

Data Analysis and Modeling in Bioinformatics

Part I of II

Module Outline

Probability Theory
Descriptive Statistics
Inferential Statistics

Shebna Rose Fabilloren

Linear Regression Analysis
Correlation
Data Structures

Dr. Jan Michael Yap

Probability Theory

- Describes the likelihood of a particular **outcome** for an **experiment**.
Probability is a number quantified between 0 and 1.
 - 0 = impossible
 - 1 = certainty
- The higher the probability, the more likely an **event** will occur.



Probability Theory

Terminologies

Experiment/Trial – is any procedure:

- a. with any well defined set of possible outcomes.
- b. whose actual outcome is not known in advance.

Sample Outcome(s) – one of the possible outcomes of an experiment/trial.

Sample Space(S) – entire set of possible outcomes

Event(E) – set of outcomes of the experiment.

- can have one outcome or more than one outcome.

Probability Theory

- Describes the likelihood of a particular outcome for an experiment.
Probability is a number quantified between 0 and 1.
 - 0 = impossible
 - 1 = certainty
- The higher the probability, the more likely an event will occur.

$$\text{Probability of an event} = \frac{\text{Number of ways it can happen}}{\text{Total number of outcomes}}$$

Probability Theory

2 Properties of Probabilities

1. Probabilities cannot be negative.
2. If we consider the probabilities of all possible outcomes, then the sum of their probabilities must equal 1.0

Probability Theory

Terminologies

Mutually Exclusive Events - events that can't happen at the same time.

Independent Events - probability that one event occurs in no way affects the probability of the other event occurring

Dependent Events - the outcome of the first event affects the outcome of the second event, so that the probability is changed



Probability Theory

Set Notation

| | | |
|--------------|-------------|--------|
| Union | \cup | OR |
| Intersection | \cap | AND |
| Subset | \subseteq | subset |

Probability Theory

Marginal Probability – the probability of an event occurring.

$$P(A)$$

Joint Probability – the probability of event A and event B occurring.

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

Marginal probabilities

| | Pass | Fail | Total |
|---------|------|------|-------|
| Males | 46 | 56 | 102 |
| Females | 68 | 30 | 98 |
| Total | 114 | 86 | 200 |

$$P(\text{male}) = 0.51$$

$$P(\text{female}) = 0.49$$

$$P(\text{passed}) = 0.57$$

$$P(\text{failed}) = 0.43$$

Probability Theory

Conditional Probability

- is a measure of the probability of an event given that another event **has already occurred**.

If the event of interest is A and the event B is known or assumed to have occurred.

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$

Probability Theory

Conditional Probability Example

| | Have pets | Do not have pets | Total |
|--------|-----------|------------------|-------|
| Male | 0.41 | 0.08 | 0.49 |
| Female | 0.45 | 0.06 | 0.51 |
| Total | 0.86 | 0.14 | 1 |

Click to add text

$$P(M|PO) = P(M \cap PO) / P(PO)$$

$$P(M \cap PO) = 0.41$$

$$P(PO) = 0.86$$

$$P(M|PO) = 0.48$$

Probability Theory

Bayes' Rule

- provides us with a way to update our beliefs based on the arrival of new, relevant pieces of evidence.

$$P(H|E) = P(E|H)P(H) / P(E)$$

Example:

$$P(\text{cancer}) = 0.05$$

$$P(\text{smoker}) = 0.10$$

$$P(\text{smoker}|\text{cancer}) = 0.20$$

Before we found out that the person smokes, our **prior** was 0.05. However, using new evidence, we can instead calculate

$$P(\text{cancer}|\text{smoker}) = \frac{P(\text{smoker}|\text{cancer}) \cdot P(\text{cancer})}{P(\text{smoker})} = \frac{(0.20)(0.05)}{0.10}$$

Probability Theory

Random Variable

- usually written “X”, is a variable that can take multiple values depending on the outcome of a random event.
- Possible outcomes are the possible values taken by the variable.

Example:

Probability of getting a 3 after throwing a die

$$P(X=3) = 1/6$$

Probability of getting a value greater than 2 after throwing a die

$$P(X > 2) = 4/6$$

Probability Theory

2 Types of Random Variables

1. **Discrete** - one which may take on only a **countable** number of distinct values such as 0,1,2,3,4,...

Ex. # of children in the family

2. **Continuous** - one which takes an **infinite** number of possible values.
- usually measurements

Ex. Height

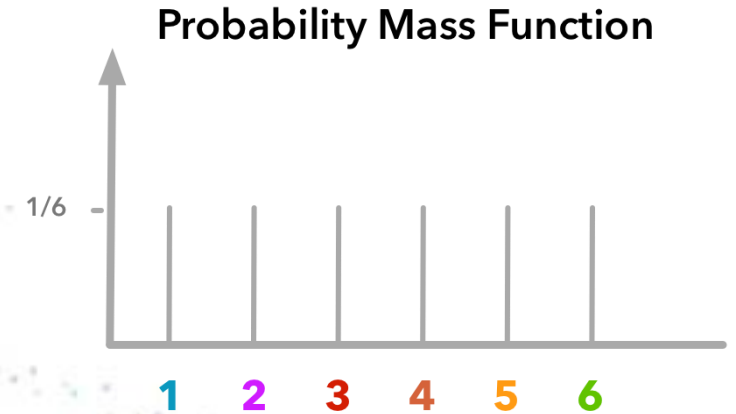
Probability Theory

Probability Distributions

- description of the probability of each possible value that a random variable can take.

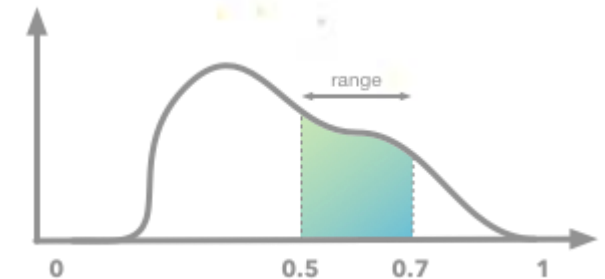
1. Probability Mass Function

- Probability distribution of a discrete variable



2. Probability Density Function

- Probability distribution of a continuous variable
- Area under the curve

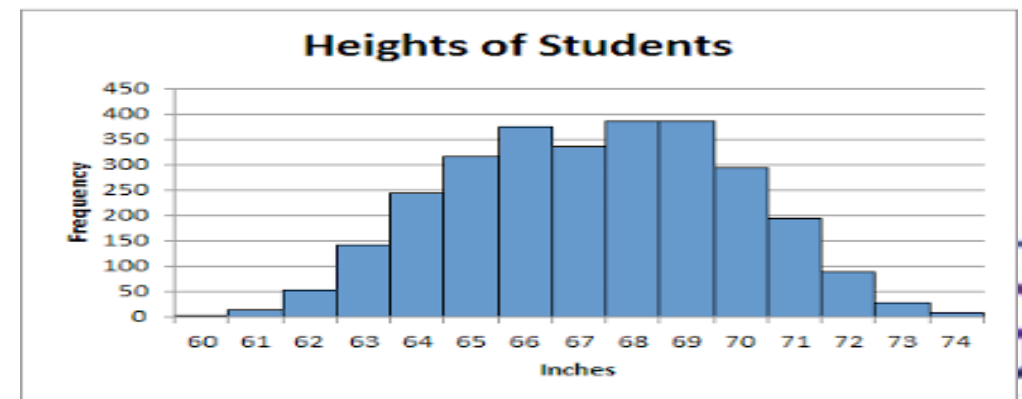
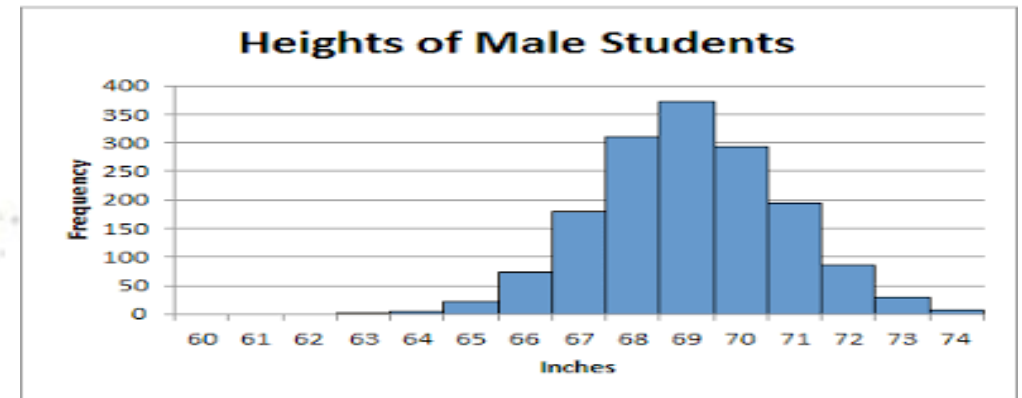
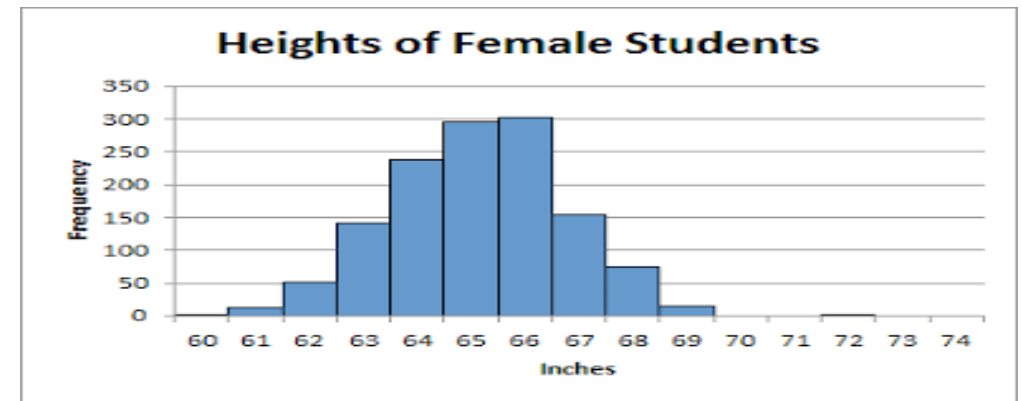


Statistics

- deals with data collection, organization, analysis, interpretation and presentation.

