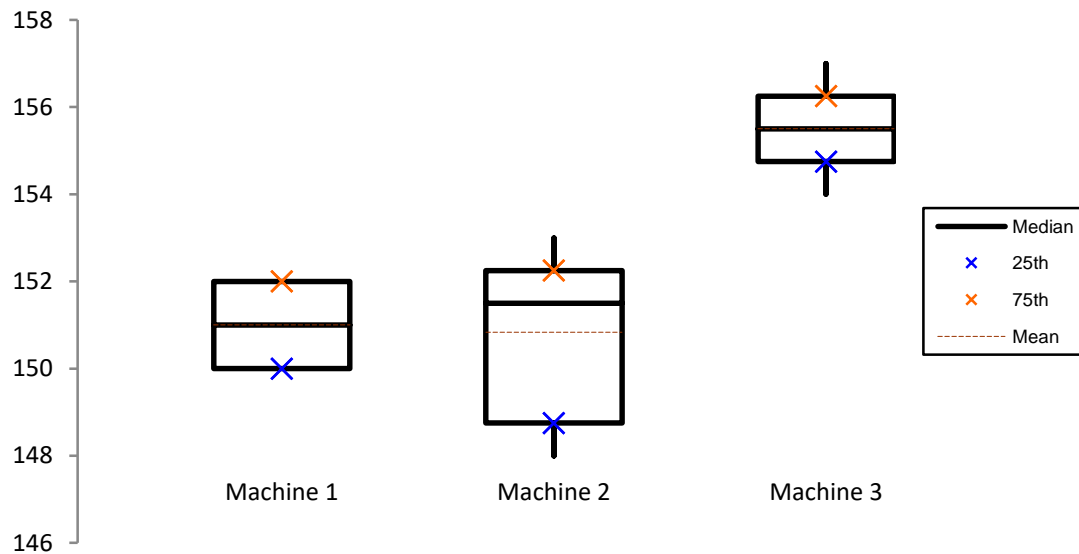
---

❖ Comparing three machines:

Machine 1	Machine 2	Machine 3
150	153	156
151	152	154
152	148	155
152	151	156
151	149	157
150	152	155
$\bar{x}_1 = 151$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$

# ANOVA

## ❖ Comparing three machines:



Machine 1	Machine 2	Machine 3
150	153	156
151	152	154
152	148	155
152	151	156
151	149	157
150	152	155
$\bar{x}_1 = 151.00$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$

# ANOVA

---

## ❖ Comparing three machines:

Machine 1	Machine 2	Machine 3
150	153	156
151	152	154
152	148	155
152	151	156
151	149	157
150	152	155
$\bar{x}_1 = 151.00$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$

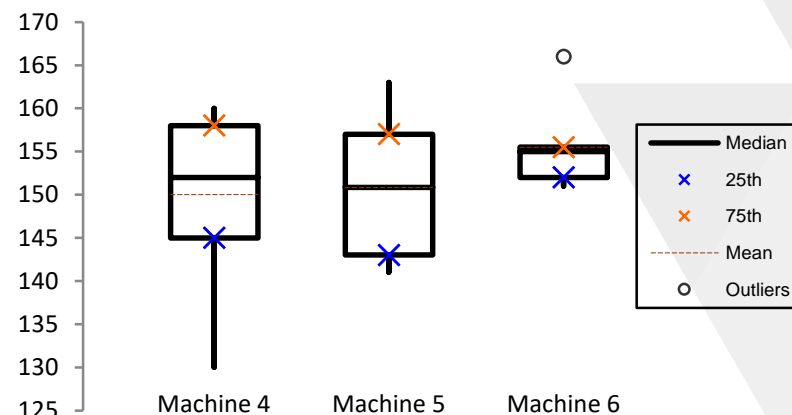
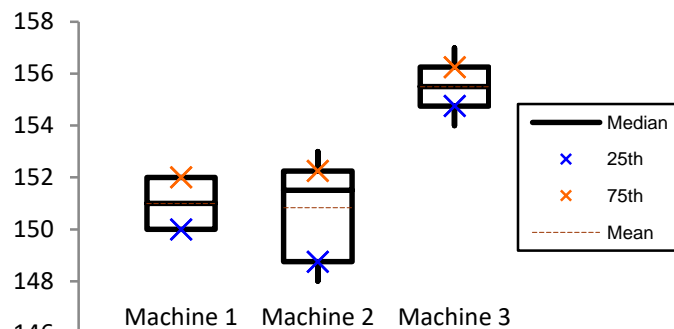
Machine 4	Machine 5	Machine 6
130	163	166
155	152	154
160	143	155
158	141	151
152	149	152
145	157	155
$\bar{x}_1 = 151.00$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$

