

**BCBB Workshop**

# Statistical Testing

Qinlu (Claire) Wang  
Statistician

Bioinformatics and Computational Biosciences Branch (BCBB)  
Office of Cyber Infrastructure and Computational Biology (OCICB)  
National Institute of Allergy and Infectious Diseases (NIAID)



NIAID

## OCICB Bioinformatics and Computational Biosciences Branch (BCBB)

- Part of NIAID
- Group of ~50
- Software developers
- Computational Biologists
- Project Management & Analysis Professionals
- Biostatistics, Phylogenetics, Genomics, Structural Biology, Proteomics, Programming



NIAID

## How to contact us?

### 1. Submit Request

*Send emails to [bioinformatics@niaid.nih.gov](mailto:bioinformatics@niaid.nih.gov) OR fill a request online at [Online Request](#)*

### 2. Attend our workshops

*[http://www.eventzilla.net/user/NIAID\\_OCICB\\_BCBB](http://www.eventzilla.net/user/NIAID_OCICB_BCBB)*

### 3. Tell us what statistical topics you want to learn about

*[BCBB Statistical Training – Suggest a class!](#)*

### 4. Join our Slack group!

*The invitation will be sent out after this workshop*



NIAID

## Learning materials

### ■ Online Statistics Education: An Interactive Multimedia Course of Study

*<http://onlinestatbook.com/2/index.html>*

### ■ Statistics Glossary v1.1

*This site is a glossary of statistical terms searchable by topic or in alphabetical order.*

*Created by Valerie J. Easton & John H. McColl*

*<http://www.stats.gla.ac.uk/steps/glossary/>*



NIAID

www.nature.com/nature/authors/gta

# nature

## Statistical checklist

Authors providing statistical analysis in their papers are asked to check the following points to ensure statistical adequacy.

It is helpful to the editors and the referees if authors confirm in their submission and resubmission cover letters that they have applied all relevant checks from the list below to the work described in their paper.

**IN THE METHODS SECTION:**

- Type and applicability of test used**
  - Comparisons of interest are clearly defined
  - Name of tests applied are clearly stated
  - All statistical methods identified unambiguously
  - Justification for use of test is given
  - Data meet all assumptions of tests applied (with particular attention paid to non-normal data sets or small sample sizes, which should be identified in the text as such)
  - Adjustments made for multiple testing is explained
- Details about the test**
  - $n$  is reported at the start of the study and for each analysis thereafter
  - Sample size calculation (or justification) is given
  - Unit of analysis is given for all comparisons
  - Alpha level is given for all statistical tests
  - Tests are clearly identified as one or two-tailed
  - Randomization procedures or other ways to eliminate bias in sampling (in particular for experiments involving animals) described
  - Actual  $P$  values are given for primary analyses
- Descriptive statistics summary**
  - $n$  for each data set is clearly stated
  - A clearly labelled measure of centre (e.g. mean or median) is given
  - A clearly labelled measure of variability (e.g. standard deviation or range) is given
  - All numbers following a  $\pm$  sign are identified as standard errors (s.e.m.) or standard deviations (s.d.)

**Anomalies**

- Any unusual or complex statistical methods are clearly defined and explained for Nature's wide readership. (Authors are encouraged to use Supplementary Information for long explanations.)
- Any data exclusions are stated and explained
- Any discrepancies in the value of  $n$  between analyses are clearly explained and justified
- Any method of treatment assignment (randomization, etc) is explained and justified
- Any data transformations are clearly described and justified

**WITHIN INDIVIDUAL GRAPHS:**

**Distortions**

- Any distorted effect sizes (e.g. by truncation of  $y$  axis) are clearly labelled and justified

**Clear labelling**

- Error bars are present on all graphs, where applicable.
- All error bars are clearly labelled

**PLEASE HELP US**

Many statistical analyses published in *Nature* are highly sophisticated and outside the scope of this checklist, particularly in the case of some studies in physical sciences disciplines. Authors and referees who have specific suggestions for additional entries to this list are encouraged to send them by e-mail to [authors@nature.com](mailto:authors@nature.com) or [referees@nature.com](mailto:referees@nature.com). *Nature* will update this checklist at intervals in an effort to ensure that papers published are statistically robust.

**NIAID**

**Outline**

## 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

## 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

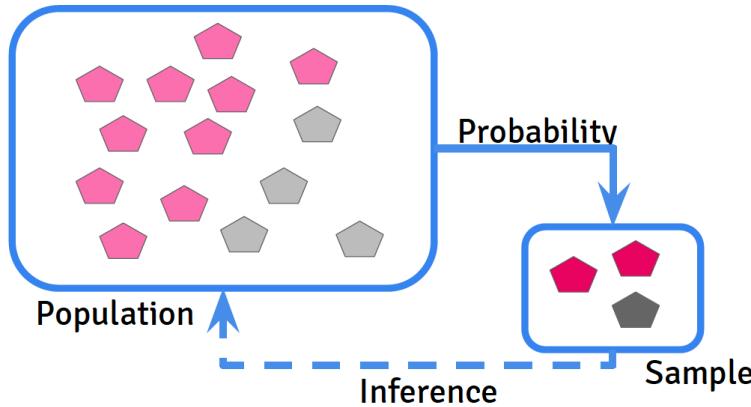
## 3. Application in Prism and R

**NIAID**

**NIH** National Institute of Allergy and Infectious Diseases

# 1. Statistical Testing Process

The best way to determine whether a statistical hypothesis is true would be to examine the entire population. Since that is often impractical, researchers typically examine a random sample from the population.



NIAID

## 1.1. Formulate Null and Alternative hypotheses

- E.g. (null)  $H_0: \mu_1 = \mu_2$  vs. (alternative)  $H_A: \mu_1 \neq \mu_2$

## 1.2. Choose the appropriate test statistics

- E.g. Student's *t*-test, Wilcoxon test, ...

## 1.3. Calculate Sample Size

- Significance Level, Power, Effect Size

## 1.4. Compute the probability of observing the test statistic under the null hypothesis

- i.e. Compute *p*-value and confidence interval

## 1.5. Make a statistical decision and a biological conclusion

- "Reject the null hypothesis" or "fail to reject the null hypothesis"
- E.g. New drug reduces viral load, vitamin C helps prevent cancer



NIAID

# Outline

## 1. Statistical Testing Process

### 1.1. Formulate Null and Alternative hypotheses

- 1.2. Calculate the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

## 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

## 3. Application in Prism and R



NIAID

## 1.1. Formulate Null And Alternative Hypotheses

- Two Types of Statistical Hypotheses
- One-Sided vs. Two-sided Tests



NIAID

## ➤ Two Types of Statistical Hypotheses

**Null hypothesis:** states that a population parameter is equal to a hypothesized value.

**Alternative hypothesis:** states that a population parameter is smaller, greater, or different than the hypothesized value in the null hypothesis.

The diagram shows two panels. The left panel, labeled  $H_1$ , shows a watering can with a green 'X' spraying water onto a small plant. A dashed box highlights a taller, more robust plant with a question mark above it. Below this is the text: "H<sub>1</sub>: Application of bio-fertilizer 'x' increase plant growth." The right panel, labeled  $H_0$ , shows a watering can with a green 'X' spraying water onto a small plant. A dashed box highlights a similar small plant with a red question mark above it. Below this is the text: "H<sub>0</sub>: Application of bio-fertilizer 'x' do not increase plant growth." At the bottom, there are two green boxes: "Alternative hypothesis" with the note "✓ The alternative hypothesis is a hypothesis which the researcher tries to prove." and "Null hypothesis" with the note "✓ The null hypothesis is a hypothesis which the researcher tries to disprove, or nullify."

NIH National Institute of Allergy and Infectious Diseases
NIAID

## ➤ One-Sided vs. Two-sided Tests

### Two-sided Test

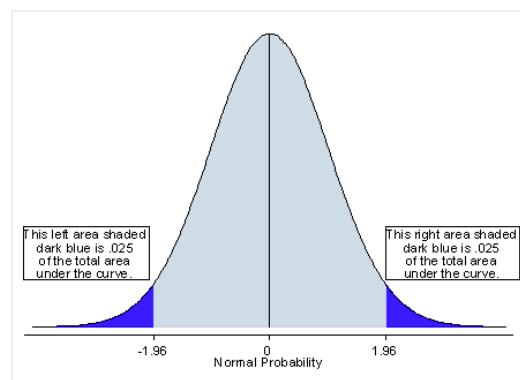
If you are using a significance level of  $\alpha$ , a two-sided test allots half of your  $\alpha$  to test the statistical significance in one direction and half of your  $\alpha$  to test statistical significance in the other direction.

For example:

$$H_0: \mu = k$$

$$H_a: \mu \neq k.$$

A two-sided test will test both if the mean is significantly greater than  $k$  **AND** if the mean significantly less than  $k$ .



### One-sided Test

If you are using a significance level of  $\alpha$ , an one-sided test allots all of your  $\alpha$  to test the statistical significance in one direction of interest.

For Example:

$$H_0: \mu = k$$

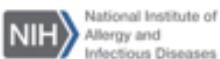
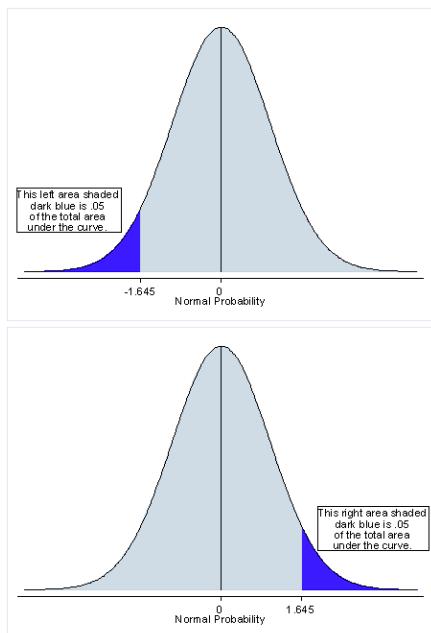
$$H_a: \mu > k$$

or

$$H_0: \mu = k$$

$$H_a: \mu < k.$$

The one-sided test provides more power to detect an effect in one direction by not testing the effect in the other direction.



NIAID

### **When two-sided test?**

When you have no idea of the direction or it is unacceptable to ignore the effect in the untested direction

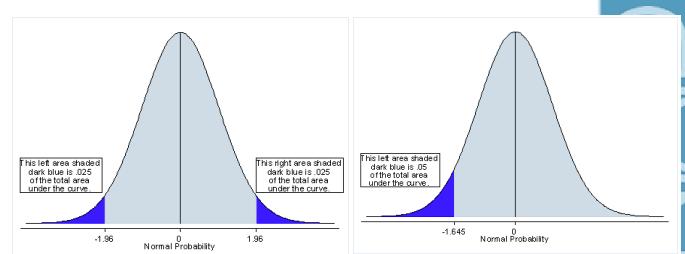
### **When one-sided test?**

If you consider the consequences of missing an effect in the untested direction are negligible and in no way irresponsible or unethical. For example, developing a new and cheaper drug.

#### **Attention!**

- Choosing a one-tailed test for the sole purpose of attaining significance is not appropriate.
- Choosing a one-tailed test after running a two-tailed test that failed to reject the null hypothesis is not appropriate

	PROS	CONS
One-Tailed Tests	<ul style="list-style-type: none"> <li>Requires less traffic</li> <li>Gains significance faster (read: <a href="#">why significance does not equal validity</a>)</li> </ul>	<ul style="list-style-type: none"> <li>Only accounts for one scenario</li> <li>Can lead to inaccurate and biased results</li> </ul>
Two-Tailed Tests	<ul style="list-style-type: none"> <li>Accounts for all three scenarios</li> <li>Leads to accurate and reliable results</li> </ul>	<ul style="list-style-type: none"> <li>Requires more traffic</li> <li>Takes longer to gain significance</li> </ul>



# Outline

## 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics**
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

## 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

## 3. Application in Prism and R



NIAID

## 1.2. Statistical Tests

- What is Statistical Test
- Variance, Standard Deviation and Spread
- Common Statistical Tests



NIAID

## ➤ What is statistical test

$$\text{Test} = \frac{\text{Difference}}{\text{Error}} = \frac{\text{Statistic} - \text{Null Value}}{\text{Error}}$$

- Almost all tests used in inferential statistics can be generalized as the ratio of a “difference” over an “error”
  - Difference between a statistic and null value (usually 0)
  - A **statistic** is nothing more than a numeric summary of the experimental data with respect to the null hypothesis
  - A **null value** is an assumption about the population under the null hypothesis
  - An **error** is an estimate of the sampling distribution error



NIAID

## Example: Two-sample Student's T-test

$$T^* = \frac{\overline{X}_1 - \overline{X}_2 - 0}{\sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2}\right) \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}}}$$

statistic →  $\overline{X}_1 - \overline{X}_2$  ← null value  
 ↓  
 standard error

- The “statistic” in a two-sample *t*-test is a difference between the two sample means and the null value is zero
  - The hypothesis  $\mu_1 = \mu_2$  implies  $\mu_1 - \mu_2 = 0$
- The standard error is an estimate of the common variance



NIAID

**Variance, Standard Deviation and Spread**

The standard deviation of the mean (SD) is the most commonly used measure of the spread of values in a distribution. SD is calculated as the square root of the variance (the average squared deviation from the mean).