Descriptive Statistics

- describe what the data show
- tells us about the middle of the data and how spread out they are



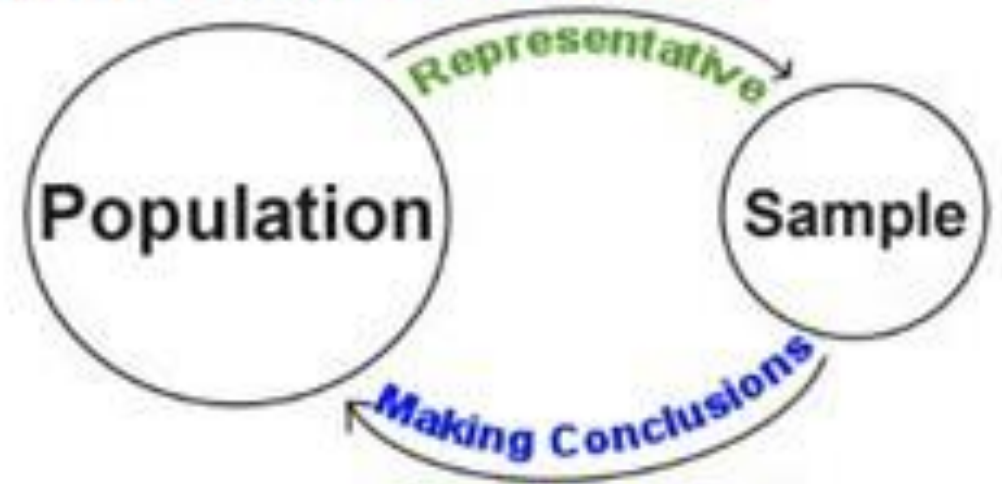
Statistics

- deals with data collection, organization, analysis, interpretation and presentation.

Inferential Statistics

- allows us to make inferences
- allows us to make conclusions that extend beyond the data we have in hand.

Inferential Statistics



Descriptive Statistics

Measures of Central Tendency

1. Mean
2. Median
3. Mode

Measures of Variation

1. Variance
2. Standard Deviation

Descriptive Statistics

Measures of Central Tendency

- Describes how a set of numbers is centered around a particular point on a line scale.

- Answers the question:

Where (around what value)

do the numbers bunch together?

1. Mean

- Is the average value that would be observed for the random variable if it could be observed over and over again an infinite number of times.

2. Median

- Value that divides the set into equal halves when all the numbers have been

Descriptive Statistics

Measures of Central Tendency

1. Mean

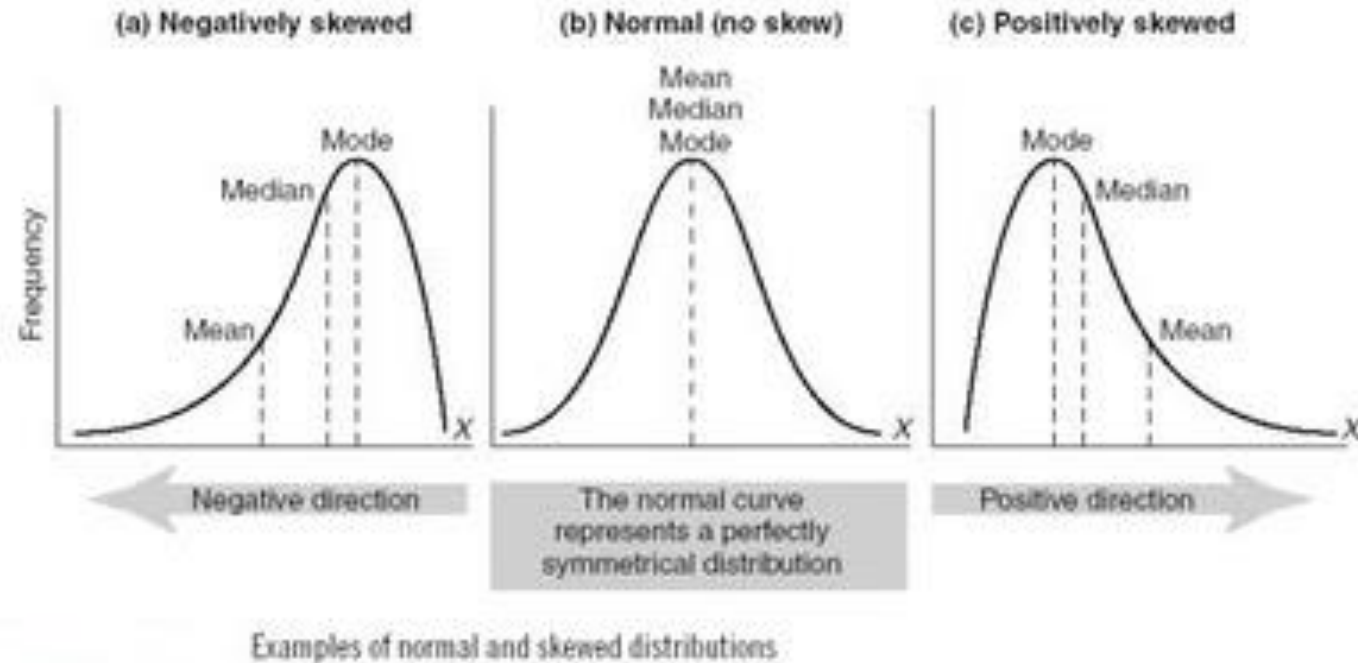
- Is the average value that would be observed for the random variable if it could be observed over and over again an infinite number of times.

2. Median

- Value that divides the set into equal halves when all the numbers have been ordered from lowest to highest.

3. Mode

- Most frequently occurring number in a set of values.



Descriptive Statistics

Measures of Dispersion

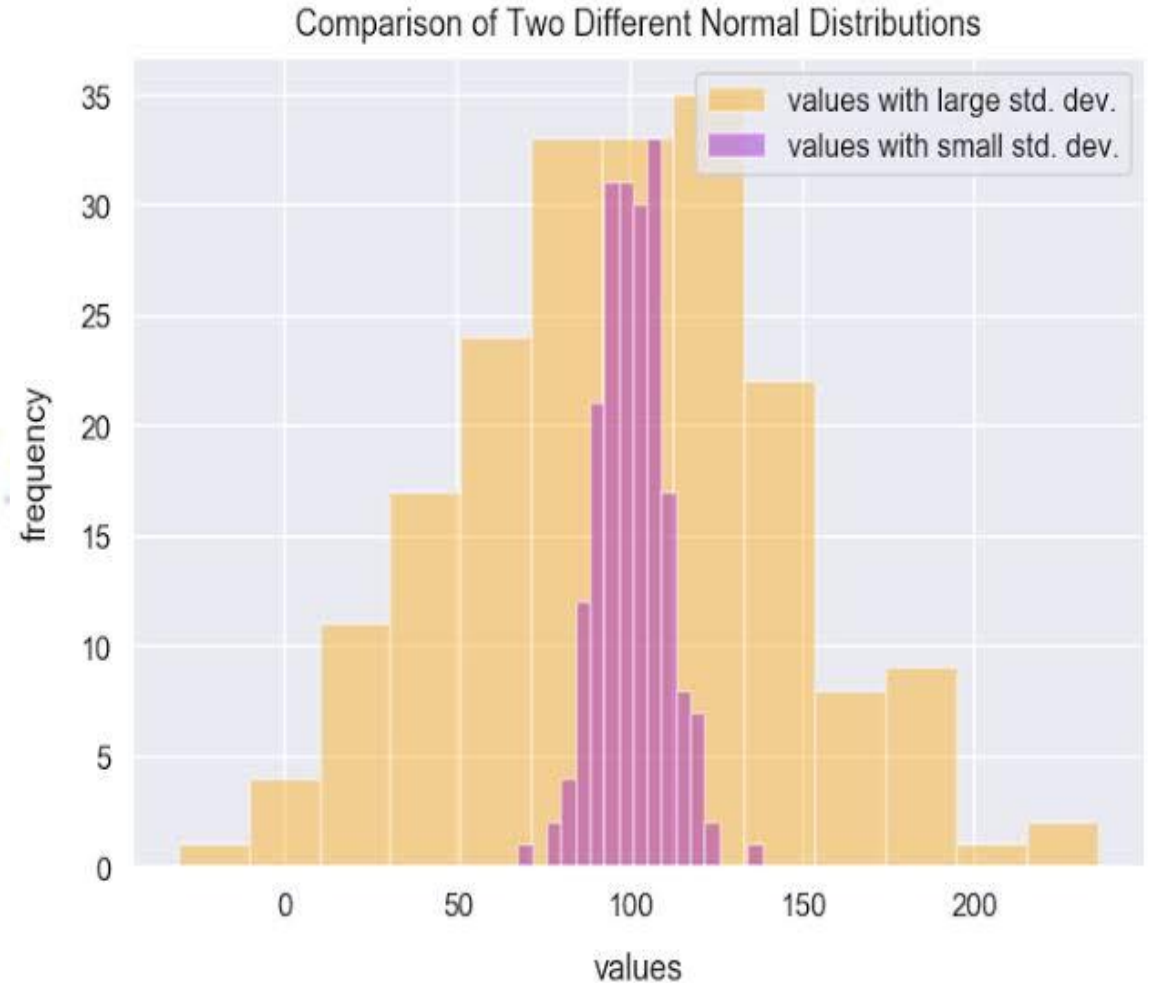
- Tells us how far the numbers are scattered about the center value of the set.

1. Variance

- is the average of the square of the deviations of a set of scores from their mean.