# ANOVA

## ❖ Comparing three machines:

Machine 1	Machine 2	Machine 3
150	153	156
151	152	154
152	148	155
152	151	156
151	149	157
150	152	155
$\bar{x}_1 = 151.00$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$

Machine 4	Machine 5	Machine 6
130	163	166
155	152	154
160	143	155
158	141	151
152	149	152
145	157	155
$\bar{x}_1 = 151.00$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$



# ANOVA

---

❖ ANOVA is Analysis of Variance

❖ Variance

$$s^2 = \frac{\sum (x_i - \bar{X})^2}{n-1}$$

❖ Numerator of this formula is Sum of Squares

❖ Total of Sum of Squares (SST) =  
 $SS_{\text{between/or treatment}} + SS_{\text{within/or error}}$

# ANOVA

---

❖  $SST = SS_{\text{between(or treatment)}} + SS_{\text{within(or error)}}$

❖ Ratio:

$$SS_{\text{between(or treatment)}} / SS_{\text{within(or error)}}$$

$$F = MS_{\text{between(or treatment)}} / MS_{\text{within(or error)}}$$

# ANOVA

---

$$\diamond SST = SS_{\text{between(or treatment)}} + SS_{\text{within(or error)}}$$

Machine 1	Machine 2	Machine 3
150	153	156
151	152	154
152	148	155
152	151	156
151	149	157
150	152	155
$\bar{x}_1 = 151.00$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$



# ANOVA



Machine 1	Machine 2	Machine 3
150	153	156
151	152	154
152	148	155
152	151	156
151	149	157
150	152	155
$\bar{x}_1 = 151.00$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$

$$\diamond SST = SS_{\text{between(or treatment)}} + SS_{\text{within(or error)}}$$

Machine 1	$x_1 - \bar{x}_1$	$Sqr(x_1 - \bar{x}_1)$	Machine 2	$x_2 - \bar{x}_2$	$Sqr(x_2 - \bar{x}_2)$	Machine 3	$x_3 - \bar{x}_3$	$Sqr(x_3 - \bar{x}_3)$	
150.00	-1.00	1.00	153.00	2.17	4.69	156.00	0.50	0.25	
151.00	0.00	0.00	152.00	1.17	1.36	154.00	-1.50	2.25	
152.00	1.00	1.00	148.00	-2.83	8.03	155.00	-0.50	0.25	
152.00	1.00	1.00	151.00	0.17	0.03	156.00	0.50	0.25	
151.00	0.00	0.00	149.00	-1.83	3.36	157.00	1.50	2.25	
150.00	-1.00	1.00	152.00	1.17	1.36	155.00	-0.50	0.25	
151.00			150.83			155.50			152.44
		4.00			18.83			5.50	

$$\diamond SS_{\text{within}} = 4.00 + 18.83 + 5.50 = 28.33$$



# ANOVA

❖  $SST = SS_{\text{between(or treatment)}} + SS_{\text{within(or error)}}$

Machine 1	Machine 2	Machine 3
150	153	156
151	152	154
152	148	155
152	151	156
151	149	157
150	152	155
$\bar{x}_1 = 151.00$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$