Variance in a population is:

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

[x is a value from the population,  $\mu$  is the mean of all x, n is the number of x in the population,  $\Sigma$  is the summation]

Variance is usually estimated from a sample drawn from a population. The unbiased estimate of population variance calculated from a sample is:

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

[ $x_i$  is the ith observation from a sample of the population,  $\bar{x}$  is the sample mean, n (sample size) - 1 is [degrees of freedom](#),  $\Sigma$  is the summation]

The spread of a distribution is also referred to as dispersion and variability. All three terms mean the extent to which values in a distribution differ from one another.

SD is the best measure of spread of an approximately normal distribution. This is not the case when there are extreme values in a distribution or when the distribution is skewed, in these situations interquartile range or semi-interquartile are preferred measures of spread. Interquartile range is the difference between the 25th and 75th centiles. Semi-interquartile range is half of the difference between the 25th and 75th centiles. For any symmetrical (not skewed) distribution, half of its values will lie one semi-interquartile range either side of the median, i.e. in the interquartile range. When distributions are approximately normal, SD is a better measure of spread because it is less susceptible to sampling fluctuation than (semi-)interquartile range.

If a variable y is a linear ( $y = a + bx$ ) transformation of x then the variance of y is  $b^2$  times the variance of x and the standard deviation of y is b times the variance of x.

The standard error of the mean is the expected value of the standard deviation of means of several samples, this is estimated from a single sample as:

$$SEM = \frac{s}{\sqrt{n}}$$

[s is standard deviation of the sample mean, n is the sample size]

**Flow Chart for Selecting Commonly Used Statistical Tests**

**How to choose appropriate statistical test?**

**NIAID**

**NIH** National Institute of Allergy and Infectious Diseases

## Parametric vs. Non-parametric Tests?

In practice:

- Parametric statistical procedures rely on assumptions about the shape of the distribution in the underlying population and about the form or parameters of the assumed distribution.
- Nonparametric statistical procedures rely on **NO** or **FEW** assumptions about the shape or parameters of the population distribution from which the sample was drawn.



NIAID

## Why don't we always use nonparametric tests?

### Drawbacks

- Nonparametric tests generally are less statistically powerful than parametric procedures when the data truly are approximately normal.
- The results of nonparametric tests are often harder to interpret than the results of parametric tests because many nonparametric tests use rankings of the values in the data rather than using the actual data.

--> In short, nonparametric procedures are useful in many cases and necessary in some, but they are not always a perfect solution.



NIAID

Don't use this approach:

~~"First perform a normality test. If the P value is low, demonstrating that the data do not follow a Gaussian distribution, choose a nonparametric test. Otherwise choose a conventional test."~~

Prism does not use this approach, because the choice of parametric vs. nonparametric is more complicated than that.

- Often, the analysis will be one of a series of experiments, you cannot rely on the results from a single normality test.
- Many biological variables follow [lognormal distributions](#). Check lognormality instead.
- Data can fail a normality test because of the presence of an [outlier](#). Removing that outlier.
- The decision of whether to use a parametric or nonparametric test is most important with small data sets. But with small data sets, [normality tests](#) have little power to detect non-gaussian distributions, so an automatic approach would give you false confidence.
- With large data sets, normality tests can be too sensitive.



NIAID

## Outline

### 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size**
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

### 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

### 3. Application in Prism and R



NIAID

## 1.3. Sample Size Calculation

The following **four quantities** have an intimate relationship:

- *Sample Size*
- ***Effect Size***
- *Significance level =  $P(\text{Type I error})$  = probability of finding an effect that is not there*
- *Power =  $1 - P(\text{Type II error})$  = probability of finding an effect that is there*

Given any three, we can determine the fourth.



NIAID

### Type I and Type II Errors

- **Type I error ("false positive")**: The error of rejecting a TRUE null hypothesis.
- **Type II error ("false negative")**: The error of not rejecting a FALSE null hypothesis

For example, when the null hypothesis is "You're not pregnant", and the alternative hypothesis is "You are pregnant":



NIAID

The chances of committing these two types of errors are inversely proportional

- To decrease the chance of committing Type I error: make your alpha ( $p$ ) value more stringent.
- To decrease the chance of committing Type II error: increase your analyses' power by either increasing your sample size or relaxing your alpha level

	H <sub>0</sub> rejected	Fail to reject H <sub>0</sub>
H <sub>0</sub> false	Correct	Type II error
H <sub>0</sub> true	Type I error	correct

$$\text{Alpha } (\alpha) = \text{Prob}(\text{Type I error})$$

$$\text{Beta } (\beta) = \text{Prob}(\text{Type II error})$$

$$\text{Power} = 1 - \beta$$



NIAID

## Effect Size

Effect size is a statistical concept that measures the strength of the relationship between two variables on a numeric scale.

- Statistic effect size helps us in determining if the difference is real or if it is due to a change of factors.
- In statistics analysis, the effect size is usually measured in three ways:
  - (1) standardized mean difference
  - (2) odd ratio
  - (3) correlation coefficient.



NIAID

## Sample Size Calculation

The following **four quantities** have an intimate relationship:

- *Sample Size*
- ***Effect Size***
- *Significance level =  $P(\text{Type I error})$  = probability of finding an effect that is not there*
- *Power =  $1 - P(\text{Type II error})$  = probability of finding an effect that is there*

Given any three, we can determine the fourth.



NIAID

## Outline

### 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis**
- 1.5. Make statistical decision and biological conclusion

### 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

### 3. Application in Prism and R



NIAID

## 1.4. Probability of Observing the Test Statistic under the Null Hypothesis

### ➤ P-value

- *What is p-value*
- *How to calculate p-value*
- *How to interpret p-value*

### ➤ Confidence Interval

- *What is confidence interval*
- *How to calculate confidence interval*
- *Z-intervals vs. T-intervals*
- *How to interpret confidence interval*



NIAID

### P-value

- ✓ P-value is a measure of strength of the evidence against null hypothesis.
- ✓ P-value is the probability of getting the observed value of test statistic, or value with even greater evidence against null hypothesis ( $H_0$ ), if the null hypothesis of a study question is true. For example, tossing a coin 10 times.
- ✗ P-value is NOT the probability of the means in the two samples being equal



NIAID

## How to calculate p-value?

1. Obtained according to the asymptotic distribution of the test statistic

- When the sample size is sufficiently large.
- E.g: Student's t-test

2. Calculate p-value empirically.

- E.g: Fisher's exact test

3. Resampling-based procedures

- When the asymptotic distribution is either unreliable due to insufficient sample size or unavailable for complicated test statistics
- E.g: permutation or bootstrap

[Comparison between resampling methods](#)



NIAID

### Example:

Suppose a pharmaceutical company manufactures ibuprofen pills. They need to perform some quality assurance to ensure they have the correct dosage, which is supposed to be 500 milligrams. This is a two-sided test because if the company's pills are deviating significantly in either direction, meaning there are more than 500 milligrams or less than 500 milligrams, this will indicate a problem.

$$H_0: \mu = 500 \text{ mg}$$

$$H_A: \mu \neq 500 \text{ mg}$$

In a random sample of 125 pills, there is an average dose of 499.3 milligrams. The standard deviation of the population is known as 6 milligrams. Because this is quantitative data, 500 mg is the population mean. We can use the following formula to calculate the z-score:

$$z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$\bar{X}$  = sample mean = 499.3 mg

$\mu$  = population mean = 500 mg

$\sigma$  = population standard deviation = 6 mg

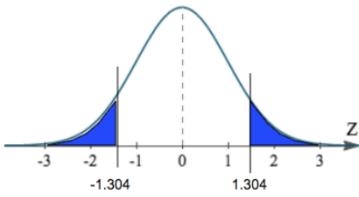
$n$  = sample size = 125

$$z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{499.3 - 500}{\frac{6}{\sqrt{125}}} = \frac{-0.7}{\frac{6}{\sqrt{125}}} = \frac{-0.7}{\frac{6}{11.18}} = \frac{-0.7}{0.53667} = -1.304$$



NIAID

We get a z-score of negative 1.304. Because this is a two-sided test, it is not enough to just look at the left tail. We also have to look at the equivalent of the right tail, or a positive 1.304.



The first way to find the p-value is to use the z-table. In the z-table, the left column will show values to the tenths place, while the top row will show values to the hundredths place. If we have a z-score of -1.304, we need to round this to the hundredths place, or -1.30. In the left column, we will first find the tenths place, or -1.3. In the top row, we will find the hundredths place, or 0.

Z-score	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
-4.4	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003
-4.3	0.0005	0.0005	0.0005	0.0004	0.0004	0.0004	0.0004	0.0004	0.0004	0.0003
-4.2	0.0007	0.0007	0.0006	0.0006	0.0006	0.0006	0.0006	0.0005	0.0005	0.0005
-4.1	0.0010	0.0009	0.0009	0.0009	0.0008	0.0008	0.0008	0.0007	0.0007	0.0007
-4.0	0.0013	0.0013	0.0013	0.0012	0.0012	0.0011	0.0011	0.0011	0.0011	0.0010
-3.9	0.0019	0.0018	0.0018	0.0017	0.0016	0.0016	0.0015	0.0014	0.0014	0.0014
-3.8	0.0026	0.0025	0.0025	0.0024	0.0023	0.0023	0.0022	0.0021	0.0021	0.0020
-3.7	0.0034	0.0033	0.0032	0.0031	0.0030	0.0029	0.0028	0.0027	0.0026	0.0025
-3.6	0.0045	0.0044	0.0043	0.0043	0.0040	0.0039	0.0038	0.0037	0.0036	0.0036
-3.5	0.0062	0.0060	0.0059	0.0057	0.0055	0.0054	0.0052	0.0051	0.0049	0.0048
-3.4	0.0082	0.0078	0.0075	0.0073	0.0071	0.0069	0.0068	0.0066	0.0066	0.0064
-3.3	0.0107	0.0104	0.0102	0.0099	0.0097	0.0094	0.0091	0.0089	0.0087	0.0084
-3.2	0.0139	0.0136	0.0132	0.0129	0.0126	0.0122	0.0119	0.0113	0.0110	0.0110
-3.1	0.0179	0.0174	0.0170	0.0166	0.0162	0.0158	0.0154	0.0150	0.0146	0.0143
-3.0	0.0228	0.0222	0.0217	0.0212	0.0207	0.0202	0.0197	0.0192	0.0188	0.0183
-2.9	0.0287	0.0281	0.0274	0.0268	0.0262	0.0256	0.0250	0.0244	0.0239	0.0233
-2.8	0.0356	0.0350	0.0343	0.0337	0.0330	0.0323	0.0316	0.0309	0.0303	0.0297
-2.7	0.0446	0.0436	0.0427	0.0418	0.0409	0.0401	0.0392	0.0384	0.0375	0.0367
-2.6	0.0548	0.0537	0.0536	0.0516	0.0505	0.0495	0.0485	0.0475	0.0465	0.0455
-2.5	0.0668	0.0655	0.0643	0.0630	0.0618	0.0606	0.0594	0.0582	0.0571	0.0559
-2.4	0.0808	0.0795	0.0778	0.0764	0.0749	0.0735	0.0721	0.0708	0.0694	0.0681
-2.3	0.0968	0.0951	0.0934	0.0918	0.0901	0.0885	0.0869	0.0853	0.0838	0.0823
-2.2	0.1151	0.1131	0.1112	0.1093	0.1075	0.1056	0.1038	0.1020	0.1003	0.0985
-2.1	0.1357	0.1335	0.1314	0.1292	0.1271	0.1251	0.1230	0.1190	0.1170	0.1170
-2.0	0.1587	0.1562	0.1539	0.1515	0.1492	0.1469	0.1446	0.1423	0.1401	0.1379
-1.9	0.1841	0.1814	0.1787	0.1759	0.1730	0.1701	0.1671	0.1640	0.1635	0.1611
-1.8	0.2119	0.2091	0.2061	0.2033	0.2005	0.1977	0.1949	0.1922	0.1892	0.1859
-1.7	0.2420	0.2389	0.2358	0.2327	0.2296	0.2266	0.2236	0.2206	0.2177	0.2148
-1.6	0.2743	0.2709	0.2676	0.2643	0.2611	0.2578	0.2546	0.2514	0.2483	0.2451
-1.5	0.3085	0.3050	0.3015	0.2981	0.2946	0.2912	0.2877	0.2843	0.2810	0.2776
-1.4	0.3446	0.3409	0.3372	0.3336	0.3300	0.3262	0.3228	0.3192	0.3156	0.3121
-1.3	0.3821	0.3783	0.3745	0.3707	0.3669	0.3632	0.3594	0.3557	0.3520	0.3483
-1.2	0.4207	0.4168	0.4129	0.4090	0.4052	0.4013	0.3974	0.3936	0.3897	0.3829
-1.1	0.4602	0.4562	0.4522	0.4483	0.4443	0.4404	0.4364	0.4325	0.4286	0.4247
0	0.5000	0.4960	0.4920	0.4880	0.4840	0.4801	0.4761	0.4721	0.4681	0.4641

This results in a p-value of 0.0968, or 9.68%, for a z-score of negative 1.304. We also need to take the positive 1.304 into account, which is the upper right tail. To calculate the true p-value, we just need to multiply 0.0968 by two, or 0.1936. This would be a p-value of 19.36%.