2. Standard Deviation

- square root of the variance



Different Probability Distributions

Bernoulli Distribution

Binomial Distribution

Discrete Uniform Distribution

Geometric Distribution

Poisson Distribution

Normal Distribution

Continuous Uniform Distribution

T Distribution

Chi-square Distribution

Common Discrete Probability Distributions

Discrete Probability Distributions

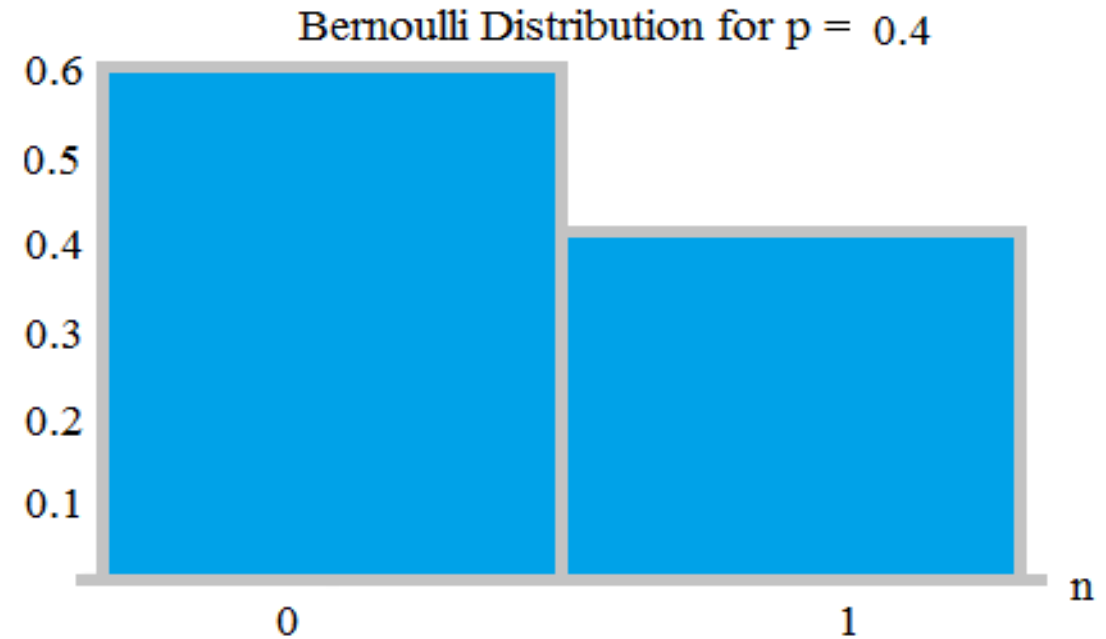
Bernoulli Distribution – probability distribution of a Bernoulli random variable.

$$P(X = x) = \begin{cases} p & \text{for } x = 1 \\ 1 - p & \text{for } x = 0 \end{cases}$$

Bernoulli Trial

One experiment with two possible outcomes.

- Success = 1
- Failure = 0



Discrete Probability Distributions

Binomial Distribution – probability distribution of a Binomial random variable.

Binomial Random Variable

- counts **how often** a particular **event** occurs in a **fixed** number of **trials**.

Conditions:

1. Fixed number of trials
2. Each trial can only have two outcomes, success or failure.
3. Trials are independent of each other. (The same probability of occurrence in each trial)

Fun fact:

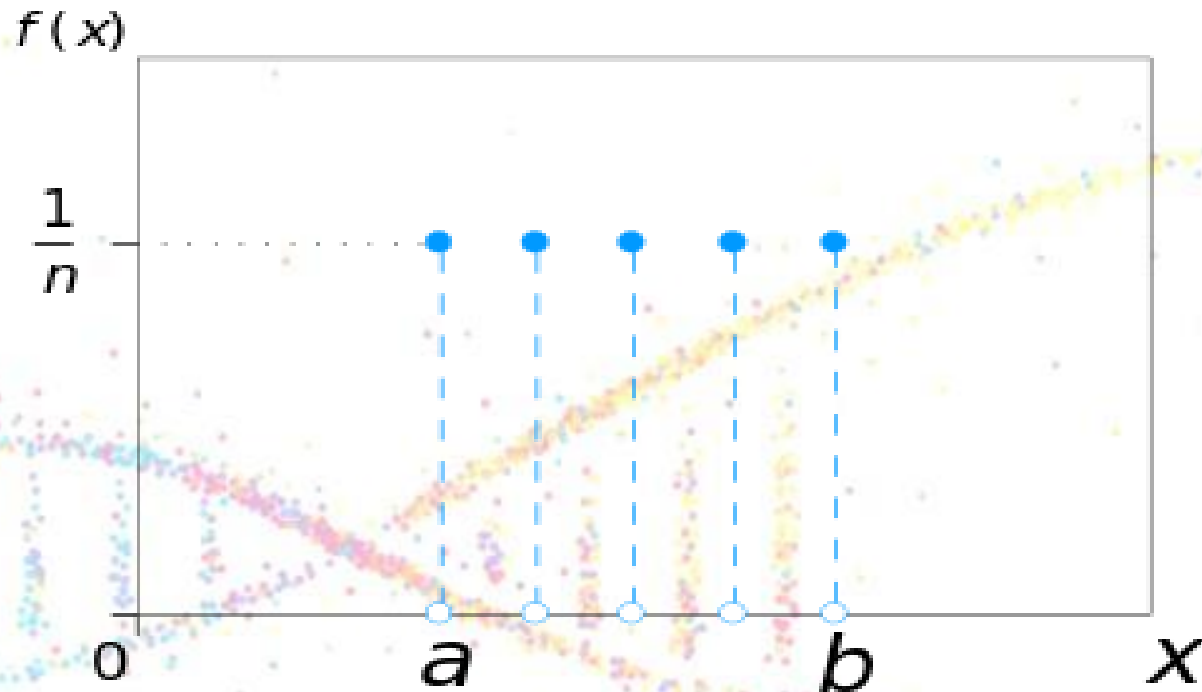
Bernoulli Distribution is the Binomial Distribution but with only one trial.

Discrete Probability Distributions

Discrete Uniform Distribution

- a known, finite number of outcomes equally likely to happen
- every one of n values has equal probability $1/n$.

$$P(X) = 1/n$$



Discrete Probability Distributions

Geometric Distribution – probability distribution of a geometric random variable

Geometric Random Variable

- represents the number of failures before you get a success in a series of Bernoulli trials.

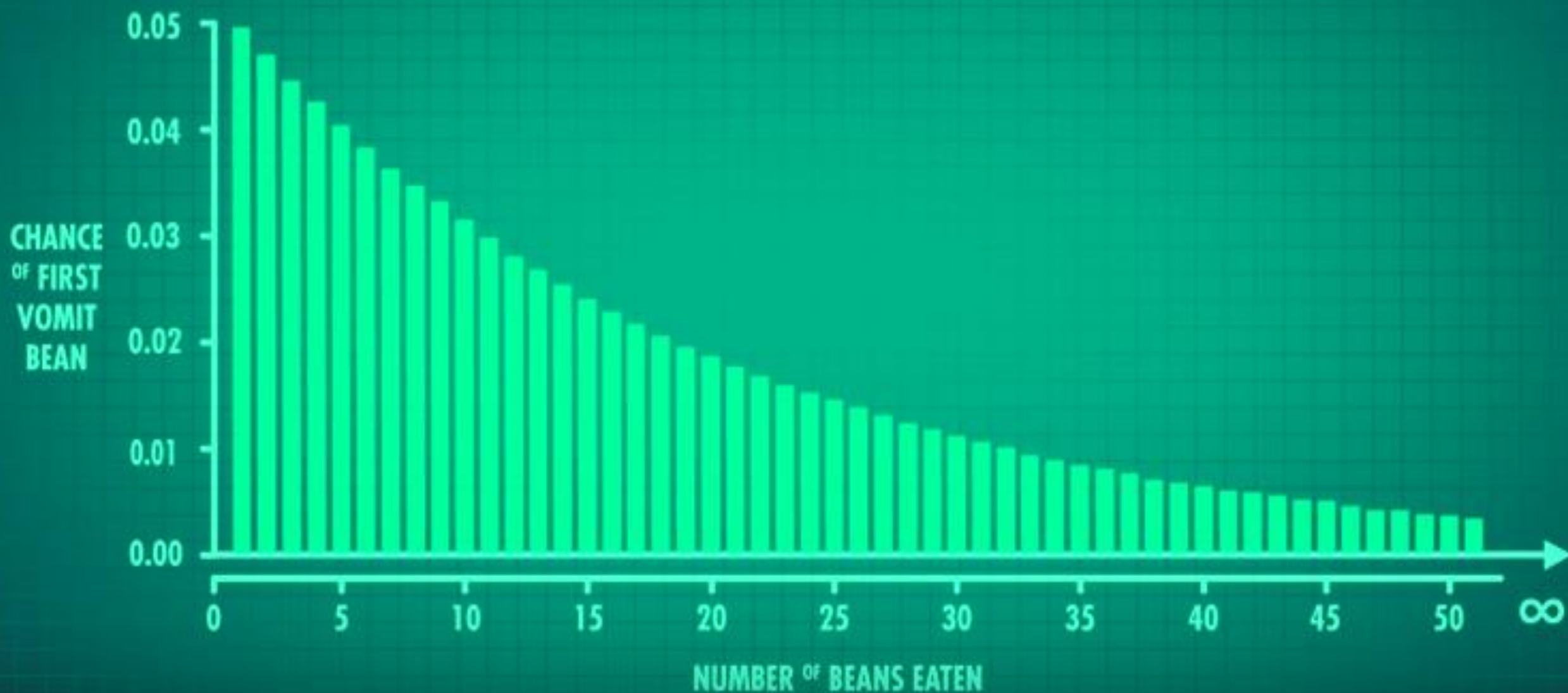
$P(X) = (1-p)^{x-1} p$, where p is success.

Example:

Probability of getting a heads on the 8th coin toss.

$$P(X=8) = (1-0.5)^{8-1} (0.5)$$

GEOMETRIC DISTRIBUTION



Discrete Probability Distributions

Poisson Distribution – probability distribution of a poisson random variable.

Poisson random variable - counts the number of events occurring in a given time period, given the average number of times the event occurs over that time period.

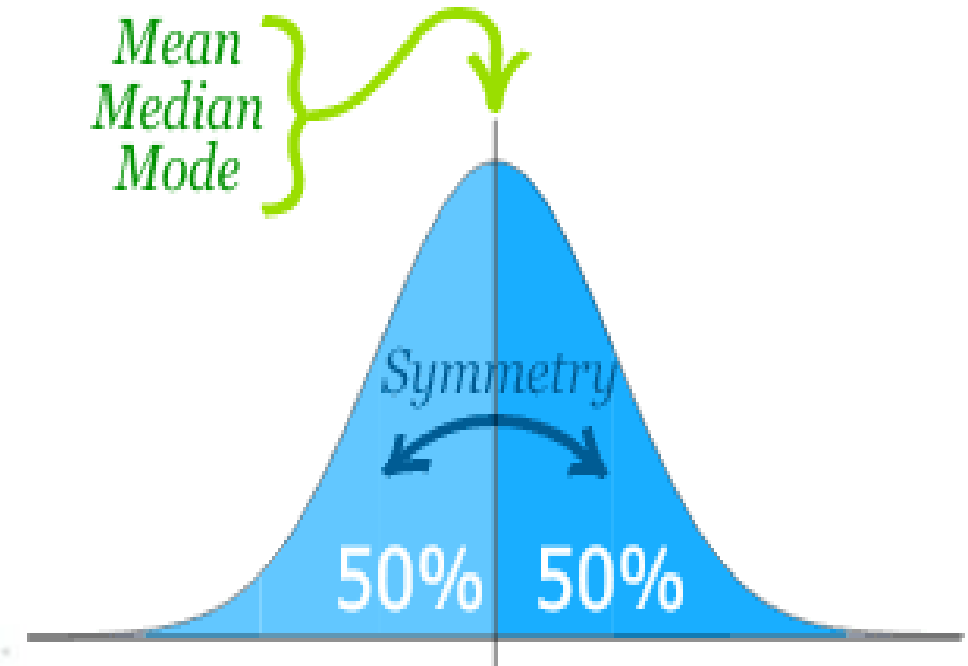
Ex. The daily average number of car accidents in Metro Manila is 299. What is the probability that exactly 142 accidents will happen tomorrow?