Machine 1	$x_1 - \bar{x}_1$	$Sqr(x_1 - \bar{x}_1)$	Machine 2	$x_2 - \bar{x}_2$	$Sqr(x_2 - \bar{x}_2)$	Machine 3	$x_3 - \bar{x}_3$	$Sqr(x_3 - \bar{x}_3)$	
150.00	-1.00	1.00	153.00	2.17	4.69	156.00	0.50	0.25	
151.00	0.00	0.00	152.00	1.17	1.36	154.00	-1.50	2.25	
152.00	1.00	1.00	148.00	-2.83	8.03	155.00	-0.50	0.25	
152.00	1.00	1.00	151.00	0.17	0.03	156.00	0.50	0.25	
151.00	0.00	0.00	149.00	-1.83	3.36	157.00	1.50	2.25	
150.00	-1.00	1.00	152.00	1.17	1.36	155.00	-0.50	0.25	
151.00			150.83			155.50			152.44
		4.00			18.83			5.50	
	-1.44	2.07		-1.61	2.58		3.06	9.36	

❖  $SS_{\text{between}} = (2.07 + 2.58 + 9.36) \times 6 = 84.06$



# ANOVA

Machine 1	Machine 2	Machine 3
150	153	156
151	152	154
152	148	155
152	151	156
151	149	157
150	152	155
$\bar{x}_1 = 151.00$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$



❖  $SST = SS_{\text{between(or treatment)}} + SS_{\text{within(or error)}}$

❖  $SST = 84.06 + 28.33 = 112.39$

❖ Degrees of freedom

❖  $\text{Total df} = df_{\text{treatment}} + df_{\text{error}}$

❖  $(N-1) = (C-1) + (N-C)$

❖  $df_{\text{treatment}} = 3-1=2, df_{\text{error}} = 18-3=15$

❖  $df_{\text{total}} = 17$



# ANOVA

Machine 1	Machine 2	Machine 3
150	153	156
151	152	154
152	148	155
152	151	156
151	149	157
150	152	155
$\bar{x}_1 = 151.00$	$\bar{x}_2 = 150.83$	$\bar{x}_3 = 155.50$



❖ Mean Sum of Square =  $SS / df$

❖  $MS_{\text{between}} = SS_{\text{between(or treatment)}} / df_{\text{treatment}}$

❖  $MS_{\text{between}} = 84.06 / 2 = 42.03$

❖  $MS_{\text{within}} = SS_{\text{within(or error)}} / df_{\text{within}}$

❖  $MS_{\text{within}} = 28.33 / 15 = 1.89$

❖  $F = MS_{\text{between}} / MS_{\text{within}} = 42.03 / 1.89 = 22.24$



# ANOVA

$$❖ F = MS_{\text{between}} / MS_{\text{within}} = 42.03/1.89 = 22.24$$

❖ Compare this with  $F_{\text{critical}}$

$$❖ F(2, 15, 0.95) = 3.68$$

❖ Reject Null Hypothesis

❖ DEMONSTRATE MS Excel

F - Distribution ( $\alpha = 0.05$  in the Right Tail)