NIAID

NIH National Institute of Allergy and Infectious Diseases

## How to interpret p-value?

In the majority of analyses, an alpha of 0.05 is used as the cutoff for significance.

### If the p-value is less than 0.05

- ✓ we reject the null hypothesis that there's no difference between the means and conclude that a significant difference does exist.

### If the p-value is larger than 0.05

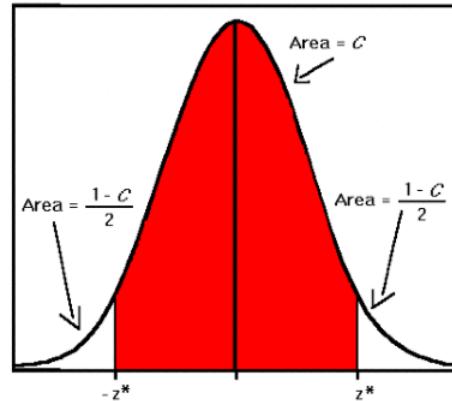
- ✓ we *cannot* conclude that a significant difference exists.
- ✗ Don't conclude that the statistically non-significant results showed "no association" / "no difference"

## Confidence Interval

*"A confidence interval gives an estimated range of values which is likely to include an unknown population parameter, the estimated range being calculated from a given set of sample data."*

----- Valerie J. Easton and John H. McColl's [Statistics Glossary v1.1](#)

Strictly speaking a 95% confidence interval means that if we were to take 100 different samples and compute a 95% confidence interval for each sample, then approximately 95 of the 100 confidence intervals will contain the true mean value ( $\mu$ ).



NIAID

## How to calculate confidence interval?

1. Obtained according to the asymptotic distribution of the test statistic
  - *When the sample size is sufficiently large.*
  
2. Resampling-based procedures
  - *When the asymptotic distribution is either unreliable due to insufficient sample size or unavailable for complicated test statistics*
  - *E.g: permutation or bootstrap*



NIAID

The **confidence interval estimate (CI)** is consist of:

- Point estimate, e.g., the sample mean
  - Investigator's desired level of confidence C (most commonly 95%)
  - Sampling variability or the standard error of the point estimate.
- 
- For a known standard deviation:  $\left( \bar{x} - z^* \frac{\sigma}{\sqrt{n}}, \bar{x} + z^* \frac{\sigma}{\sqrt{n}} \right)$
  - For an unknown standard deviation:  $\left( \bar{x} - t^* \frac{s}{\sqrt{n}}, \bar{x} + t^* \frac{s}{\sqrt{n}} \right)$



NIAID

## Z-Intervals

*Use a z-interval when:  
the sample size is greater than or equal to 30 and population standard deviation known OR Original population normal with the population standard deviation known.*

### Formula for the z-interval

If these conditions hold, we will use this formula for calculating the confidence interval:

$$\bar{x} \pm z_c \left( \frac{\sigma}{\sqrt{n}} \right)$$

where  $z_c$  is a critical value from the normal distribution (see below) and  $n$  is the sample size.

Common values of  $z_c$  are:

CONFIDENCE LEVEL	CRITICAL VALUE
90%	1.645
95%	1.96
99%	2.575



❖ [Understanding the t-distribution and its normal approximation](#)

NIAID

### Example using a z-interval

Suppose that in a sample of 50 college students in Illinois, the mean credit card debt was \$346. Suppose that we also have reason to believe (from previous studies) that the population standard deviation of credit card debts for this group is \$108. Use this information to calculate a 95% confidence interval for the mean credit card debt of all college students in Illinois.

#### Solution

Since we wish to estimate the mean, we immediately know we will be using either a t-interval or a z-interval. Looking a bit closer, we see that we have a large sample size ( $n = 50$ ) and we know the population standard deviation. Therefore, we will use a z-interval with  $z_c = 1.96$ . From reading the problem, we also have:

- Mean is \$346:  $\bar{x} = 346$
- Population standard deviation is 108:  $s = 108$

Applying the formula:

$$\bar{x} \pm z_c \left( \frac{\sigma}{\sqrt{n}} \right)$$

$$346 \pm 1.96 \left( \frac{108}{\sqrt{50}} \right)$$

The  $\pm$  indicates that we need to perform two different operations: a subtraction and an addition.

Left hand endpoint:

$$346 - 1.96 \left( \frac{108}{\sqrt{50}} \right) = 316.1$$

Right hand endpoint:

$$346 + 1.96 \left( \frac{108}{\sqrt{50}} \right) = 375.9$$

This gives our 95% confidence interval for  $\mu$ , the population mean, as  $(316.1, 375.9)$ .



NIAID

## T-Intervals:

### **Use a t-interval when:**

Population standard deviation UNKNOWN and original population normal OR sample size greater than or equal to 30 and Population standard deviation UNKNOWN.

### **Formula for the t-interval**

The formula for a t-interval is:

$$\bar{x} \pm t_c \left( \frac{s}{\sqrt{n}} \right)$$

where  $t_c$  is a critical value from the t-distribution,  $s$  is the sample standard deviation and  $n$  is the sample size.

### **Finding $t_c$**

The value of  $t_c$  depends on the sample size through the use of “degrees of freedom” where  $df = n - 1$ . We will use this to look up the value of  $t_c$  in a table (a nice [free version of that table can be found here](#), or typically in the back of your textbook if you are currently taking a class).



NIAID

### Example using a t-interval

Suppose that a sample of 38 employees at a large company were surveyed and asked how many hours a week they thought the company wasted on unnecessary meetings. The mean number of hours these employees stated was 12.4 with a standard deviation of 5.1. Calculate a 99% confidence interval to estimate the mean amount of time all employees at this company believe is wasted on unnecessary meetings each week.

#### Solution

As before, since we are estimating a mean with a confidence interval, we know it will either be a t-interval or a z-interval. In this case, we have a large sample ( $n = 38$ ), but we only have the sample standard deviation. If you aren't sure of that - read closely. The standard deviation of 5.1 was in the context of the sample, so  $s = 5.1$ . Thus, we will go ahead and use a t-interval since  $\sigma$  is unknown.

Before we can do that however, we need to look up the critical value. To know which row in the t-table to look at, we find the degrees of freedom which is  $n - 1 = 38 - 1 = 37$ . Using the table linked [here](#):

df (degrees of freedom)	Confidence Intervals, c					
	80%	90%	95%	98%	99%	99.9%
	Level of Significance for One-Tailed Test, $\alpha$					
0.100	0.050	0.025	0.010	0.005	0.001	
0.200	0.10	0.05	0.02	0.01	0.001	
36	2.000	1.989	2.028	2.051	2.719	3.582
37	1.990	1.987	2.026	2.431	2.717	3.574
38	1.984	1.986	2.024	2.429	2.712	3.566
39	1.980	1.985	2.023	2.426	2.708	3.558
40	1.977	1.984	2.021	2.423	2.704	3.551

Now that we have that, we plug the values into the formula and do the calculations to get our two endpoints. Remember that we have:

- Sample mean:  $\bar{x} = 12.4$
- Sample size:  $n = 38$
- Sample standard deviation:  $s = 5.1$
- Critical value:  $t_c = 2.715$



NIAID

Therefore the interval is:

$$\bar{x} \pm t_c \left( \frac{s}{\sqrt{n}} \right)$$

$$12.4 \pm 2.715 \left( \frac{5.1}{\sqrt{38}} \right)$$

This gives us the following two endpoints for our interval.

Left hand endpoint:

$$12.4 - 2.715 \left( \frac{5.1}{\sqrt{38}} \right) = 10.2$$

Right hand endpoint:

$$12.4 + 2.715 \left( \frac{5.1}{\sqrt{38}} \right) = 14.6$$

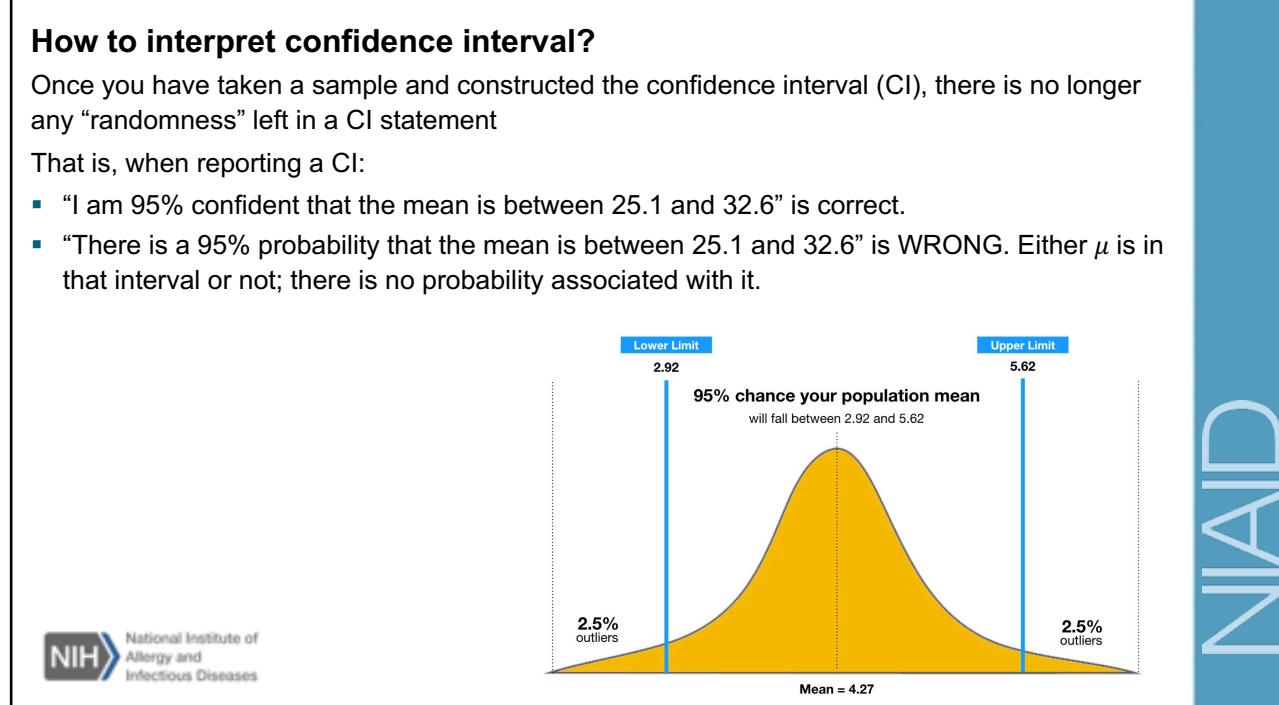
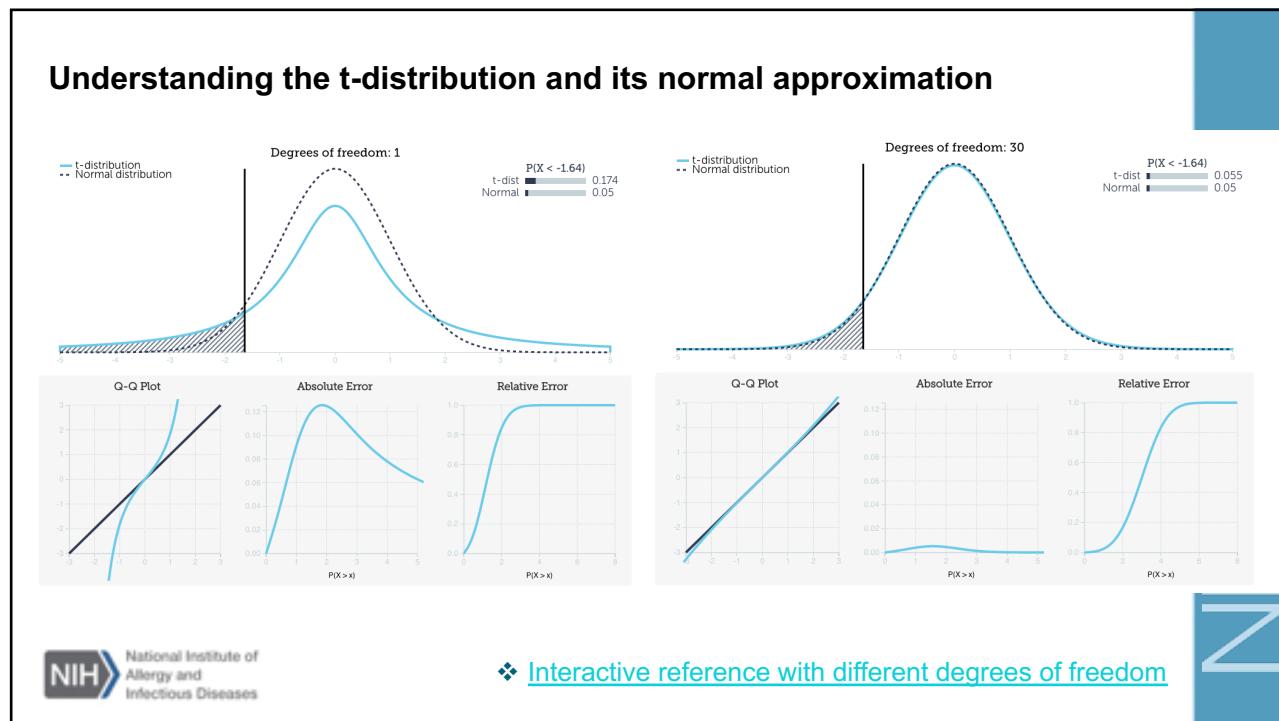
99% Confidence Interval for  $\mu$ :  $(10.2, 14.6)$

#### Interpretation

"We are 99% confident that the mean amount of time that all employees at this company think is wasted on meetings each week is between 10.2 and 14.6 hours."