Common Continuous Probability Distributions

Continuous Probability Distributions

Normal Distribution

- bell-shaped
- symmetric
- most of the observations cluster around the central peak



Important in statistics, because of **Central Limit Theorem**.

Continuous Probability Distributions

Normal Distribution

Central Limit Theorem

States that the sampling distribution of the sample means approaches a normal distribution as the sample size gets larger — no matter what the shape of the population distribution.

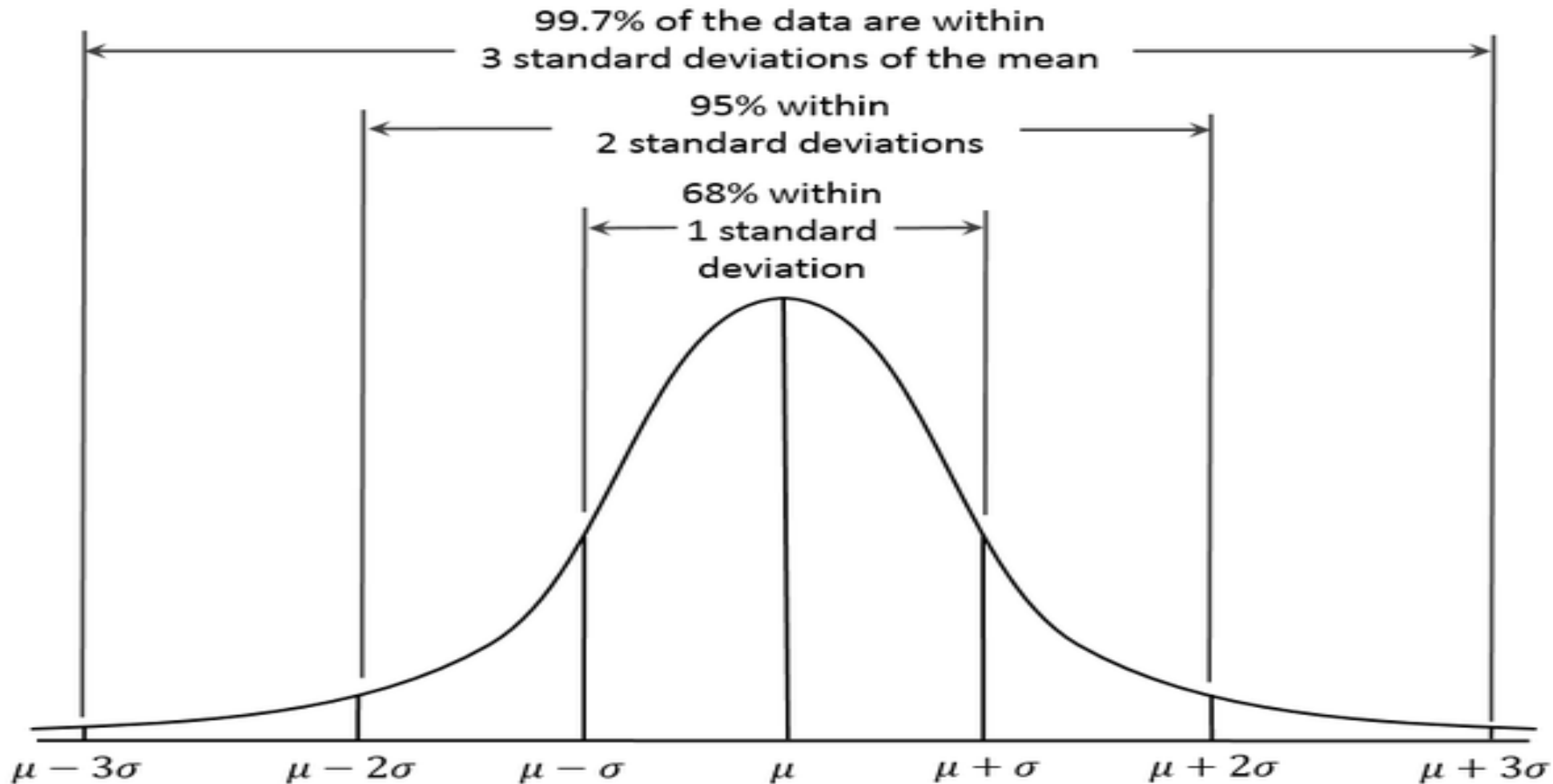
Demo:

<https://www.youtube.com/watch?v=dlbkaurTAUg>

Continuous Probability Distributions

Normal Distribution

Empirical Rule



Continuous Probability Distributions

Standard Normal Distribution

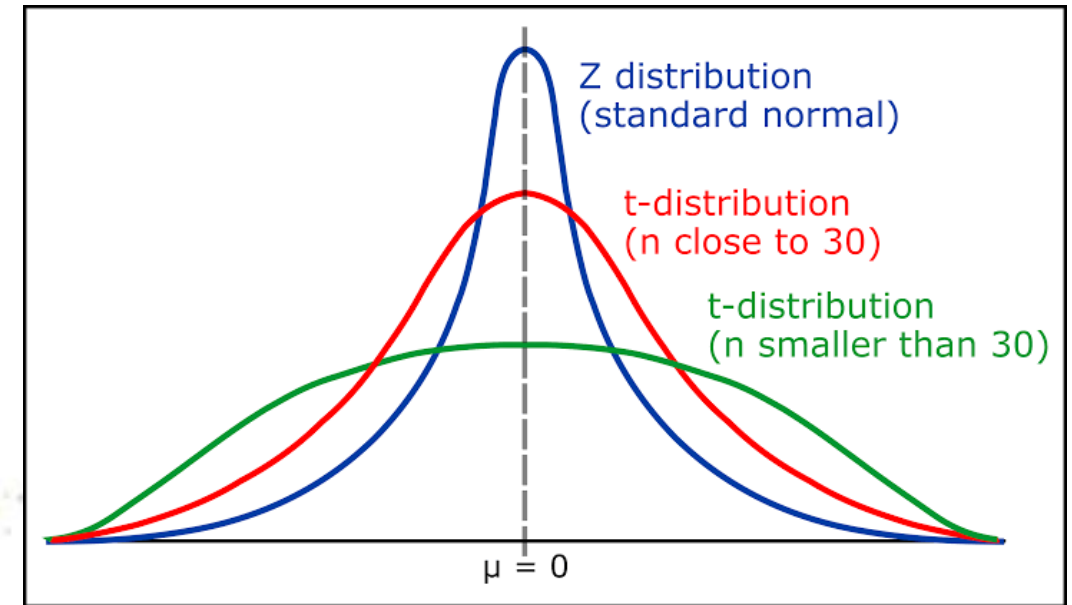
- Specific instance of normal distribution that has a
mean = 0
standard deviation = 1.
- Normal distribution is converted into standard normal distribution in order to utilize the standard normal distribution table to find areas under the normal curve.



Continuous Probability Distributions

T-Distribution

- Distribution of the T-test statistic
- Looks almost identical to the normal distribution(except that its shorter, and fatter)
- Used instead of the normal distribution when you have **small samples(t-score)**.
- The larger the sample size, the more the t distribution looks like the normal distribution.