		Numerator Degrees of Freedom								
df <sub>2</sub>	df <sub>1</sub>	1	2	3	4	5	6	7	8	9
1		161.45	199.50	215.71	224.58	230.16	233.99	236.77	238.88	240.54
2		18.513	19.000	19.164	19.247	19.296	19.330	19.353	19.371	19.385
3		10.128	9.5521	9.2766	9.1172	9.0135	8.9406	8.8867	8.8452	8.8123
4		7.7086	7.9443	6.5914	6.3882	6.2561	6.1631	6.0942	6.0410	6.9988
5		6.6079	5.7861	5.4095	5.1922	5.0503	4.9503	4.8759	4.8183	4.7725
6		5.9874	5.1433	4.7571	4.5337	4.3874	4.2839	4.2067	4.1468	4.0990
7		5.5914	4.7374	4.3468	4.1203	3.9715	3.8660	3.7870	3.7257	3.6767
8		5.3177	4.4590	4.0662	3.8379	3.6875	3.5806	3.5005	3.4381	3.3881
9		5.1174	4.2565	3.8625	3.6331	3.4817	3.3738	3.2927	3.2296	3.1789
10		4.9646	4.1028	3.7083	3.4780	3.3258	3.2172	3.1355	3.0717	3.0204
11		4.8443	3.9823	3.5874	3.3567	3.2039	3.0946	3.0123	2.9480	2.8962
12		4.7472	3.8853	3.4903	3.2592	3.1059	2.9961	2.9134	2.8486	2.7964
13		4.6672	3.8056	3.4105	3.1791	3.0254	2.9153	2.8321	2.7669	2.7144
14		4.6001	3.7389	3.3439	3.1122	2.9582	2.8477	2.7642	2.6987	2.6458
15		4.5431	3.6823	3.2874	3.0556	2.9013	2.7905	2.7066	2.6408	2.5876
16		4.4940	3.6337	3.2389	3.0069	2.8524	2.7413	2.6572	2.5911	2.5377
17		4.4513	3.5915	3.1968	2.9647	2.8100	2.6987	2.6143	2.5480	2.4943
18		4.4139	3.5546	3.1599	2.9277	2.7729	2.6613	2.5767	2.5102	2.4563
19		4.3807	3.5219	3.1274	2.8951	2.7401	2.6283	2.5435	2.4768	2.4227
20		4.3512	3.4928	3.0984	2.8661	2.7109	2.5990	2.5140	2.4471	2.3928
21		4.3248	3.4668	3.0725	2.8401	2.6848	2.5727	2.4876	2.4205	2.3660
22		4.3009	3.4434	3.0491	2.8167	2.6613	2.5491	2.4638	2.3965	2.3419
23		4.2793	3.4221	3.0280	2.7955	2.6400	2.5277	2.4422	2.3748	2.3201
24		4.2597	3.4028	3.0088	2.7763	2.6207	2.5082	2.4226	2.3551	2.3002
25		4.2417	3.3852	2.9912	2.7587	2.6030	2.4904	2.4047	2.3371	2.2821
26		4.2252	3.3690	2.9752	2.7426	2.5868	2.4741	2.3883	2.3205	2.2655
27		4.2100	3.3541	2.9604	2.7278	2.5719	2.4591	2.3732	2.3053	2.2501
28		4.1960	3.3404	2.9467	2.7141	2.5581	2.4453	2.3593	2.2913	2.2360
29		4.1830	3.3277	2.9340	2.7014	2.5454	2.4324	2.3463	2.2783	2.2229
30		4.1709	3.3158	2.9223	2.6896	2.5336	2.4205	2.3343	2.2662	2.2107
40		4.0847	3.2317	2.8387	2.6060	2.4495	2.3359	2.2490	2.1802	2.1240
60		4.0012	3.1504	2.7581	2.5252	2.3683	2.2541	2.1665	2.0970	2.0401
120		3.9201	3.0718	2.6802	2.4472	2.2899	2.1750	2.0868	2.0164	1.9588
∞		3.8415	2.9957	2.6049	2.3719	2.2141	2.0986	2.0096	1.9384	1.8799

# Goodness of Fit Test (Chi Square)

- ❖ To test if the sample is coming from a population with specific distribution.
- ❖ Other goodness-of-fit tests are
  - ❖ Anderson-Darling
  - ❖ Kolmogorov-Smirnov
- ❖ Chi Square Goodness of Fit can be used for any type of data: Continuous or Discrete.

# Goodness of Fit Test (Chi Square)

- ❖  $H_0$ : The data follow a specified distribution.
- ❖  $H_a$ : The data do not follow the specified distribution.
- ❖ Calculated Statistic:  $\chi^2 = \sum_{i=1}^k \frac{(O-E)^2}{E}$
- ❖ Critical Statistic: Chi square for k-1 degrees of freedom for specific alpha.