NIAID



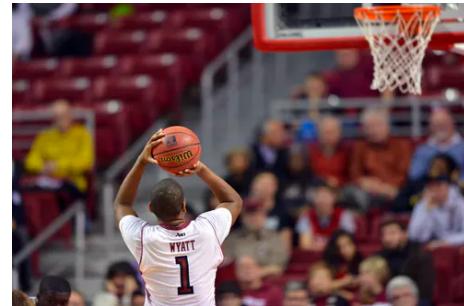
95% Confidence Interval simply means that if I take out samples a large number of times (tending to infinity) from a population with 0.05 significance level then 95% of those intervals will encompass the true parameter value.

There are two cases throwing a ball into a mug or throwing mug over a fixed lying ball. Most people think Confidence Interval as the first case but the things are other way round.

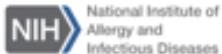
Randomness is not associated with parameter but the randomness is associated with interval. True parameters are fixed. But intervals estimated are random.

Roughly you can understand it like this. It is like, a ball (true parameter) is lying and you are throwing a mug(CI with 0.05 significance level) over it a large number of times. And asymptotically 95% of those throws will cover the ball.

So if I say in terms of frequentist then in future if you take out a sample with 95% CI then there is 0.95 probability that it will encompass true parameter value.



Ring  
Toss  
Game



❖ [More Details](#)

NIAID

## Outline

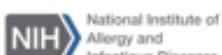
### 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion**

### 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

### 3. Application in Prism and R



NIAID

## 1.5. Statistical Decision and Biological Conclusions

- Statistical Decision
- Biological Conclusions
- How to report statistical results



NIAID

- A **statistical decision** is a choice to “reject the null hypothesis” or “fail to reject the null hypothesis”
  - The decision is based on a critical value or decision rule
  - E.g. Reject the null hypothesis if p-value < 0.05
- A **biological conclusion** is the final interpretation of the statistical testing process in plain language
  - E.g. Vitamin C prevents cancer, drug reduced viral loads, ...
  - Make sure conclusion can be justified by the hypotheses



NIAID

## How to report statistical results

### Overall:

- Every statistical paper should report all methods completely enough so someone else could reproduce the work exactly.
- Every figure and table should present the data clearly
- All the results should be reported completely enough that no one wonders what you actually did.

### Statistical Methods:

- State the full name of the test.
- Identify the program of the program that did the calculations
- State all options you selected. Repeated measures? Report enough detail so anyone could start with your data and get precisely the same results you got.

### Graphing Data:

- Present data clearly. Focus on letting the reader see the data, and not only your conclusions.
- When possible, graph the individual data, not a summary of the data.



NIAID

## Outline

### 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

### 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

### 3. Application in Prism and R



NIAID

# Outline

## 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

## 2. Common Statistical Tests

- 2.1. Choose appropriate tests**
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

## 3. Application in Prism and R



NIAID

## 2.1. Choose Appropriate Tests

### We might need to answer some questions:

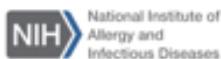
1. *How many independent variables / factors?*
  2. *How many levels do the independent variables have? i.e, how many groups do we have?*
  3. *What is the nature of the independent variables / factors? Continuous? Categorical?*
  4. *What is the nature of the dependent variable? ordinal or nominal?*
  5. *Is the dependent variable normally distributed?*
- .....



NIAID

## Useful References for choosing appropriate tests

1. Prism 8 Statistics Guide ([link](#))
  - This guide examines general principles of statistical analysis, looks at how to conduct those analyses in Prism, and how to interpret results of these analyses.
2. Choosing the Correct Statistical Test in SAS, Stata, SPSS and R ([link](#))
  - A flow chart to choose the best statistical test and the related codes with SAS, Stata, SPSS and R
3. Summary and Analysis of Extension Program Evaluation in R ([link](#))
  - A dictionary of statistical testing methods and the related R codes
4. Analysis Data Model (ADaM) Examples in Commonly Used Statistical Analysis Methods ([link](#))
  - Commonly used statistical analysis methods in clinical trials (proved by FDA)



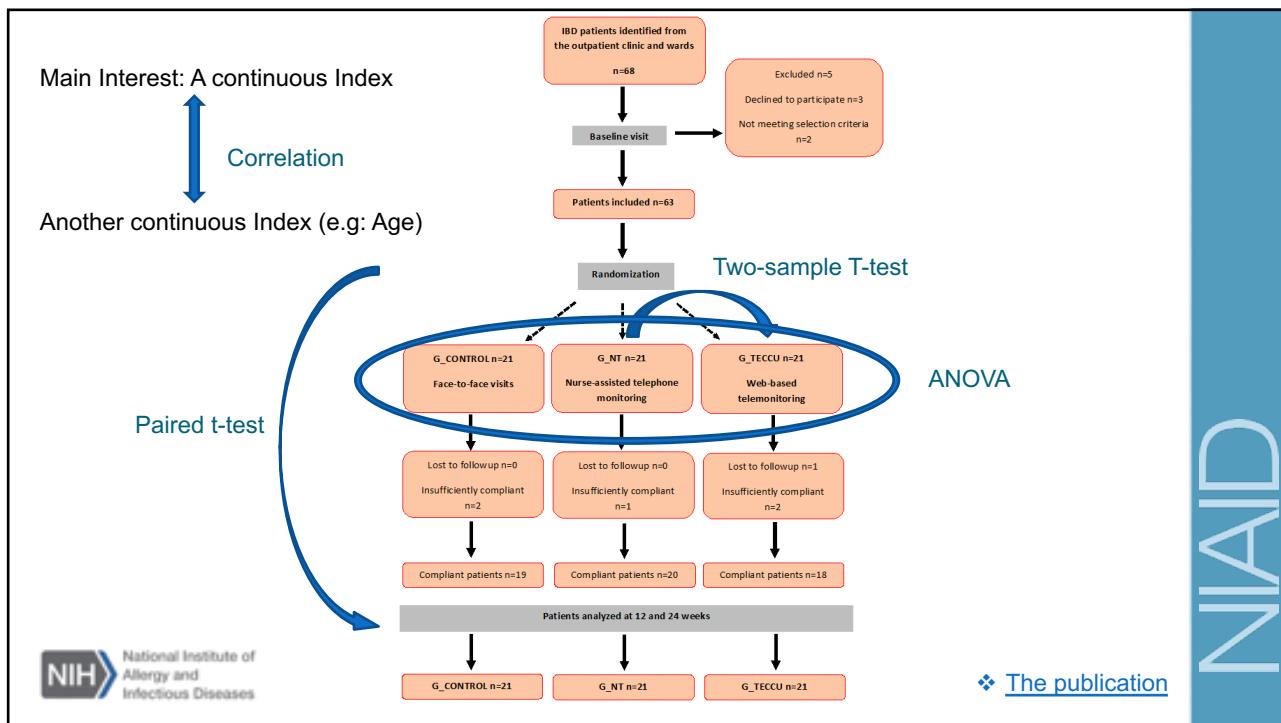
NIAID

### Common Statistical Tests

Type of Test:	Use:
<b>Correlational</b>	These tests look for an association between variables
<b>Pearson correlation</b>	Tests for the strength of the association between two continuous variables
<b>Spearman correlation</b>	Tests for the strength of the association between two ordinal variables (does not rely on the assumption of normal distributed data)
<b>Chi-square</b>	Tests for the strength of the association between two categorical variables
<b>Comparison of Means:</b> <i>look for the difference between the means of variables</i>	
<b>Paired T-test</b>	Tests for difference between two related variables
<b>Independent T-test</b>	Tests for difference between two independent variables
<b>ANOVA</b>	Tests the difference between group means after any other variance in the outcome variable is accounted for
<b>Regression:</b> <i>assess if change in one variable predicts change in another variable</i>	
<b>Simple regression</b>	Tests how change in the predictor variable predicts the level of change in the outcome variable
<b>Multiple regression</b>	Tests how change in the combination of two or more predictor variables predict the level of change in the outcome variable
<b>Non-parametric:</b> <i>are used when the data does not meet assumptions required for parametric tests</i>	
<b>Wilcoxon rank-sum test</b>	Tests for difference between two independent variables - takes into account magnitude and direction of difference
<b>Wilcoxon sign-rank test</b>	Tests for difference between two related variables - takes into account magnitude and direction of difference
<b>Sign test</b>	Tests if two related variables are different – ignores magnitude of change, only takes into account direction



NIAID



Analysis Type	Example	Parametric Procedure	Nonparametric Procedure
Compare means between two distinct/independent groups	Is the mean systolic blood pressure (at baseline) for patients assigned to placebo different from the mean for patients assigned to the treatment group?	Two-sample t-test	Wilcoxon rank-sum test
Compare two quantitative measurements taken from the same individual	Was there a significant change in systolic blood pressure between baseline and the six-month follow-up measurement in the treatment group?	Paired t-test	Wilcoxon signed-rank test
Compare means between three or more distinct/independent groups	If our experiment had three groups (e.g., placebo, new drug #1, new drug #2), we might want to know whether the mean systolic blood pressure at baseline differed among the three groups?	Analysis of variance (ANOVA)	Kruskal-Wallis test
Estimate the degree of association between two quantitative variables	Is systolic blood pressure associated with the patient's age?	Pearson coefficient of correlation	Spearman's rank correlation

National Institute of Allergy and Infectious Diseases

NIAID

# Outline

## 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Calculate the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

## 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing**
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

## 3. Application in Prism and R



NIAID

## 2.2.Data Preprocessing

- Descriptive Statistics
- Data Transformation
- Missing Data



NIAID

## Descriptive Statistics



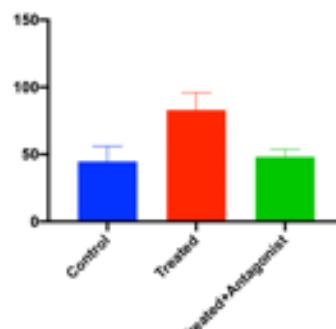
National Institute of  
Allergy and  
Infectious Diseases

Metrics	Meaning
<u>Minimum</u>	The smallest value.
<u>25<sup>th</sup> Percentile</u>	25% of values are lower than this.
<u>Median</u>	Half the values are lower; half are higher.
<u>75<sup>th</sup> Percentile</u>	75% of values are lower than this.
<u>Maximum</u>	The largest value.
<u>Mean</u>	The average.
<u>Standard Deviation</u>	Quantifies variability or scatter.
<u>Standard Error of Mean</u>	Quantifies how precisely the mean is known.
<u>95% Confidence Interval</u>	Given some assumptions, there is a 95% chance that this range includes the true overall mean.
<u>Coefficient of Variation</u>	The standard deviation divided by the mean.
<u>Geometric Mean</u>	Compute the logarithm of all values, compute the mean of the logarithms, and then take the antilog of that mean. It is a better measure of central tendency when data follow a lognormal distribution (long tail).
<u>Harmonic Mean</u>	Compute the reciprocal of all values, compute the mean of the reciprocals, and then take the reciprocal of that mean.
<u>Quadratic Mean</u>	Compute the square of all values, compute the mean of the squares, and then take the square root of that mean.
<u>Skewness</u>	Quantifies how symmetrical the distribution is. A distribution that is symmetrical has a skewness of 0.
<u>Kurtosis</u>	Quantifies whether the tails of the data distribution matches the Gaussian distribution. A Gaussian distribution has a kurtosis of 0.