Hypothesis Testing



Hypothesis Testing

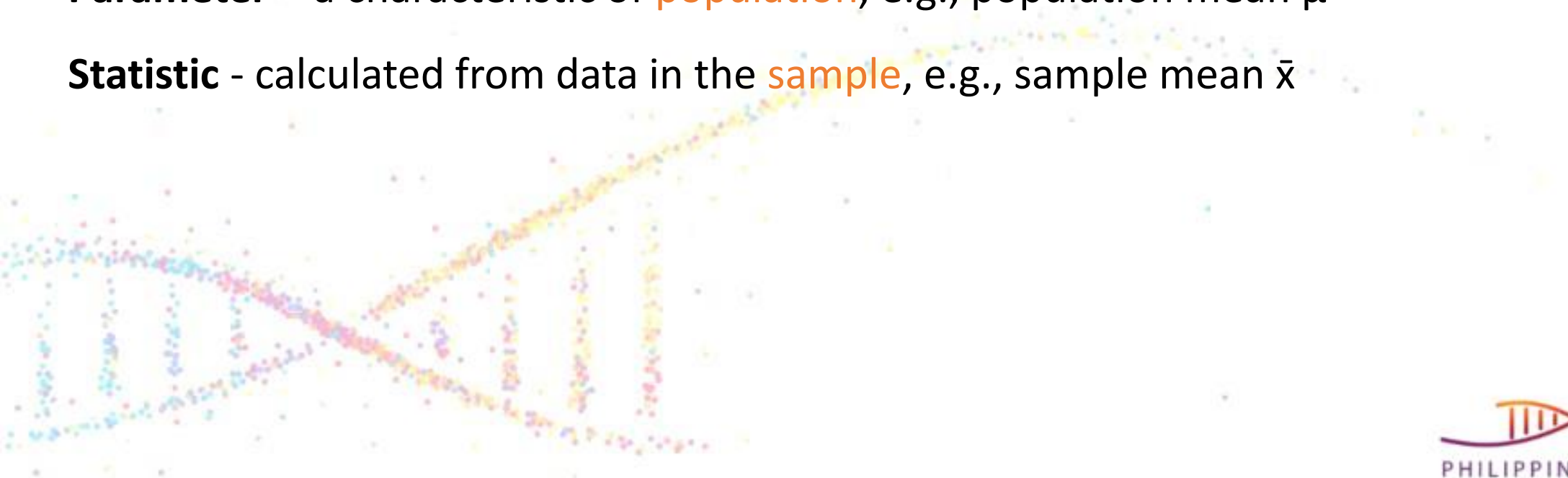
Population - all possible values

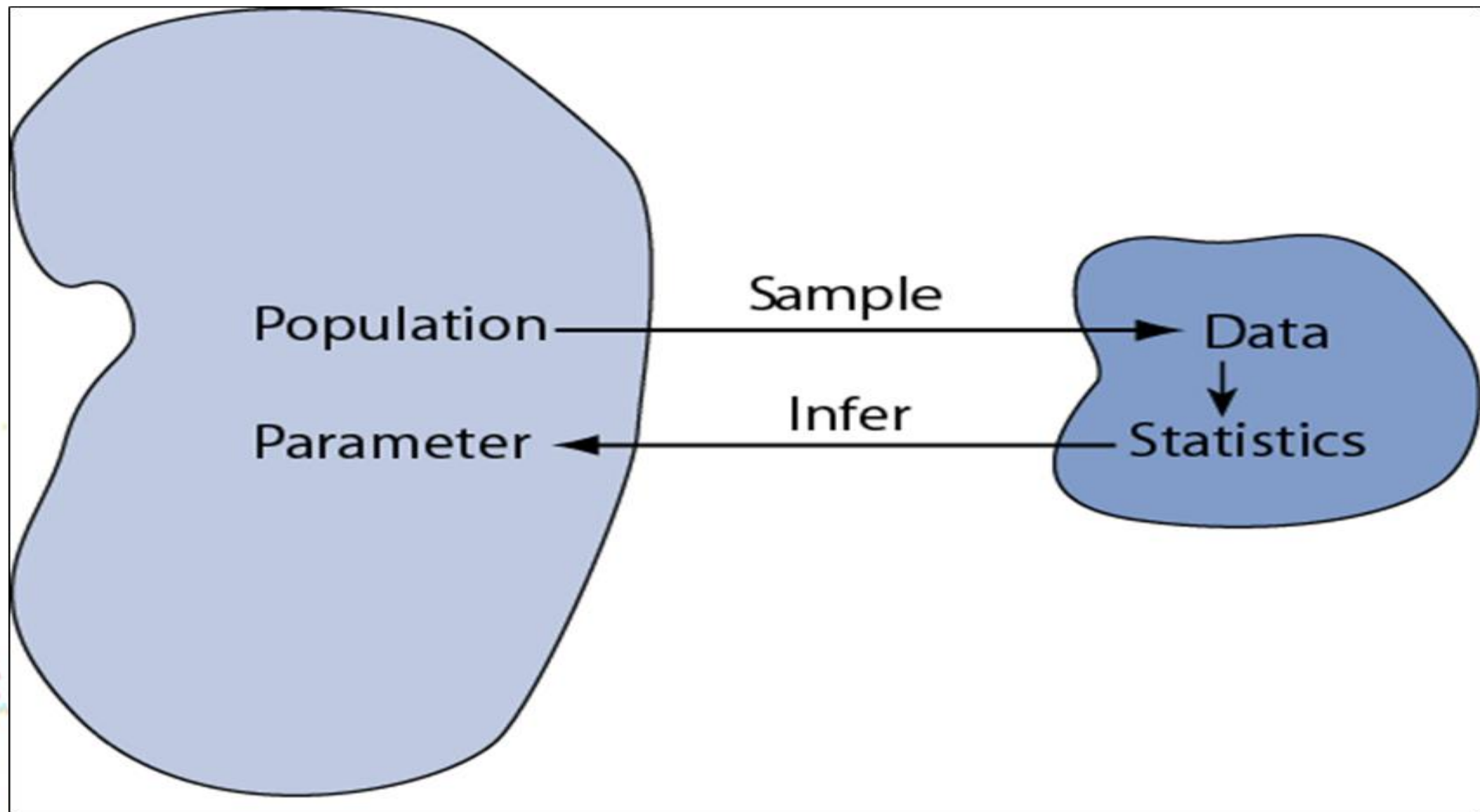
Sample - a portion of the population

Statistical inference - generalizing from a sample to a population with calculated degree of certainty

Parameter - a characteristic of **population**, e.g., population mean μ

Statistic - calculated from data in the **sample**, e.g., sample mean \bar{x}





Hypothesis Testing

6 Steps in Statistical Hypothesis Testing

Step 1: State the Null **Hypothesis**

Step 2: State the Alternative **Hypothesis**

Step 3: Set α (level of significance)

Step 4: Calculate a **test** statistic

Step 5: Compare calculated test statistic with cut-off p-value

Step 6: Draw a conclusion about H_0

Hypothesis Testing

2 Types of Hypothesis

1. Null hypothesis (H_0)

- is a claim of “no difference in the population”
- currently accepted value for a parameter

2. Alternative hypothesis (H_a)

- claims “ H_0 is false”
- research hypothesis

Hypothesis Testing

It is believed that a candy machine makes chocolate bars that are on average 5g. A worker claims that the machine after maintenance no longer makes 5g bars. Write the H_0 and H_a .

$$H_0 : \mu = 5$$

$$H_a : \mu \neq 5$$

Note:

Null and alternative hypothesis are mathematical opposites.

Hypothesis Testing

Confidence Level

- “C”
- 90%, 95% , 99%
- How confident are we in our decision?



Hypothesis Testing

Significance Level(α)

- Probability of erroneously rejecting the H_0
- Test statistic falls in the critical region when the H_0 is actually true



Hypothesis Testing

Errors in Hypothesis Testing

| | In Reality | |
|----------------------|--|---|
| Decision | H_0 is TRUE | H_0 is FALSE |
| Fail to Reject H_0 | OK | Type II Error β = probability of Type II Error |
| Reject H_0 | Type I Error α = probability of Type I Error | OK |

Hypothesis Testing

Errors in Hypothesis Testing

Type I error
(false positive)



Type II error
(false negative)



Hypothesis Testing

Calculate test statistic

Critical Value

- is a point on the test distribution that is compared to the test statistic to determine whether to reject the null hypothesis.
- Where do we draw the line to make a decision?

Test Statistic

- is *calculated from sample data* and used to decide.
- Equation depends on the type of problem you have.

Hypothesis Testing

Two-tailed, Left-tailed, Right-tailed Tests

The *tails* in a distribution are the extreme regions bounded by critical values

Two-Tailed Versus One-Tailed Hypothesis Tests

Figure A:
Two-Tailed Test

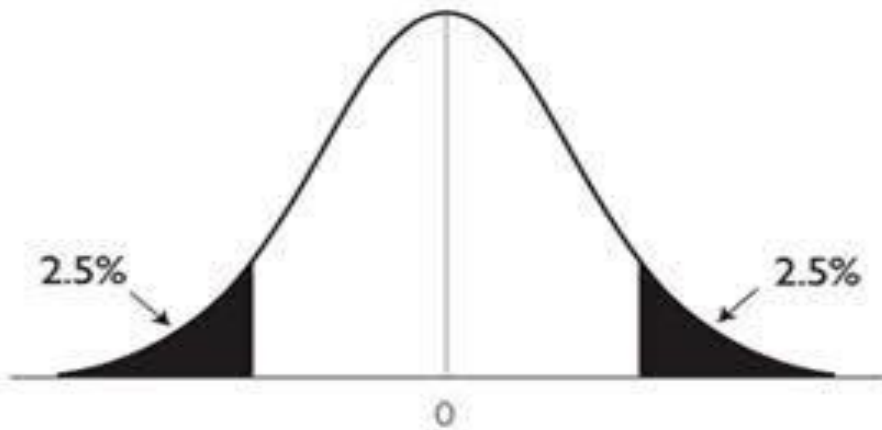
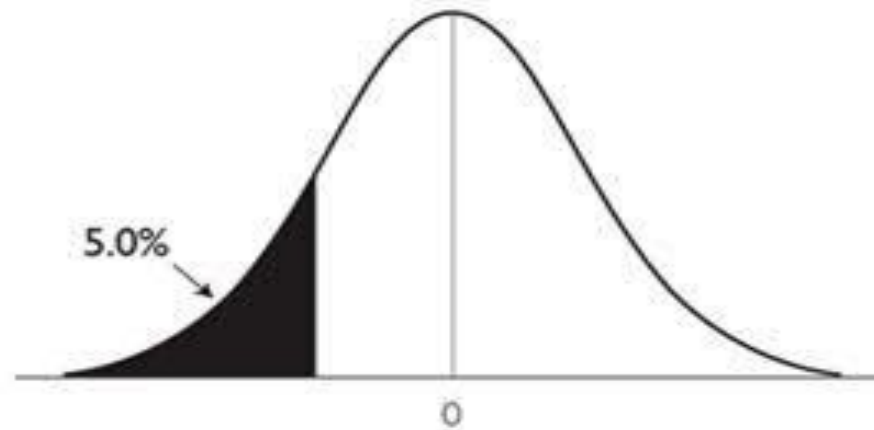


Figure B:
One-Tailed Test
(Left-Tailed Test)



Source: The Heritage Foundation.

Hypothesis Testing

- **Hypothesis** - assumption about a population parameter
- Also called *significance testing*
- Tests a claim about a parameter using data in a sample



Hypothesis Testing

Null Hypothesis Significance Testing (NHST)

- tries to discredit an idea by assuming an idea is true and then showing that if you make that assumption, something contradictory happens.
- In short, proof by contradiction.



Hypothesis Testing

Concept of P-Value

Probability of obtaining a sample more extreme than the ones observed in your data, assuming that H_0 is true.

Answers the question: How rare is your data?

How? By telling you the probability of getting data as extreme as the data you observed if the null hypothesis is true.

Ex. P-value = 0.06

Your data is in the top 6% most extreme samples we'd expect to see based on the distribution of sample means.

Hypothesis Testing

Hypothesis Testing and p-value

p-values need a cut-off.

Usually set as 0.05.

If the p-value is less than our cut-off, it's sufficient evidence to reject the H_0 .

Rejecting H_0 means result is statistically significant.

Statistically significant = unlikely due to random chance alone

Note: Always report your p-value!

Hypothesis Testing

Decision Criterion

Reject H_0 if the p-value $\leq \alpha$ (where α is the significance level, such as 0.05).

Fail to reject H_0 if the p-value $> \alpha$.



Start

What type of test ?

Left-tailed

Right-tailed

Two-tailed

Is the test statistic to the right or left of center ?