# Goodness of Fit Test (Chi Square)

- ❖ A coin is flipped 100 times. Number of heads are noted. Is this coin biased?

$$\chi^2 = \sum_{i=1}^k \frac{(O-E)^2}{E}$$

Expected	Observed
50	51
50	52
50	56
50	82
50	65



# Goodness of Fit Test (Chi Square)

- ❖ A coin is flipped 100 times. Number of heads are noted. Is this coin biased?

$$\chi^2 = \sum_{i=1}^k \frac{(O-E)^2}{E}$$

Expected	Observed	O-E	(O-E) <sup>2</sup>	(O-E) <sup>2</sup> /E
50	51	1	1	0.02
50	52	2	4	0.08
50	56	6	36	0.72
50	82	32	1024	20.48
50	65	15	225	4.5
				$\chi^2 = 25.8$

# Goodness of Fit Test (Chi Square)

- ❖ A coin is flipped 100 times. Number of heads are noted. Is this coin biased?

$$\chi^2 = \sum_{i=1}^k \frac{(O-E)^2}{E}$$

Percentage Points of the Chi-Square Distribution

Degrees of Freedom	Probability of a larger value of $\chi^2$								
	0.99	0.95	0.90	0.75	0.50	0.25	0.10	0.05	0.01
1	0.000	0.004	0.016	0.102	0.455	1.32	2.71	3.84	6.63
2	0.020	0.103	0.211	0.575	1.386	2.77	4.61	5.99	9.21
3	0.115	0.352	0.584	1.212	2.366	4.11	6.25	7.81	11.34
4	0.297	0.711	1.064	1.923	3.357	5.39	7.78	9.49	13.28
5	0.554	1.145	1.610	2.675	4.351	6.63	9.24	11.07	15.09
6	0.872	1.635	2.204	3.455	5.348	7.84	10.64	12.59	16.81
7	1.239	2.167	2.833	4.255	6.346	9.04	12.02	14.07	18.48
8	1.647	2.733	3.490	5.071	7.344	10.22	13.36	15.51	20.09
9	2.088	3.325	4.168	5.899	8.343	11.39	14.68	16.92	21.67
10	2.558	3.940	4.865	6.737	9.342	12.55	15.99	18.31	23.21
11	3.053	4.575	5.578	7.584	10.341	13.70	17.28	19.68	24.72
12	3.571	5.226	6.304	8.438	11.340	14.85	18.55	21.03	26.22
13	4.107	5.892	7.042	9.299	12.340	15.98	19.81	22.36	27.69
14	4.660	6.571	7.790	10.165	13.339	17.12	21.06	23.68	29.14
15	5.229	7.261	8.547	11.037	14.339	18.25	22.31	25.00	30.58
16	5.812	7.962	9.312	11.912	15.338	19.37	23.54	26.30	32.00
17	6.408	8.672	10.085	12.792	16.338	20.49	24.77	27.59	33.41
18	7.015	9.390	10.865	13.675	17.338	21.60	25.99	28.87	34.80
19	7.633	10.117	11.651	14.562	18.338	22.72	27.20	30.14	36.19
20	8.260	10.851	12.443	15.452	19.337	23.83	28.41	31.41	37.57
22	9.542	12.338	14.041	17.240	21.337	26.04	30.81	33.92	40.29
24	10.856	13.848	15.659	19.037	23.337	28.24	33.20	36.42	42.98
26	12.198	15.379	17.292	20.843	25.336	30.43	35.56	38.89	45.64
28	13.565	16.928	18.939	22.657	27.336	32.62	37.92	41.34	48.28
30	14.953	18.493	20.599	24.478	29.336	34.80	40.26	43.77	50.89
40	22.164	26.509	29.051	33.660	39.335	45.62	51.80	55.76	63.69
50	27.707	34.764	37.689	42.942	49.335	56.33	63.17	67.50	76.15
60	37.485	43.188	46.459	52.294	59.335	66.98	74.40	79.08	88.38

$$X^2_{cal} = 25.8$$

$$X^2_{(4,0.95)} = 9.49$$

# Goodness of Fit Test (Chi Square)

- ❖ A coin is flipped 100 times. Number of heads are noted. Is this coin biased?