NIAID

Descriptive statistics	A	B	C
	Control	Treated	Treated+Antagonist
	Y	Y	Y
1 Number of values	6	5	6
2			
3 Minimum	23.00	64.00	39.00
4 25% Percentile	39.50	70.50	43.50
5 Median	46.00	87.00	49.50
6 75% Percentile	54.00	93.50	52.00
7 Maximum	54.00	98.00	55.00
8 Range	31.00	34.00	16.00
9			
10 Mean	44.67	83.00	48.17
11 Std. Deviation	11.40	12.98	5.529
12 Std. Error of Mean	4.652	5.805	2.257

One-way ANOVA data



National Institute of  
Allergy and  
Infectious Diseases

NIAID

## Data Transformation

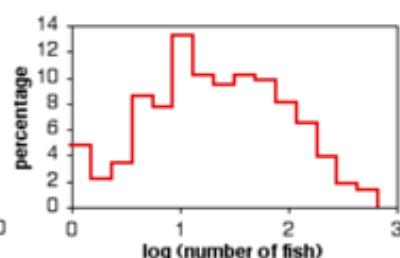
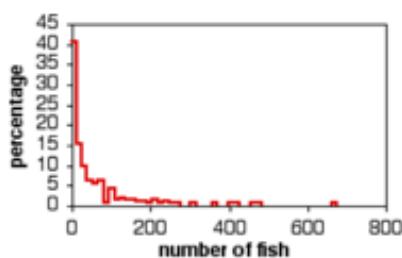
### Assumption: Normal distribution?

If your data do not come from Normal distributions, you have three options:

- Your best option is to transform the values to make the distributions more normal.
- Use the nonparametric procedures instead of parametric procedures.
- Use parametric procedures anyway, knowing that some parametric procedures are fairly robust to violations of a Normal distribution with large samples.



NIAID



Eastern mudminnow (*Gasterosteus aculeatus*).

Histograms of number of Eastern mudminnows per 75 m section of stream (samples with 0 mudminnows excluded). Untransformed data on left, log-transformed data on right.

#### Steps:

- Transform Data
- Do Statistical Test on a transformed data
- Back-transform your results!!!

	Mean	Confidence Interval
Log-transformed	1.044	(0.7, 1.388)
Back-transformed	$10^{1.044}$	$(10^{(0.7)}, 10^{(1.388)})$

❖ More Info



NIAID

## How to choose the right transformation?

- Use a transformation that other researchers commonly use in your field
- Decide which transformation to use before you do the statistical test

Function	Comments
$Y = Y * K$	Enter K in the box provided.
$Y = Y + K$	"
$Y = Y - K$	"
$Y = Y / K$	"
$Y = Y^2$	Enter K in the box provided.
$Y = Y^K$	"
$Y = \log(Y)$	Log base 10
$Y = -\log(Y)$	Natural logarithm (base e)
$Y = \ln(Y)$	Ten to the $Y^{\text{th}}$ power (inverse of log).
$Y = 10^Y$	$e^Y$ (inverse of ln)
$Y = \exp(Y)$	"
$Y = 1/Y$	Square root.
$Y = \sqrt{Y}$	$\ln(Y/1-Y)$
$Y = \logit(Y)$	$Y$ must be between 0.0 and 1.0. See notes below.
$Y = \text{probit}(Y)$	Column rank. Smallest $Y$ value gets rank of 1.
$Y = \text{rank}(Y)$	Number of SDs from the column mean.
$Y = \text{zscore}(Y)$	$Y$ is in radians.
$Y = \sin(Y)$	"
$Y = \cos(Y)$	Result is in radians.
$Y = \tan(Y)$	Absolute value.
$Y = \arcsin(Y)$	Gaussian. Mean=0. SD=K (you enter).
$Y = \text{ABS}(Y)$	"
$Y = Y + \text{Random}$	"
$Y = X / Y$	Enter K in the box provided.
$Y = Y / X$	"
$Y = Y - X$	Logarithm base 2
$Y = Y * X$	2.0 to the power of $Y$
$Y = Y^* X$	"
$Y = X - Y$	"
$Y = K - Y$	"
$Y = K/Y$	"
$Y = \log_2(Y)$	"
$Y = 2^Y$	"

National Institute of  
Allergy and  
Infectious Diseases

## Missing Data

How does Prism handle missing values?

- Some statistical tests work fine with unequal sample size/missing values
  - E.g: Unpaired t, the Mann-Whitney nonparametric test, ordinary one-way ANOVA.
- Some statistical tests require matched sample sizes without missing values
  - If one value is missing, that subject (row) is ignored
  - E.g: Paired t-test, Wilcoxon matched pairs test, ANOVA with repeated measures.

**Complete Case (CC) Analysis**

❖ [More Imputation Methods](#)

NIAID

National Institute of  
Allergy and  
Infectious Diseases

# Outline

## 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

## 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation**
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

## 3. Application in Prism and R



NIAID

## 2.3. Correlation

When two variables vary together, statisticians say that there is a lot of covariation or correlation.



NIAID

## The Difference between Correlation and Regression

	Correlation	Regression
<b>Goal</b>	Quantifies the degree to which two variables are related	Finds the best line that predicts Y from X
<b>Data</b>	Both variables are measured	X is a variable you manipulate
<b>Does it matter which variable is X and which is Y?</b>	No	Yes
<b>Assumptions</b>	Both X and Y are measured and sampled from Gaussian distributions ( <b>Pearson correlation</b> ). OR makes no assumption about the distribution, as the calculation are based on ranks ( <b>Spearman nonparametric correlation</b> )	The X values are not assumed to be sampled from a Gaussian distribution. The distances of the points from the best-fit line is assumed to follow a Gaussian distribution, with the SD of the scatter not related to the X or Y values.
<b>Results</b>	Correlation computes the value of the Pearson correlation coefficient, r. Its value ranges from -1 to +1	Quantifies goodness of fit with $r^2$ , sometimes shown in uppercase as $R^2$



NIAID

## Correlation Coefficient

Value of $r/r_s$	Interpretation
1.0	Perfect correlation
0 to 1	The two variables tend to increase or decrease together.
0.0	The two variables do not vary together at all.
-1 to 0	One variable increases as the other decreases.
-1.0	Perfect negative or inverse correlation.

If  $r$  or  $r_s$  is far from zero, there are four possible explanations:

- Changes in the X variable causes a change the value of the Y variable.
- Changes in the Y variable causes a change the value of the X variable.
- Changes in another variable influence both X and Y.
- X and Y don't really correlate at all, and you just happened to observe such a strong correlation by chance. The P value quantifies the likelihood that this could occur.



NIAID

## r<sup>2</sup>

r<sup>2</sup> is a value that ranges from zero to one, and is the fraction of the variance in the two variables that is “shared”.

### How to explain r<sup>2</sup>?

- 59% of the variance in X can be explained by variation in Y
- 59% of the variance in Y can be explained by variation in X
- 59% of the variance is shared between X and Y

if  $r^2=0.59$

### P-value

- If the P value is small, you can reject the idea that the correlation is due to random sampling.
- If the P value is large, the data do not give you any reason to conclude that the correlation is real.



NIAID

## Try It!

~ 5 mins

- Try to run correlation in Prism
- Try to run correlation in R
- Compare the results of Prism and R



NIAID

# Outline

## 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

## 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test**
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

## 3. Application in Prism and R



NIAID

## 2.4.Two-sample T-test

Two-sample t-test measures whether the means of the measurement variable are different in the two groups. There are two kinds of two-sample t-test: paired and unpaired.

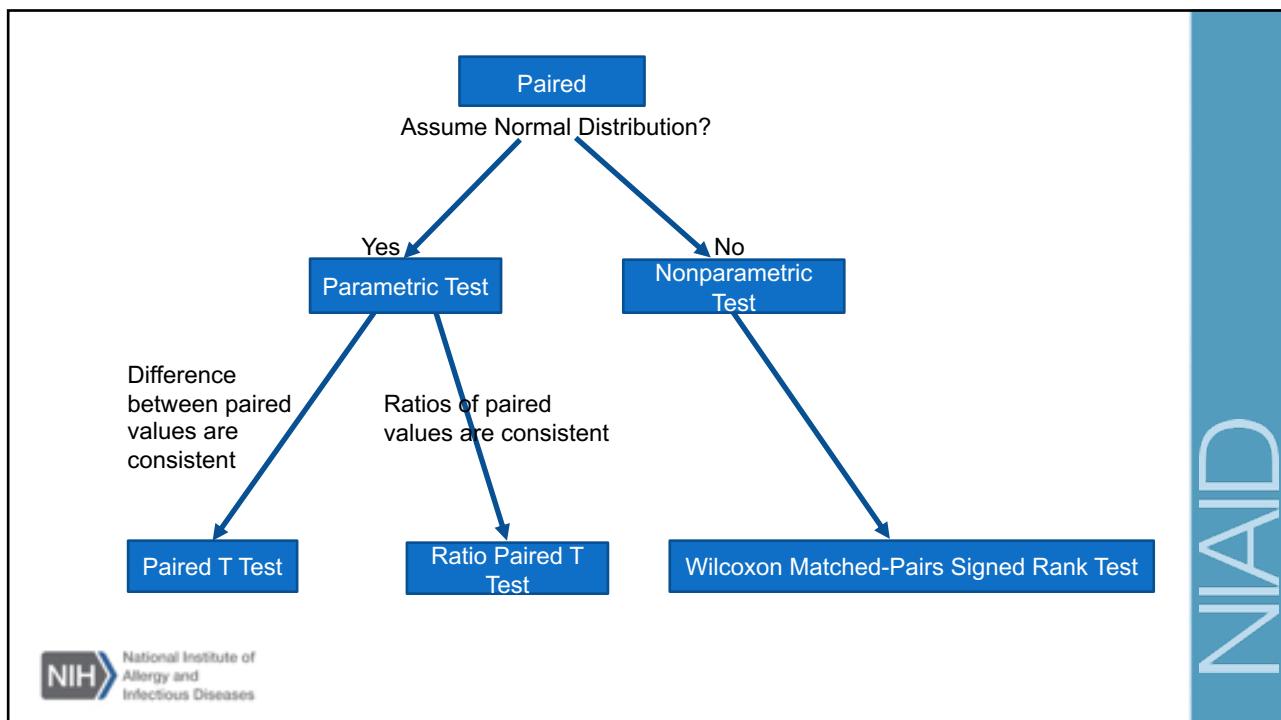
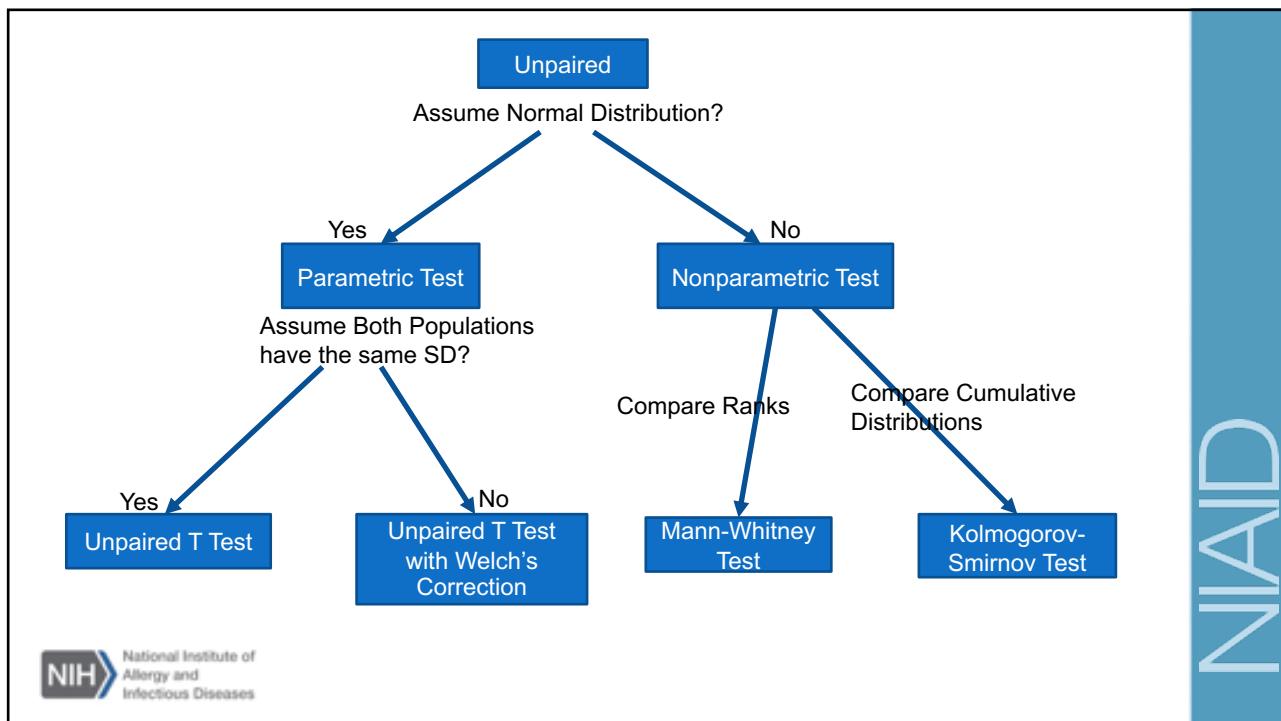
### Paired

Choose a paired test when the columns of data are matched. That means that values on the same row are related to each other. Here are some examples:

- You measure a variable in each subject before and after an intervention.
- You recruit subjects as pairs, matched for variables such as age, ethnic group, and disease severity. One of the pair gets one treatment; the other gets an alternative treatment.
- You run a laboratory experiment several times, each time with a control and treated preparation handled in parallel.
- You measure a variable in twins or child/parent pairs.