Left

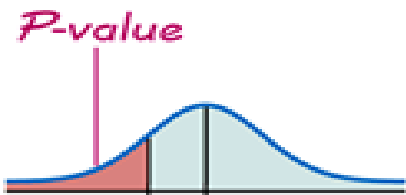
Right

P -value = area to the left of the test statistic

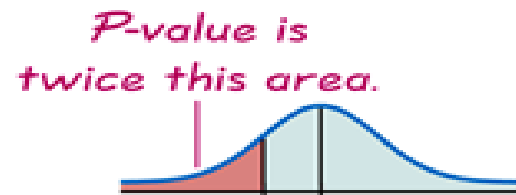
P -value = twice the area to the left of the test statistic

P -value = twice the area to the right of the test statistic

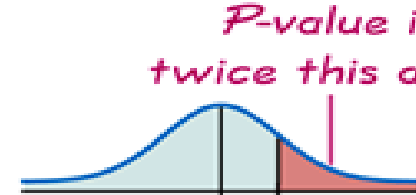
P -value = area to the right of the test statistic



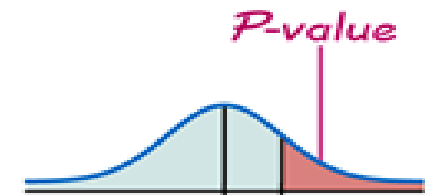
Test statistic



Test statistic

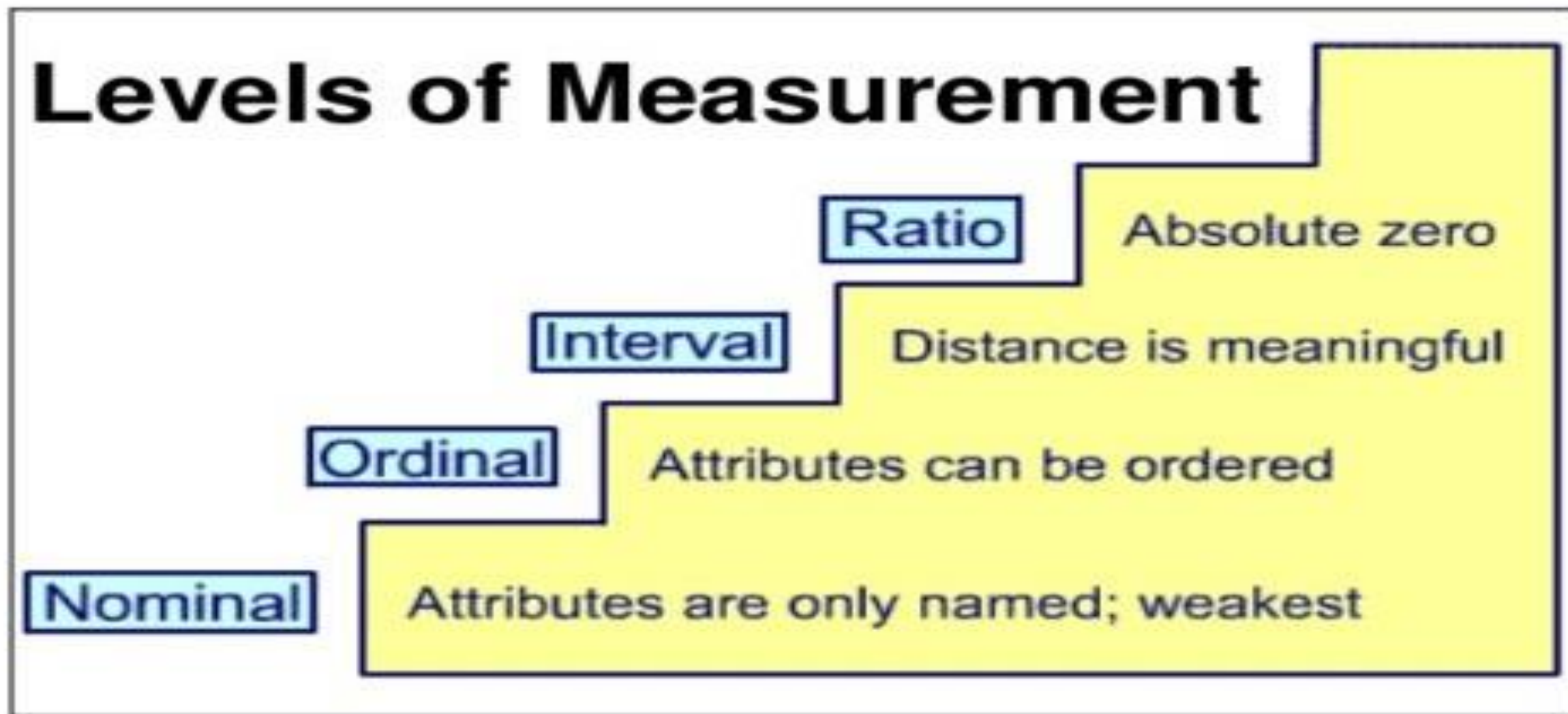


Test statistic



Test statistic

Levels of Measurement (Types of Data)



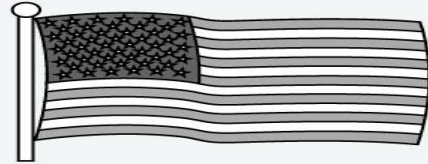
Levels of Measurement (Types of Data)

Exhibit 4.3

Levels of Measurement

Qualitative

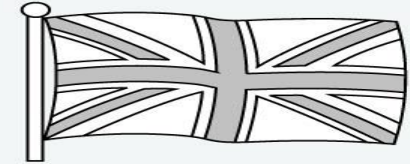
Nominal or categorical level of measurement:
Nationality



American



Canadian



British

Quantitative

Ordinal level of measurement:
Level of conflict

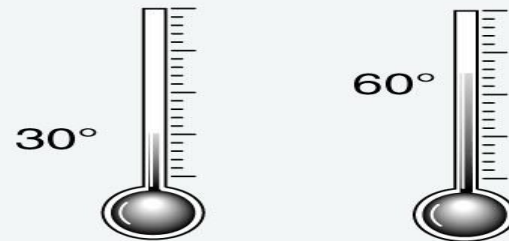


Low



High

Interval level of measurement:
Temperature in degrees Fahrenheit



Ratio level of measurement:
Group size



5



7

Statistical Tests

- Aim to identify a null hypothesis, collect some data, then estimate the probability of getting the observed data if the null hypothesis were true.

2 Types of Statistical Tests

1. Parametric Test
2. Nonparametric Test



Statistical Tests

Parametric Test

- Assumes data is from a **normal** distribution or is **approximately** a normal distribution.
- Uses the mean value for central tendency
- Use if type of data is interval or ratio
- With the help of the *central limit theorem*, nonnormal distributed data can still use parametric tests given that it meets the sample size requirements:

Parametric analyses

Sample size requirements for nonnormal data

1-sample t-test

Greater than 20

2-sample t-test

Each group should have more than 15 observations

One-Way ANOVA

- For 2-9 groups, each group should have more than 15 observations
- For 10-12 groups, each group should have more than 20 observations

Statistical Tests

Nonparametric Test

- Does not assume anything about the underlying distribution.
- Uses the **median** value for central tendency
- Use this if your data is **nominal** or **ordinal**.
- **Small** number of samples



Statistical Tests

Parametric Vs. Nonparametric Test

Which one to use?

If normally or approximately normally distributed, use parametric.



| Statistical Tests | | | | | | | |
|-----------------------------|------------------------|--------------------------|-----------------------------------|------------------------------|----------------------|------------------------|----------------|
| Type of test | Level of measurement | Sample characteristics | | | | | Correlation |
| | | One sample | Two sample | | K samples (i.e., >2) | | |
| | | | Independent | Dependent | Independent | Dependent | |
| Parametric | Interval or ratio | Z-test or <i>t</i> -test | Independent sample <i>t</i> -test | Paired sample <i>t</i> -test | One-way ANOVA | Repeated measure ANOVA | Pearson's test |
| Nonparametric | Categorical or nominal | Chi-square test | Chi-square test | Mc-Nemar test | Chi-square test | Cochran's Q | |
| | Rank or ordinal | Chi-square test | Mann-Whitney U-test | Wilcoxon signed rank test | Kruskal-Wallis | Friedman's ANOVA | Spearman's rho |
| ANOVA: Analysis of variance | | | | | | | |


| Type of test | Level of measurement | Sample characteristics | | | | | Correlation |
|---------------|------------------------|--------------------------|-----------------------------------|------------------------------|----------------------|------------------------|----------------|
| | | One sample | Two sample | | K samples (i.e., >2) | | |
| | | | Independent | Dependent | Independent | Dependent | |
| Parametric | Interval or ratio | Z-test or <i>t</i> -test | Independent sample <i>t</i> -test | Paired sample <i>t</i> -test | One-way ANOVA | Repeated measure ANOVA | Pearson's test |
| Nonparametric | Categorical or nominal | Chi-square test | Chi-square test | Mc-Nemar test | Chi-square test | Cochran's Q | |
| | Rank or ordinal | Chi-square test | Mann-Whitney U-test | Wilcoxon signed rank test | Kruskal-Wallis | Friedman's ANOVA | Spearman's rho |

ANOVA: Analysis of variance

| Type of test | Level of measurement | Sample characteristics | | | | | Correlation |
|---------------|------------------------|--------------------------|-----------------------------------|------------------------------|----------------------|------------------------|----------------|
| | | One sample | Two sample | | K samples (i.e., >2) | | |
| | | | Independent | Dependent | Independent | Dependent | |
| Parametric | Interval or ratio | Z-test or <i>t</i> -test | Independent sample <i>t</i> -test | Paired sample <i>t</i> -test | One-way ANOVA | Repeated measure ANOVA | Pearson's test |
| Nonparametric | Categorical or nominal | Chi-square test | Chi-square test | Mc-Nemar test | Chi-square test | Cochran's Q | |
| | Rank or ordinal | Chi-square test | Mann-Whitney U-test | Wilcoxon signed rank test | Kruskal-Wallis | Friedman's ANOVA | Spearman's rho |

ANOVA: Analysis of variance

Common Statistical Tests



Common Statistical Tests

Independent Samples T-Test

- compares the means **between two unrelated groups** on the same continuous, dependent variable.

Hypotheses:

$H_0 : \mu_1 = \mu_2$ ("the two population means are equal")

$H_a : \mu_1 \neq \mu_2$ ("the two population means are not equal")

OR

$H_0 : \mu_1 - \mu_2 = 0$ ("the difference between the two population means is equal to 0")

$H_a : \mu_1 - \mu_2 \neq 0$ ("the difference between the two population means is not 0")

where μ_1 and μ_2 are the population means for group 1 and group 2, respectively.

Example:

Heart rate of 10 people before drinking a cup of energy drink and then measure the heart rate of **some other group of people** who have drank energy drinks.