$$\chi^2 = \sum_{i=1}^k \frac{(O-E)^2}{E}$$

- ❖  $\chi^2_{\text{cal}} = 25.8$

- ❖  $\chi^2_{(4,0.95)} = 9.49$

- ❖ Reject Null Hypothesis

- ❖ Coin is biased

# Contingency Tables

- ❖ To find relationship between two discrete variables.

	Smoker	Non Smoker	
Male	60	40	100
Female	35	40	75
	95	80	175

	Operator 1	Operator 2	Operator 3	
Shift 1	22	26	23	71
Shift 2	28	62	26	116
Shift 3	72	22	66	160
	122	112	115	347

# Contingency Tables

---

- ❖ Null hypothesis is that there is no relationship between the row and column variables.
- ❖ Alternate hypothesis is that there is a relationship. Alternate hypothesis does not tell what type of relationship exists.

	Operator 1	Operator 2	Operator 3	
Shift 1	22	26	23	71
Shift 2	28	62	26	116
Shift 3	72	22	66	160
	122	112	115	347

# Contingency Tables

- ❖ Calculate Chi square statistic.

$$\chi^2 = \sum_{i=1}^k \frac{(O-E)^2}{E}$$

	Operator 1	Operator 2	Operator 3	
Shift 1	22	26	23	71
Shift 2	28	62	26	116
Shift 3	72	22	66	160
	122	112	115	347

# Contingency Tables

❖ Calculate Chi square statistic.

$$\chi^2 = \sum_{i=1}^k \frac{(O-E)^2}{E}$$

<u>OBSERVED</u>	Operator 1	Operator 2	Operator 3	
Shift 1	22	26	23	71
Shift 2	28	62	26	116
Shift 3	72	22	66	160
	122	112	115	347

<u>EXPECTED</u>	Operator 1	Operator 2	Operator 3	
Shift 1	122x71/347	112x71/347	115x71/347	71
Shift 2	122x116/347	112x116/347	115x116/347	116
Shift 3	122x160/347	112x160/347	115x160/347	160
	122	112	115	347

# Contingency Tables

❖ Calculate Chi square statistic.

$$\chi^2 = \sum_{i=1}^k \frac{(O-E)^2}{E}$$

<u>EXPECTED</u>	Operator 1	Operator 2	Operator 3	
Shift 1	122x71/347	112x71/347	115x71/347	71
Shift 2	122x116/347	112x116/347	115x116/347	116
Shift 3	122x160/347	112x160/347	115x160/347	160
	122	112	115	347

<u>EXPECTED</u>	Operator 1	Operator 2	Operator 3	
Shift 1	24.96	22.91	23.53	71
Shift 2	40.78	37.44	38.44	116
Shift 3	56.25	51.64	53.02	160
	122	112	115	347



# Contingency Tables

❖ Calculate Chi square statistic.

$$\chi^2 = \sum_{i=1}^k \frac{(O-E)^2}{E}$$

<u>OBSERVED</u>	Operator 1	Operator 2	Operator 3	
Shift 1	22	26	23	71
Shift 2	28	62	26	116
Shift 3	72	22	66	160
	122	112	115	347

<u>EXPECTED</u>	Operator 1	Operator 2	Operator 3	
Shift 1	24.96	22.91	23.53	71
Shift 2	40.78	37.44	38.44	116
Shift 3	56.25	51.64	53.02	160
	122	112	115	347