NIAID



## Unpaired t test Example in Prism:

There are two groups, the male group which includes 5 subjects while the female group includes 6 subjects:

	Group A		Group B	
	Male	Female	Male	Female
	Y	Y	Y	Y
1	54		43	
2	23		34	
3	45		65	
4	54		77	
5	45		46	
6			65	



NIAID

## Normal distribution?

Prism offers four normality tests. In this case, the null hypothesis is that all the values were sampled from a population that follows a Normal distribution.

- When p-value is high: The data are not inconsistent with a Normal distribution
- When p-value is low: Reject that null hypothesis that the data are not sampled from a Gaussian population.

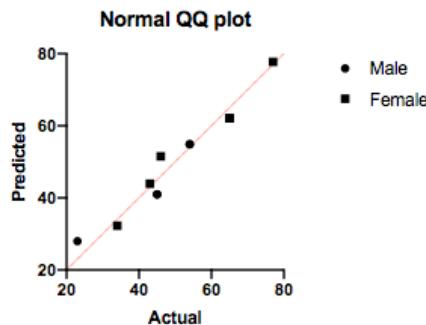
### The difference between different normality tests



	Normality and Lognormality Tests	
	A Male	B Female
1	Test for normal distribution	
2	Anderson-Darling test	
3	A <sup>2*</sup>	N too small
4	P value	N too small
5	Passed normality test (alpha=0.05)?	
6	P value summary	
7		
8	D'Agostino & Pearson test	
9	K <sup>2</sup>	N too small
10	P value	N too small
11	Passed normality test (alpha=0.05)?	
12	P value summary	
13		
14	Shapiro-Wilk test	
15	W	0.8070
16	P value	0.0923
17	Passed normality test (alpha=0.05)?	Yes
18	P value summary	ns
19		
20	Kolmogorov-Smirnov test	
21	KS distance	0.3252
22	P value	>0.1000
23	Passed normality test (alpha=0.05)?	Yes
24	P value summary	ns
25		
26	Number of values	5
		6

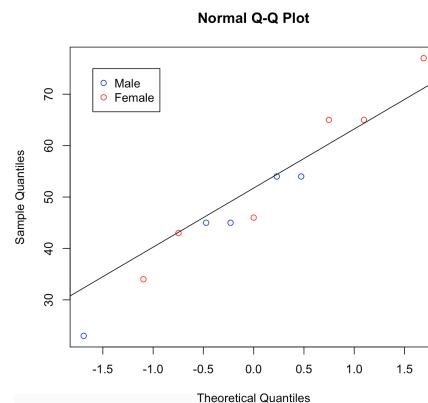
## QQ Normality Plots

[QQ plot in Prism](#)



National Institute of  
Allergy and  
Infectious Diseases

[QQ plot in R](#)



NIAID

## How to check Same Standard Deviation?

### Formula of F-test

The test statistic can be obtained by computing the ratio of the two variances  $S_A^2$  and  $S_B^2$ .

$$F = \frac{S_A^2}{S_B^2}$$

The degrees of freedom are  $n_A - 1$  (for the numerator) and  $n_B - 1$  (for the denominator).

✓ Note that, the more this ratio deviates from 1, the stronger the evidence for unequal population variances.

✗ Note that, the F-test requires the two samples to be normally distributed.



National Institute of  
Allergy and  
Infectious Diseases

NIAID

Unpaired t test		Welch's t test	
1	Table Analyzed	Unpaired t test data	1
2		2	Table Analyzed
3	Column B	Female	3
4	vs.	vs.	4
5	Column A	Male	5
6		6	Column A
7	Unpaired t test		7
8	P value	0.2613	8
9	P value summary	ns	9
10	Significantly different ( $P < 0.05$ )?	No	10
11	One- or two-tailed P value?	Two-tailed	11
12	t, df	t=1.199, df=9	12
13		Welch-corrected t, df	13
14	How big is the difference?		14
15	Mean of column A	44.20	15
16	Mean of column B	55.00	16
17	Difference between means (B - A) $\pm$ SEM	10.80 $\pm$ 9.010	17
18	95% confidence interval	-9.583 to 31.18	18
19	R squared (eta squared)	0.1377	19
20		20	
21	F test to compare variances		21
22	F, DFn, Dfd	1.680, 5, 4	22
23	P value	0.6354	23
24	P value summary	ns	24
25	Significantly different ( $P < 0.05$ )?	No	25
26		26	
27	Data analyzed		27
28	Sample size, column A	5	28
29	Sample size, column B	6	29

> res.ftest

F test to compare two variances

```

data: Y by Gender
F = 1.6801, num df = 5, denom df = 4, p-value = 0.6354
alternative hypothesis: true ratio of variances is not equal to 1
95 percent confidence interval:
 0.1794174 12.4127514
sample estimates:
ratio of variances
 1.680149

```

National Institute of  
Allergy and  
Infectious Diseases

Unpaired t test		Welch's t test	
1	Table Analyzed	Unpaired t test data	1
2		2	Table Analyzed
3	Column B	Female	3
4	vs.	vs.	4
5	Column A	Male	5
6		6	Column A
7	Unpaired t test		7
8	P value	0.2613	8
9	P value summary	ns	9
10	Significantly different ( $P < 0.05$ )?	No	10
11	One- or two-tailed P value?	Two-tailed	11
12	t, df	t=1.199, df=9	12
13		Welch-corrected t, df	13
14	How big is the difference?		14
15	Mean of column A	44.20	15
16	Mean of column B	55.00	16
17	Difference between means (B - A) $\pm$ SEM	10.80 $\pm$ 9.010	17
18	95% confidence interval	-9.583 to 31.18	18
19	R squared (eta squared)	0.1377	19
20		20	
21	F test to compare variances		21
22	F, DFn, Dfd	1.680, 5, 4	22
23	P value	0.6354	23
24	P value summary	ns	24
25	Significantly different ( $P < 0.05$ )?	No	25
26		26	
27	Data analyzed		27
28	Sample size, column A	5	28
29	Sample size, column B	6	29

> res

Two Sample t-test

```

data: Y by Gender
t = 1.1986, df = 9, p-value = 0.2613
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -9.583046 31.183046
sample estimates:
mean in group Female mean in group Male
 55.0               44.2

```

National Institute of  
Allergy and  
Infectious Diseases

## How to explain the results?

- When the **p-value** of the test is less than the significance level alpha = 0.05, We can conclude that male's average is significantly different from female's average.
- When the **p-value** of the test is larger than the significance level alpha = 0.05, we *cannot* conclude that a significant difference exists between male and female.



NIAID

## Try It!

~ 5 mins

- Try the same process of unpaired two-sample t-test in Prism
- Try to run the R codes
- Compare the results of Prism and R



NIAID

# Outline

## 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

## 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons**
- 2.6. Survival Analysis

## 3. Application in Prism and R



NIAID

# 2.5.ANOVA

## One-Way ANOVA

The one-way ANOVA compares the means between the groups you are interested in and determines whether any of those means are statistically significantly different from each other.

Specifically, it tests the null hypothesis:

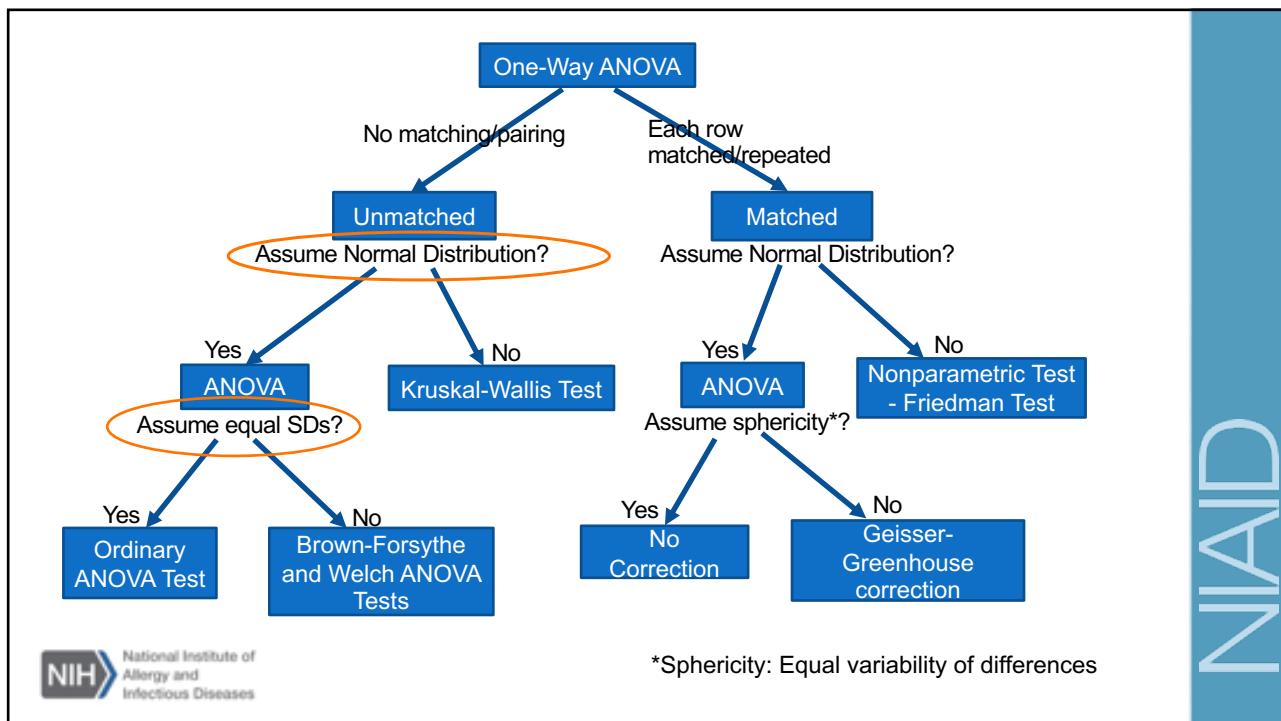
$$H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_k$$

Where  $\mu$  = group mean and  $k$  = number of groups. If the one-way ANOVA returns a statistically significant result, we reject the null hypothesis and conclude that at least two group means are statistically significantly different from each other.

Note: One-way ANOVA is an **omnibus** test statistic



NIAID



## Sphericity

Briefly sphericity means that you waited long enough between treatments for any treatment effect to wash away.

		A	B	C	D
		Control	Treatment 1	Treatment 2	Treatment 3
		Y	Y	Y	Y
1	GS	54	43	78	111
2	JM	23	34	65	99
3	HM	45	65	99	78
4	DR	54	77	79	90
5	PS	45	46	87	95

### Should we assume sphericity?

- If your experimental design relies on matching rather than repeated measurements, then you can assume sphericity, as violations are essentially impossible.
- If your experiment design is repeated measures, we recommend that you do not assume sphericity.

In this case, it is ANOVA with repeated measures.



## Unmatched One-Way ANOVA example

The following is an example provided by Prism. In this case, there are three groups: *Control* (n = 6), *Treated* (n = 5) and *Treated + Antagonist* (n = 6).