Common Statistical Tests

Paired Sample T-test

- compares two means that are from the same individual, object, or related units.
- two means:
 - represent two different times(e.g pre-test and post-test with an intervention between the two time points)
 - two different but related conditions or units

Example:

Heart rate of 10 people before drinking the energy drink and then measure the heart rate of the same 10 people after drinking the energy drink.

Common Statistical Tests

Chi-square Test of Independence

- used to determine whether two categorical variables from one single population are related or not.

Is there a significant association between the two variables?

Hypotheses:

H_0 : The two variables are independent

H_a : The two variables are not independent.

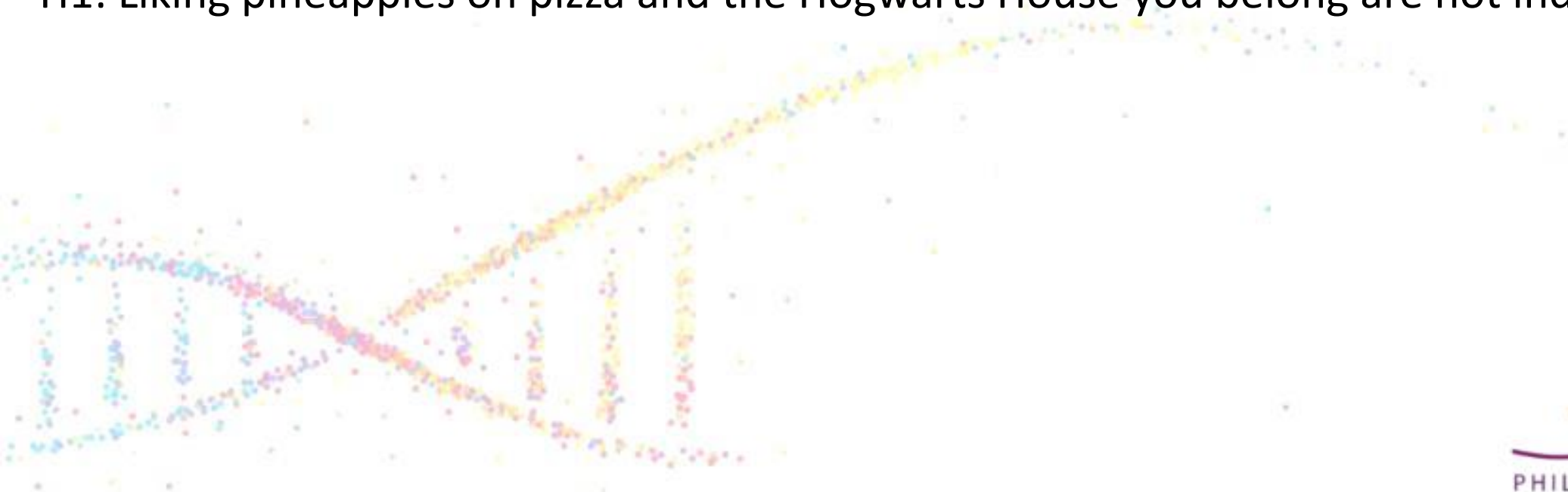
Common Statistical Tests

Chi-square Test of Independence Example

Does liking pineapples on pizza affect the probabilities of you identifying which of the Hogwarts Houses you belong?

H0: Liking pineapples on pizza and the Hogwarts House you belong are independent.

H1: Liking pineapples on pizza and the Hogwarts House you belong are not independent.



Which statistical test should be used?

That depends on the:

- a. Type of data you have
Continuous or Discrete
- b. Type of question you want to answer:
Relationships or Differences.
- c. Number of treatment groups
- d. Parametric or Nonparametric

Use [this](#) flowchart as your guide.

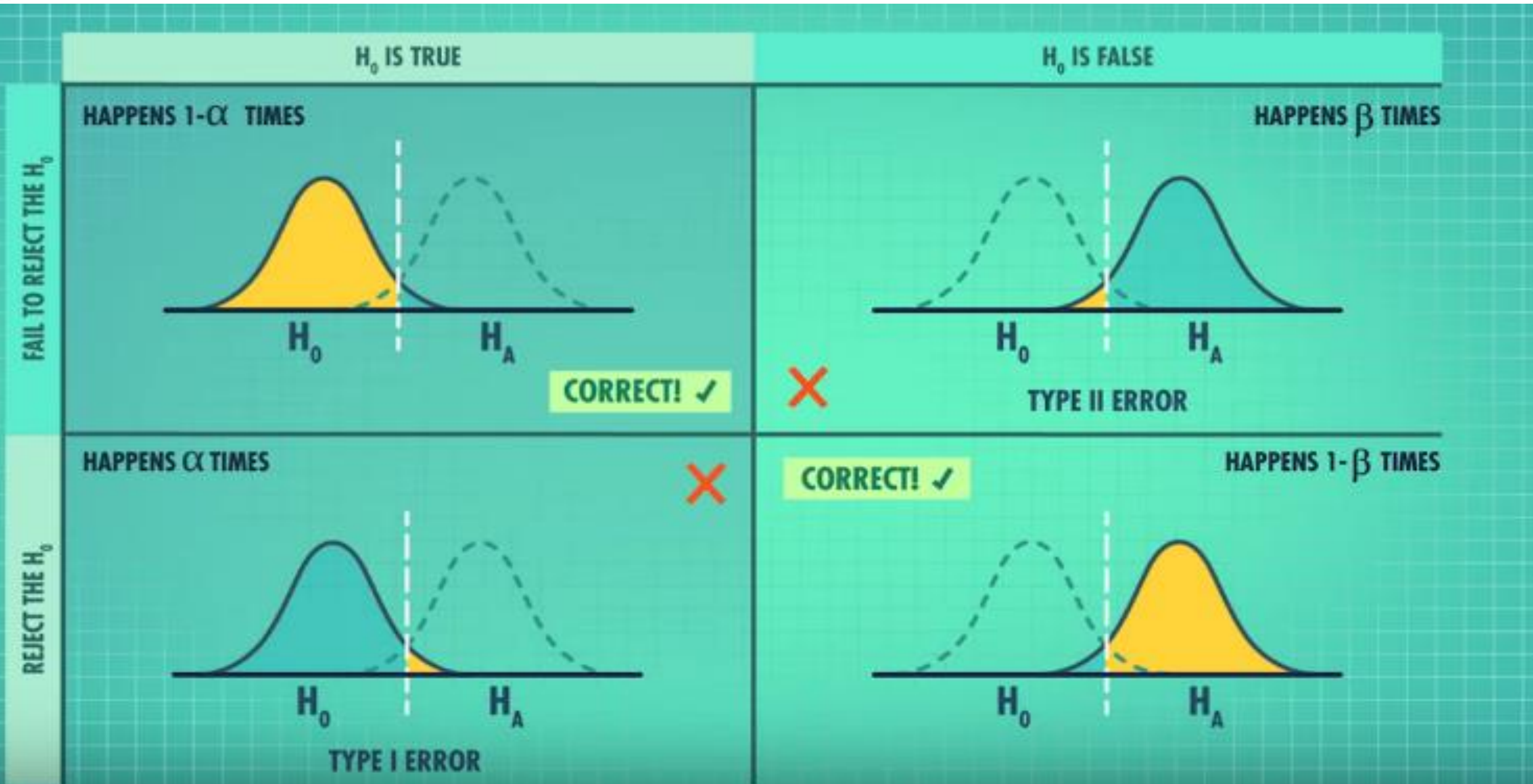
Type I and Type II Errors

Type I and II Errors

| | | Actual Situation “Truth” | |
|----------|---------------------|--|--|
| | | H_0 True | H_0 False |
| Decision | Do Not Reject H_0 | Correct Decision $1 - \alpha$ | Incorrect Decision Type II Error β |
| | Reject H_0 | Incorrect Decision Type I Error α | Correct Decision $1 - \beta$ |

$$\alpha = P(\text{Type I Error}) \quad \beta = P(\text{Type II Error})$$

Concept of Statistical Power



Concept of Statistical Power

Power of a hypothesis test is the **probability** that the test **correctly rejects the null hypothesis**.

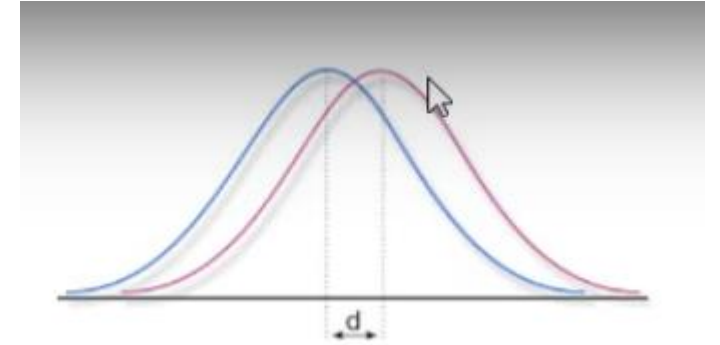
$$\text{power} = P(\text{reject } H_0 \mid H_0 \text{ is false})$$

- Used to avoid Type II error

Factors affecting Statistical Power

Effect Size

- how strong is the difference between two distributions.