<u>(O-E)<sup>2</sup>/E</u>	Operator 1	Operator 2	Operator 3	
Shift 1	(22-24.96) <sup>2</sup> /24.96 = 0.35	0.42	0.01	71
Shift 2	(28-40.78) <sup>2</sup> /40.78 = 4.00	16.11	4.03	116
Shift 3	(72-56.25) <sup>2</sup> /56.25 = 4.41	17.01	3.18	160
	122	112	115	347

$$\chi^2 = 49.52$$

# Contingency Tables

- ❖ Calculate Chi square statistic = 49.52
- ❖ Degrees of freedom =  $(r-1)(c-1) = 4$
- ❖ Chi square critical = 9.49
- ❖ Reject null hypothesis
- ❖ There is a relationship between the shift and the operator.

Degrees of Freedom	Probability of a larger value of $\chi^2$								
	0.99	0.95	0.90	0.75	0.50	0.25	0.10	0.05	0.01
1	0.000	0.004	0.016	0.102	0.455	1.32	2.71	3.84	6.63
2	0.020	0.103	0.211	0.575	1.386	2.77	4.61	5.99	9.21
3	0.115	0.352	0.584	1.212	2.366	4.11	6.25	7.81	11.34
4	0.297	0.711	1.064	1.923	3.357	5.39	7.78	9.49	13.28
5	0.554	1.145	1.610	2.675	4.351	6.63	9.24	11.07	15.09
6	0.872	1.635	2.204	3.455	5.348	7.84	10.64	12.59	16.81
7	1.239	2.167	2.833	4.255	6.346	9.04	12.02	14.07	18.48
8	1.647	2.733	3.490	5.071	7.344	10.22	13.36	15.51	20.09
9	2.088	3.325	4.168	5.899	8.343	11.39	14.68	16.92	21.67
10	2.558	3.940	4.865	6.737	9.342	12.55	15.99	18.31	23.21
11	3.053	4.575	5.578	7.584	10.341	13.70	17.28	19.68	24.72
12	3.571	5.226	6.304	8.438	11.340	14.85	18.55	21.03	26.22
13	4.107	5.892	7.042	9.299	12.340	15.98	19.81	22.36	27.69
14	4.660	6.571	7.790	10.165	13.339	17.12	21.06	23.68	29.14
15	5.229	7.261	8.547	11.037	14.339	18.25	22.31	25.00	30.58
16	5.812	7.962	9.312	11.912	15.338	19.37	23.54	26.30	32.00
17	6.408	8.672	10.085	12.792	16.338	20.49	24.77	27.59	33.41
18	7.015	9.390	10.865	13.675	17.338	21.60	25.99	28.87	34.80
19	7.633	10.117	11.651	14.562	18.338	22.72	27.20	30.14	36.19
20	8.260	10.851	12.443	15.452	19.337	23.83	28.41	31.41	37.57
22	9.542	12.338	14.041	17.240	21.337	26.04	30.81	33.92	40.29
24	10.856	13.848	15.659	19.037	23.337	28.24	33.20	36.42	42.98
26	12.198	15.379	17.292	20.843	25.336	30.43	35.56	38.89	45.64
28	13.565	16.928	18.939	22.657	27.336	32.62	37.92	41.34	48.28
30	14.953	18.493	20.599	24.478	29.336	34.80	40.26	43.77	50.89
40	22.164	26.509	29.051	33.660	39.335	45.62	51.80	55.76	63.69
50	27.707	34.764	37.689	42.942	49.335	56.33	63.17	67.50	76.15
60	37.485	43.188	46.459	52.294	59.335	66.98	74.40	79.08	88.38

# Correlation

---

- ❖  $Y = f(X)$ ,
  - ❖ where Y is Dependent variable or the result (output)
  - ❖ X is Independent variable, input or the controllable variable
  
- ❖ For example in the study of marks obtained by students in a subject (Y) vs hours of study (X)

# Correlation

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n(\sum x^2) - (\sum x)^2][n(\sum y^2) - (\sum y)^2]}}$$