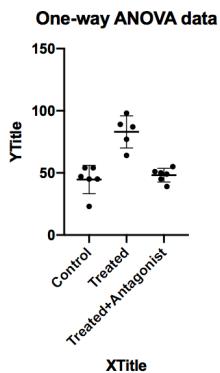
	Group A	Group B	Group C
	Control	Treated	Treated+Antagonist
	Y	Y	Y
1	54	87	45
2	23	98	39
3	45	64	51
4	54	77	49
5	45	89	50
6	47		55



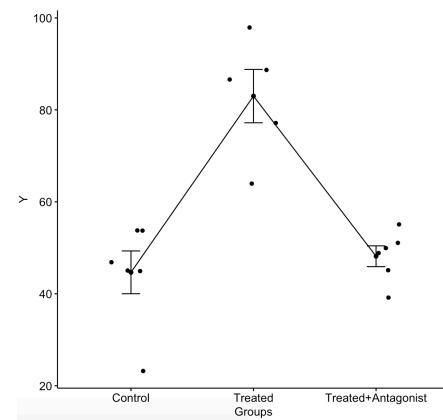
NIAID

## Plots

box-plot in Prism



box-plot in R



NIAID

Considering that the first group, the control group doesn't pass the normality test, it is better to use nonparametric test – Kruskal-Wallis Test

Normality and Lognormality Tests Tabular results		A	B	C
		Control	Treated	Treated+Antagonist
1	Test for normal distribution		Y	
2	Anderson-Darling test		Y	
3	A2*		Y	
4	P value			
5	Passed normality test (alpha=0.05)?			
6	P value summary			
7				
8	D'Agostino & Pearson test			
9	K2		N too small	N too small
10	P value			
11	Passed normality test (alpha=0.05)?			
12	P value summary			
13				
14	Shapiro-Wilk test			
15	W	0.7890	0.9643	0.9517
16	P value	0.0467	0.8376	0.7537
17	Passed normality test (alpha=0.05)?	No	Yes	Yes
18	P value summary	*	ns	ns
19				
20	Kolmogorov-Smirnov test			
21	KS distance	0.3450	0.2210	0.2266
22	P value	0.0245	>0.1000	>0.1000
23	Passed normality test (alpha=0.05)?	No	Yes	Yes
24	P value summary	*	ns	ns
25				
26	Number of values	6	5	6

 National Institute of  
Allergy and  
Infectious Diseases

## Kruskal-Wallis Test?

Kruskal-Wallis test by ranks is also called as one-way ANOVA on ranks, which extends the Mann–Whitney U test, which is used for comparing only two groups

To perform this test, Prism first ranks all the values from low to high, paying no attention to which group each value belongs. The smallest number gets a rank of 1. The largest number gets a rank of N, where N is the total number of values in all the groups. The discrepancies among the rank sums are combined to create a single value called the Kruskal-Wallis statistic (some books refer to this value as H). A large Kruskal-Wallis statistic corresponds to a large discrepancy among rank sums.

 National Institute of  
Allergy and  
Infectious Diseases

As the p-value is less than the significance level 0.05, we can conclude that there are significant differences between these treatment groups.

```
> kruskal.test(Y ~ Group, data = unmatched_one_way_anova)

Kruskal-Wallis rank sum test

data: Y by Group
Kruskal-Wallis chi-squared = 10.18, df = 2, p-value = 0.006158
```

Kruskal-Wallis test	
1	Table Analyzed
2	One-way ANOVA data
3	Kruskal-Wallis test
4	P value
5	Exact or approximate P value?
6	Exact
7	P value summary
8	**
9	Do the medians vary signif. (P < 0.05)?
10	Yes
11	Number of groups
12	3
13	Kruskal-Wallis statistic
	10.18
	Data summary
	Number of treatments (columns)
	3
	Number of values (total)
	17



NIAID

### Assumption: Normal distribution?

The assumption - Normal distribution, is not too important with large samples, but it is important with small sample sizes (especially with unequal sample sizes).

If your data do not come from Normal distributions, you have three options:

- Your best option is to transform the values (perhaps to logs or reciprocals) to make the distributions more normal.
- Use the Kruskal-Wallis nonparametric test instead of ANOVA.
- Use ANOVA anyway, knowing that it is fairly robust to violations of a Normal distribution with large samples.



NIAID

## Assumption: same standard deviation?

- The assumption, same standard deviation, is not very important when all the groups have the same (or almost the same) number of subjects, but is very important when sample sizes differ.
- Prism tests for equality of variance with two tests: The Browne-Forsythe test and Bartlett's test.
- Don't base your conclusion solely on these tests. Also think about data from other similar experiments. If you have plenty of previous data that convinces you that the variances are really equal, ignore these tests (unless the P value is really tiny) and interpret the ANOVA results as usual.



NIAID

Let's run ANOVA anyway to see the output provided by Prism and R.

\*\* BUT in this case, ANOVA might not be acceptable because of the violation of Normal distribution.

```
> summary(res.aov)
   Df Sum Sq Mean Sq F value    Pr(>F)
Group     2   4760   2380.2   22.57 4.16e-05 ***
Residuals 14   1476    105.4
---
Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Ordinary one-way ANOVA					
1	Table Analyzed				One-way ANOVA data
2	Data sets analyzed				A-C
3					
4	ANOVA summary				
5	F	22.57			
6	P value	<0.0001			
7	P value summary	***			
8	Significant diff. among means (P < 0.05)?	Yes			
9	R square	0.7633			
10					
11	Brown-Forsythe test				
12	F (DFn, DFd)	0.7307 (2, 14)			
13	P value	0.4991			
14	P value summary	ns			
15	Are SDs significantly different (P < 0.05)?	No			
16					
17	Bartlett's test				
18	Bartlett's statistic (corrected)	2.986			
19	P value	0.2247			
20	P value summary	ns			
21	Are SDs significantly different (P < 0.05)?	No			
22					
23	ANOVA table	SS	DF	MS	F (DFn, DFd)    P value
24	Treatment (between columns)	4760	2	2380	F (2, 14) = 22.57 P<0.0001
25	Residual (within columns)	1476	14	105.4	
26	Total	6236	16		
27					
28	Data summary				
29	Number of treatments (columns)	3			
30	Number of values (total)	17			



## Try It!

~ 5 mins

- Try the same process of unmatched one-way ANOVA in Prism
- Try to run the R codes
- Compare the results of Prism and R



NIAID

### Repeated measures one-way ANOVA example:

This example is provided by Prism. The four columns represent four sequential treatments. Each row represents a different subject (or a different set of matched data).

Repeated-measures ANOVA compares the means of three or more matched groups.

Table format: Column		Group A	Group B	Group C	Group D
		Control	Treatment 1	Treatment 2	Treatment 3
1	GS	54	43	78	111
2	JM	23	34	37	41
3	HM	45	65	99	78
4	DR	31	33	36	35
5	PS	15	25	30	26



NIAID

Check if it follows normal distribution

And the answer is YES!

Normality and Lognormality Tests Tabular results		A	B	C	D
		Control	Treatment 1	Treatment 2	Treatment 3
1	<b>Test for normal distribution</b>		Y		Y
2	<b>Anderson-Darling test</b>				
3	A <sup>2</sup> *	N too small	N too small	N too small	N too small
4	P value				
5	Passed normality test (alpha=0.05)?				
6	P value summary				
7					
8	<b>D'Agostino &amp; Pearson test</b>				
9	K <sub>2</sub>	N too small	N too small	N too small	N too small
10	P value				
11	Passed normality test (alpha=0.05)?				
12	P value summary				
13					
14	<b>Shapiro-Wilk test</b>				
15	W	0.9646	0.8908	0.8317	0.8854
16	P value	0.8393	0.3613	0.1434	0.3345
17	Passed normality test (alpha=0.05)?	Yes	Yes	Yes	Yes
18	P value summary	ns	ns	ns	ns
19					
20	<b>Kolmogorov-Smirnov test</b>				
21	KS distance	0.1649	0.2519	0.3320	0.2858
22	P value	>0.1000	>0.1000	0.0750	>0.1000
23	Passed normality test (alpha=0.05)?	Yes	Yes	Yes	Yes
24	P value summary	ns	ns	ns	ns
25					
26	<b>Number of values</b>	5	5	5	5

**NIH** National Institute of Allergy and Infectious Diseases

- The P value answers this question: If all the populations really have the same mean, what is the chance that random sampling would result in means as far apart (or more so) as observed in this experiment?
- If the overall P value is large: the data does not give you any reason to conclude that the means differ.
- If the overall P value is small: it is unlikely that the differences you observed are due to random sampling. You can reject the idea that all the populations have identical means.

RM one-way ANOVA					
	Table Analyzed	Repeated measures one-way ANOVA data			
1	Repeated measures ANOVA summary				
2	Assume sphericity?	No			
3	F	3.618			
4	P value	0.0873			
5	P value summary	ns			
6	Statistically significant (P < 0.05)?	No			
7	Geisser-Greenhouse's epsilon	0.5793			
8	R square	0.4750			
9					
10	Was the matching effective?				
11	F	10.45			
12	P value	0.0007			
13	P value summary	***			
14	Is there significant matching (P < 0.05)?	Yes			
15	R square	0.6464			
16					
17	ANOVA table	SS	DF	MS	F (DFn, DFd)
18	Treatment (between columns)	2175	3	725.0	F 1.738, 6.951 = 3.618
19	Individual (between rows)	8372	4	2093	P=0.0873
20	Residual (random)	2404	12	200.4	P=0.0007
21	Total	12951	19		
22					
23	Data summary				
24	Number of treatments (columns)	4			
25	Number of subjects (rows)	5			
26	Number of missing values	0			

**NIH** National Institute of Allergy and Infectious Diseases

```
> summary(model1,multivariate=F)

Univariate Type III Repeated-Measures ANOVA Assuming Sphericity

  Sum Sq num Df Error SS den Df F value Pr(>F)
(Intercept) 44086      1   8371.7      4 21.0643 0.01011 *
Trials       2175      3   2404.3     12  3.6184 0.04549 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Mauchly Tests for Sphericity

  Test statistic p-value
Trials      0.057644 0.19416

Greenhouse-Geisser and Huynh-Feldt Corrections
for Departure from Sphericity

  GG eps Pr(>F[GG])
Trials 0.57928  0.08727 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