Higher effect size = higher statistical power

Lower effect size = lower statistical power

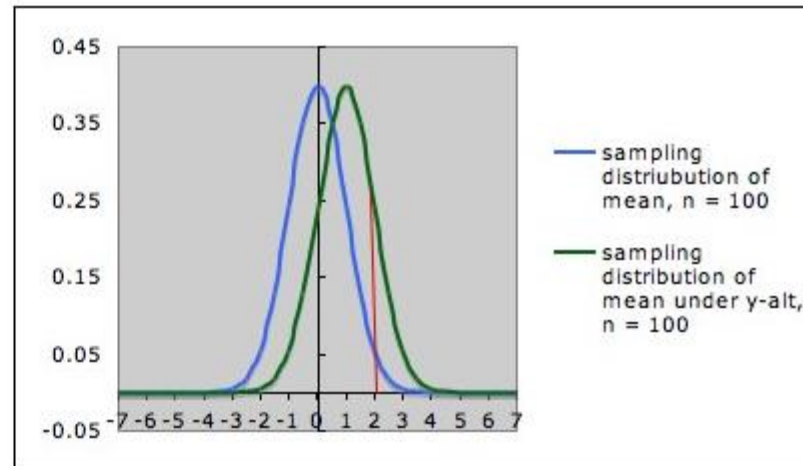
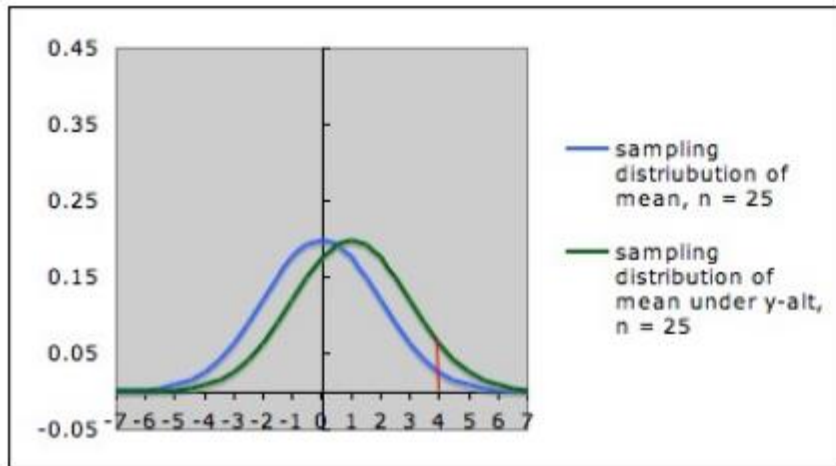
Factors affecting Statistical Power

Sample Size

- number of samples

Higher sample size = higher statistical power

Lower sample size = lower statistical power



Factors affecting Statistical Power

Significance Level(α)

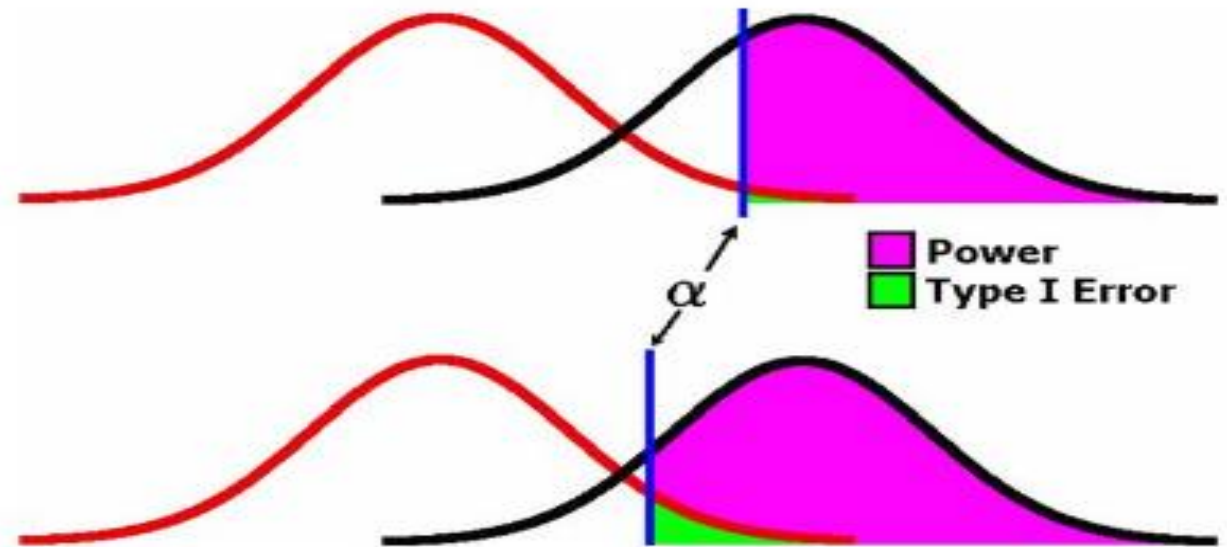
- how unlikely a positive result must be for the H_0 to be rejected.

Higher α = higher statistical power

Lower α = lower statistical power

NOTE:

If you increase α , you increase the chance of type 1 error.



Concept of Statistical Power

Experiments cost money, so if you're going to go through the process of growing cells in a petri dish or testing a new software tool you want to be relatively confident you'll be able to detect an effect if there is one.



Concept of Statistical Power

Across many fields, a statistical power of 80% or more is considered sufficient.

When researchers design studies, they'll decide how many subjects they need based on estimates of effect size and power.

Statistical Power tells us our chance of detecting an effect if there is one.

Measures of Test Accuracy

1. **Accuracy** - How close a value is to its true value.
2. **Precision** - How repeatable a measurement is.
3. **Specificity** - Describes how well the test is detecting non-diseased individuals as truly not having the disease.
4. **Sensitivity** - Describes how well the test detects disease in all who truly have the disease.

P = non-diseased

N = disease

Multiple Testing Correction

Multiple Testing - refers to any instance that involves the simultaneous *testing* of more than one hypothesis.

Fishing Expedition

A group of clinical investigators conduct a research study that collects hundreds or maybe even thousands of variables on a sample of patients. With a rich dataset in hand, they set out indiscriminately examining variables and conducting hypothesis tests with the expectation that this endeavor may yield a bountiful harvest of meaningful scientific results. However, this kind of fishing expedition can have unexpected complications.

Multiple Testing Correction

Why Multiple Testing Matters

In general, if we perform m hypothesis tests, what is the probability of at least 1 false positive?

$$P(\text{Making an error}) = \alpha$$

$$P(\text{Not making an error}) = 1 - \alpha$$

$$P(\text{Not making an error in } m \text{ tests}) = (1 - \alpha)^m$$

$$P(\text{Making at least 1 error in } m \text{ tests}) = 1 - (1 - \alpha)^m$$

Multiple Testing Correction

Type I Error Rate Control

Family-wise Error Rate (FWER)

- Control the probability that there is *a single type I error* in the entire set (family) of hypotheses tested.

Bonferroni Correction

- single step procedure makes equal adjustments to each p-value.

Multiple Testing Correction

Bonferroni Correction

1. Divide the alpha level by the number of tests you're running and apply that alpha level to each individual test.

Example:

For example, if your overall alpha level is .05 and you are running 5 tests, then each test will have an alpha level of $.05/5 = .01$

2. Apply the new alpha level to each test for finding p-values. In this example, the p-value would have to be .01 or less for statistical significance.

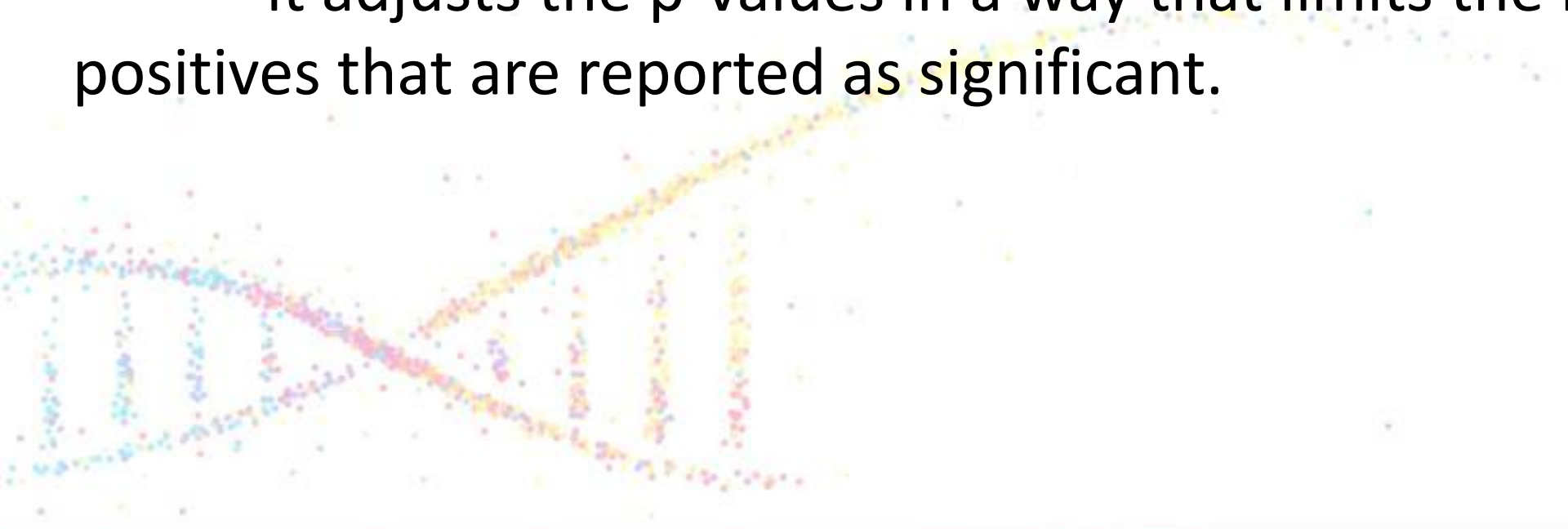
Multiple Testing Correction

False Discovery Rate(FDR)

- designed to control the proportion of false positives among the set of rejected hypotheses.

Benjamini Hochberg Correction

- it adjusts the p-values in a way that limits the number of false positives that are reported as significant.



Multiple Testing Correction

Benjamini Hochberg Correction

- it adjusts the p-values in a way that limits the number of false positives that are reported as significant.

Steps:

1. Put the individual p-values in ascending order.
2. Assign ranks to the p-values. For example, the smallest has a rank of 1, the second smallest has a rank of 2.
3. Pick a desired FDR level, q (e.g 5%)
4. Starting from the top of the list, Calculate each individual p-value's Benjamini-Hochberg critical value, using the formula $(i/m)q$, where:
 - i = the individual p-value's rank,
 - m = total number of tests,
 - q = the false discovery rate (a percentage, chosen by you).
5. Accept all genes with $p \leq (i/m)q$

B&H FDR Example

Controlling the FDR at $\delta = 0.05$

| Rank (j) | P-value | $(j/m) \times \delta$ | Reject H_0 ? |
|----------|---------|-----------------------|----------------|
| 1 | 0.0008 | 0.005 | 1 |
| 2 | 0.009 | 0.010 | 1 |
| 3 | 0.165 | 0.015 | 0 |
| 4 | 0.205 | 0.020 | 0 |
| 5 | 0.396 | 0.025 | 0 |
| 6 | 0.450 | 0.030 | 0 |
| 7 | 0.641 | 0.035 | 0 |
| 8 | 0.781 | 0.040 | 0 |
| 9 | 0.900 | 0.045 | 0 |
| 10 | 0.993 | 0.050 | 0 |

Multiple Testing Correction

Family-wise Error Rate – control probability of single type 1 error.

Bonferroni Correction

False Discovery Rate – control the probability of false positives among the rejected hypotheses.

Benjamini Hochberg Correction