Hours Studied (X)	Test Score % (Y)
20	40
24	55
46	69
62	83
22	27
37	44
45	61
27	33
65	71
23	37

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

# Correlation

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n(\sum x^2) - (\sum x)^2][n(\sum y^2) - (\sum y)^2]}}$$

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

	Hours Studied (X)	Test Score % (Y)	XY	X2	Y2
	20	40	800	400	1600
	24	55	1320	576	3025
	46	69	3174	2116	4761
	62	83	5146	3844	6889
	22	27	594	484	729
	37	44	1628	1369	1936
	45	61	2745	2025	3721
	27	33	891	729	1089
	65	71	4615	4225	5041
	23	37	851	529	1369
<b>SUM</b>	<b>371</b>	<b>520</b>	<b>21764</b>	<b>16297</b>	<b>30160</b>

# Correlation Coefficient

---

- ❖ **Correlation**
- ❖ Measures the strength of linear relationship between Y and X
- ❖ Pearson Correlation Coefficient,  $r$  ( $r$  varies between -1 and +1)
  - ❖ Perfect positive relationship:  $r = 1$
  - ❖ No relationship:  $r = 0$
  - ❖ Perfect negative relationship:  $r = -1$

# Correlation Coefficient

---



# Correlation vs Causation

---

- ❖ **Correlation does not imply causation**
  - ❖ a correlation between two variables does not imply that one causes the other



# Coefficient of Determination

- ❖ Coefficient of Determination,  $r^2$
- ❖ Proportion of the variance in the dependent variable that is predictable from the independent variable
- ❖ (varies from 0.0 to 1.0 or zero to 100%)
  - ❖ None of the variation in Y is explained by X,  $r^2 = 0.0$
  - ❖ All of the variation in Y is explained by X,  $r^2 = 1.0$
- ❖  $r = 0.88, r^2 = 0.77$

# Regression Analysis

---

- ❖ Quantifies the relationship between Y and X ( $Y = a + bX$ )

# Regression Analysis

- ❖ Quantifies the relationship between Y and X ( $Y = a + bX$ )

	Hours Studied (X)	Test Score % (Y)	XY	X <sup>2</sup>	Y <sup>2</sup>
	20	40	800	400	1600
	24	55	1320	576	3025
	46	69	3174	2116	4761
	62	83	5146	3844	6889
	22	27	594	484	729
	37	44	1628	1369	1936
	45	61	2745	2025	3721
	27	33	891	729	1089
	65	71	4615	4225	5041
	23	37	851	529	1369
SUM	371	520	21764	16297	30160

$$Y = a + bX$$

$$b = \frac{N\sum XY - (\sum X)(\sum Y)}{N\sum X^2 - (\sum X)^2} \quad a = \frac{\sum Y - b\sum X}{N}$$

# Regression Analysis

- ❖ Quantifies the relationship between Y and X ( $Y = 15.79 + 0.97.X$ )

$$Y = a + bX$$

$$b = \frac{N\sum XY - (\sum X)(\sum Y)}{N\sum X^2 - (\sum X)^2} \quad a = \frac{\sum Y - b\sum X}{N}$$

	Hours Studied (X)	Test Score % (Y)	XY	X2	Y2
	20	40	800	400	1600
	24	55	1320	576	3025
	46	69	3174	2116	4761
	62	83	5146	3844	6889
	22	27	594	484	729
	37	44	1628	1369	1936
	45	61	2745	2025	3721
	27	33	891	729	1089
	65	71	4615	4225	5041
	23	37	851	529	1369
SUM	371	520	21764	16297	30160

# Regression Analysis

---

- ❖ For a student studying 50 hrs what is the expected test score %?

$$Y = a + bX$$

$$b = \frac{N\sum XY - (\sum X)(\sum Y)}{N\sum X^2 - (\sum X)^2} \quad a = \frac{\sum Y - b\sum X}{N}$$