  HF eps Pr(>F[HF])
Trials 0.9856654 0.04648165
```



NIAID

## Try It!

~ 5 mins

- Try the same process of Repeated measures one-way ANOVA
- Try to run the R codes
- Compare the results of Prism and R



NIAID

## Two-way ANOVA

Two-way ANOVA, also called two-factor ANOVA, is used to evaluate simultaneously the effect of two grouping variables (A and B) on a response variable.

### Two-way ANOVA test hypotheses

1. There is no difference in the means of factor A
2. There is no difference in means of factor B
3. There is no interaction between factors A and B

The alternative hypothesis for cases 1 and 2 is: the means are not equal.

The alternative hypothesis for case 3 is: there is an interaction between A and B.



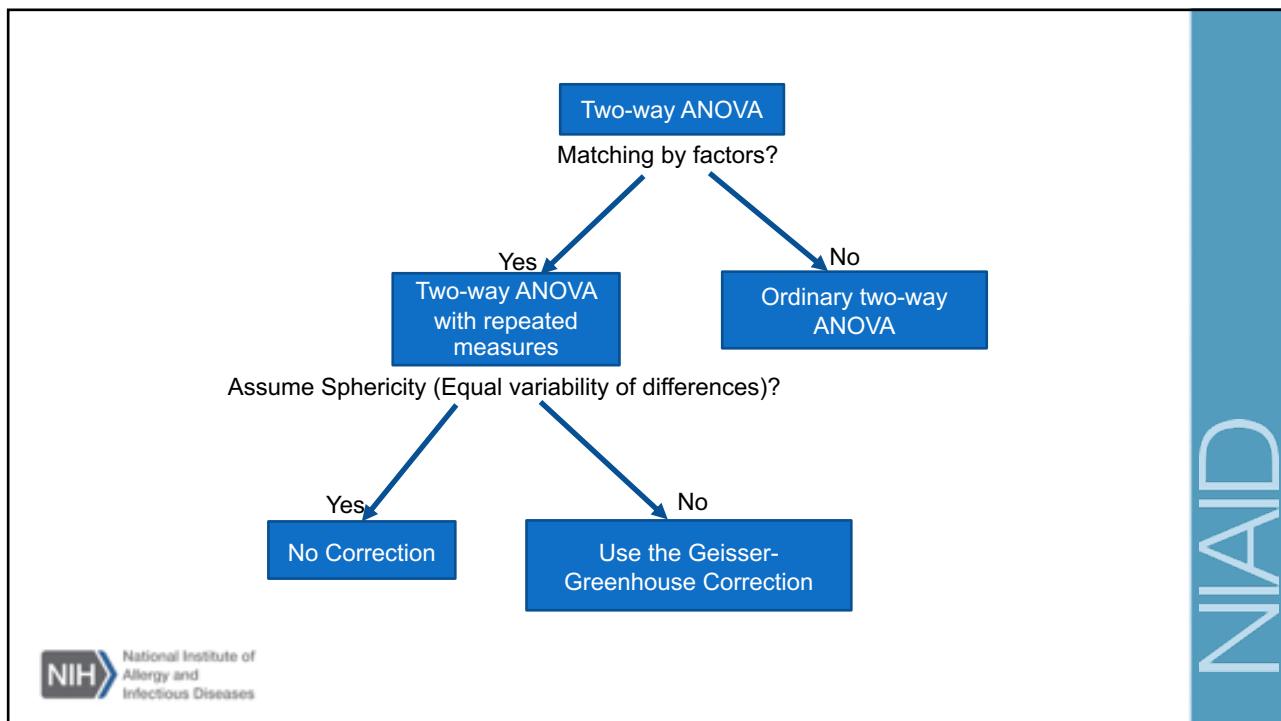
NIAID

One-Way vs Two-Way ANOVA Differences Chart

	One-Way ANOVA	Two-Way ANOVA
Definition	A test that allows one to make comparisons between the means of three or more groups of data.	A test that allows one to make comparisons between the means of three or more groups of data, where two independent variables are considered.
Number of Independent Variables	One.	Two.
What is Being Compared?	The means of three or more groups of an independent variable on a dependent variable.	The effect of multiple groups of two independent variables on a dependent variable and on each other.
Number of Groups of Samples	Three or more.	Each variable should have multiple samples.



NIAID



### Ordinary Two-way ANOVA Example

The columns represent two cell lines. The rows represent two treatments. Within each treatment for each cell line, five replicate values are entered into subcolumns. This experiment has no matching or repeated measures. For this reason, it is OK that one of the values is missing

Table format: Grouped		Group A					Group B				
		Wild-type cells					GPP5 cell line				
		A:Y1	A:Y2	A:Y3	A:Y4	A:Y5	B:Y1	B:Y2	B:Y3	B:Y4	B:Y5
1	Serum starved	34	36	41		43	98	87	95	99	88
2	Normal culture	23	19	26	29	25	32	29	26	33	30



**ANOVA table**

The ANOVA table breaks down the overall variability between measurements (expressed as the sum of squares) into four components:

- *Interactions between row and column.*
- *Variability among columns.*
- *Variability among rows.*
- *Residual or error. Variation among replicates not related to systematic differences between rows and columns.*

ANOVA results					
2way ANOVA ANOVA results					
1	Table Analyzed	Two-way ANOVA , not RM			
2		Ordinary			
3	Two-way ANOVA				
4	Alpha	0.05			
5					
6	Source of Variation	% of total variation	P value	P value summary	Significant?
7	Interaction	18.67	<0.0001	****	Yes
8	Row Factor	46.14	<0.0001	****	Yes
9	Column Factor	28.12	<0.0001	****	Yes
10					
11	ANOVA table	SS (Type III)	DF	MS	F (DFn, DFd) P value
12	Interaction	2859	1	2859	F (1, 15) = 162.8 P<0.0001
13	Row Factor	7066	1	7066	F (1, 15) = 402.4 P<0.0001
14	Column Factor	4306	1	4306	F (1, 15) = 245.2 P<0.0001
15	Residual	263.4	15	17.56	
16					
17	Difference between column means				
18	Predicted (LS) mean of Wild-type cells	31.45			
19	Predicted (LS) mean of GPP5 cell line	61.70			
20	Difference between predicted means	-30.25			
21	SE of difference	1.932			
22	95% CI of difference	-34.37 to -26.13			
23					
24	Difference between row means				
25	Predicted (LS) mean of Serum starved	65.95			
26	Predicted (LS) mean of Normal culture	27.20			
27	Difference between predicted means	38.75			
28	SE of difference	1.932			
29	95% CI of difference	34.63 to 42.87			
30					
31	Interaction CI				
32	Mean diff, A1 - B1	-54.90			
33	Mean diff, A2 - B2	-5.600			
34	(A1 -B1) - (A2 - B2)	-49.30			
35	95% CI of difference	-57.53 to -41.07			
36	(B1 -A1) - (B2 - A2)	49.30			
37	95% CI of difference	41.07 to 57.53			

National Institute of Allergy and Infectious Diseases

### Interaction P value

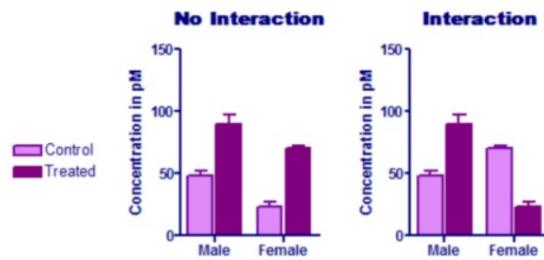
The null hypothesis is that there is no interaction between columns and rows. Often the test of interaction is the most important of the three tests. For example, If columns represent drugs and rows represent gender, then the null hypothesis is that the differences between the drugs are consistent for men and women.

```
> summary(res.aov2)
Df Sum Sq Mean Sq F value Pr(>F)
Treatment      1   8276   8276  471.3 9.58e-13 ***
Cell           1    3917   3917  223.1 2.06e-10 ***
Treatment:Cell 1   2859   2859  162.8 1.86e-09 ***
Residuals     15    263     18
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
1 observation deleted due to missingness
```

National Institute of Allergy and Infectious Diseases

The graph on the left below shows no interaction. The treatment has about the same effect in males and females. The graph on the right, in contrast, shows a huge interaction. the effect of the treatment is completely different in males (treatment increases the concentration) and females (where the treatment decreases the concentration).

In this two-way ANOVA example, the test for interaction leads to statistically significant results, we probably won't learn anything of interest from the other two P values. In the example above, a statistically significant interaction means that the effect of the treatment (difference between treated and control) differs between cell groups.



NIAID

## Try It!

~ 5 mins

- Try the same process of Ordinary Two-way ANOVA
- Try to run the R codes
- Compare the results of Prism and R

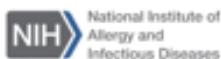


NIAID

## Multiple Comparisons

Once ANOVA test has been completed, what's next?

- Subgroup differences among the different experimental and control groups
- The subgroup differences are called "pairwise" differences

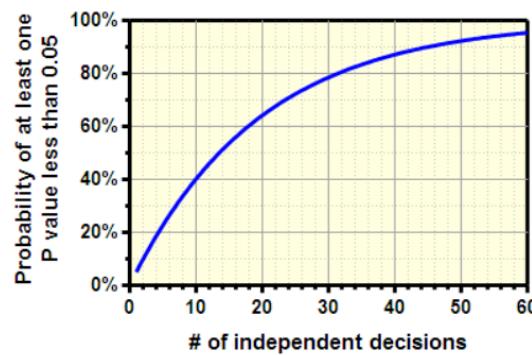


National Institute of  
Allergy and  
Infectious Diseases

❖ [Multiple Comparisons in Prism](#)

NIAID

Interpreting multiple P values is difficult. If you test several independent null hypotheses and leave the threshold at 0.05 for each comparison, the chance of obtaining at least one "statistically significant" result is greater than 5% (even if all null hypotheses are true).

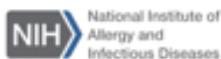


National Institute of  
Allergy and  
Infectious Diseases

NIAID

## Three Approaches to Multiple Comparisons

1. Don't correct for multiple comparisons
2. Control the Type I error rate for the family of comparisons
3. Control the False Discovery Rate (FDR)



NIID

### Approach 1: Don't correct for multiple comparisons

#### **When it makes sense to not correct for multiple comparisons?**

1. Account for multiple comparisons when interpreting the results rather than in the calculations
2. Corrections for multiple comparisons may not be needed if you make only a few planned comparisons
3. Corrections for multiple comparisons are not needed when the comparisons are complementary

Fisher's Least Significant Difference (LSD)



NIID

## Approach 2: Control Type I error rate for the family of comparisons

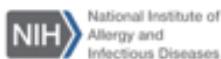
### What are multiplicity adjusted P values?

The P value from a single comparison is the answer to two equivalent questions:

- If the null hypothesis were true, what is the chance that random sampling would result in a difference this large or larger?
- What is the smallest definition of the threshold (alpha) of statistical significance at which this result would be statistically significant?

The adjusted P-value answers this question:

- What is the smallest significance level, when applied to the entire family of comparisons, at which this particular comparison will be deemed statistically significant?



NIAID

### Methods:

	Bonferroni	Sidak	Holm-Sidak	Tukey	Dunnett	Newman-Keuls
Only as follow-up test to ANOVA?	No	No	No	Yes	Yes	
Compare every mean with every other mean?	Yes	Yes	Yes	Yes	No	Yes
Compare means within rows/columns	Yes	Yes	Yes	No	No	No
Compare every mean to a control mean	Yes	Yes	Yes	No	Yes	No
CI?	Yes	Yes	No	Yes	Yes	No
Multiplicity adjusted P values	Yes	Yes	Yes	Yes	Yes	No
Assume independence?	No	Yes				

❖ [More details](#)



NIAID

### Approach 3: Control the False Discovery Rate (FDR)

	"Discovery"	"Not a discovery"	Total
No difference. Null hypothesis true	A	B	A+B
A difference truly exists	C	D	C+D
Total	A+C	B+D	A+B+C+D

**Discovery:** Statistically significant

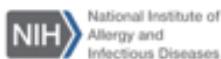
**False Discovery rate:** The ratio  $A/(A+C)$  Usually called Q

The usual approach to statistical significance and multiple comparisons asks the question:

- If the null hypothesis is true what is the chance of getting "statistically significant" results?

The False Discovery Rate (FDR) answers a different question:

- If a comparison is a "discovery", what is the chance that the null hypothesis is true?



NIAID

### Three Methods used to Control the FDR:

1. Original Method of Benjamin and Hochberg
  - Assumes "test statistics are independent or positive dependent"
2. Two-stage step-up method of Benjamini, Krieger and Yekutieli
  - has more power than the Benjamini and Hochberg method, while making the same assumptions
3. Corrected method of Benjamini & Yakutiel
  - Requires no assumption about how the various comparisons correlate with each other.
  - has less power; The method is very conservative.

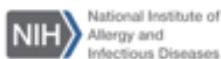


NIAID

## Ordinary Two-way ANOVA Example (with multiple comparisons)

The columns represent two cell lines. The rows represent two treatments. Within each treatment for each cell line, five replicate values are entered into subcolumns. This experiment has no matching or repeated measures. For this reason, it is OK that one of the values is missing

Table format: Grouped		Group A					Group B				
		Wild-type cells					GPP5 cell line				
		A:Y1	A:Y2	A:Y3	A:Y4	A:Y5	B:Y1	B:Y2	B:Y3	B:Y4	B:Y5
1	Serum starved	34	36	41		43	98	87	95	99	88
2	Normal culture	23	19	26	29	25	32	29	26	33	30



NIAID

Uncorrected Fisher's LSD	Mean Diff.	95.00% CI of diff.	Significant?	Summary	Individual P Value
Serum starved:Wild-type cells vs. Serum starved:GPP5 cell line	-54.90	-60.89 to -48.91	Yes	****	<0.0001
Serum starved:Wild-type cells vs. Normal culture:Wild-type cells	14.10	8.108 to 20.09	Yes	***	0.0002
Serum starved:Wild-type cells vs. Normal culture:GPP5 cell line	8.500	2.508 to 14.49	Yes	**	0.0085
Serum starved:GPP5 cell line vs. Normal culture:Wild-type cells	69.00	63.35 to 74.65	Yes	****	<0.0001
Serum starved:GPP5 cell line vs. Normal culture:GPP5 cell line	63.40	57.75 to 69.05	Yes	****	<0.0001
Normal culture:Wild-type cells vs. Normal culture:GPP5 cell line	-5.600	-11.25 to 0.04894	No	ns	0.0518

Tukey's multiple comparisons test	Mean Diff.	95.00% CI of diff.	Significant?	Summary	Adjusted P Value
Serum starved:Wild-type cells vs. Serum starved:GPP5 cell line	-54.90	-63.00 to -46.80	Yes	****	<0.0001
Serum starved:Wild-type cells vs. Normal culture:Wild-type cells	14.10	5.998 to 22.20	Yes	***	0.0008
Serum starved:Wild-type cells vs. Normal culture:GPP5 cell line	8.500	0.3981 to 16.60	Yes	*	0.0383
Serum starved:GPP5 cell line vs. Normal culture:Wild-type cells	69.00	61.36 to 76.64	Yes	****	<0.0001
Serum starved:GPP5 cell line vs. Normal culture:GPP5 cell line	63.40	55.76 to 71.04	Yes	****	<0.0001
Normal culture:Wild-type cells vs. Normal culture:GPP5 cell line	-5.600	-13.24 to 2.039	No	ns	0.1936



NIAID

## Tukey multiple pairwise-comparisons

```
> TukeyHSD(res.aov2)
  Tukey multiple comparisons of means
  95% family-wise confidence level

  Fit: aov(formula = Y ~ Treatment + Cell + Treatment:Cell, data = two_way_anova)

  $Treatment
    diff      lwr      upr p adj
Serum Starved-Normal Culture 41.8 37.69614 45.90386   0

  $Cell
    diff      lwr      upr p adj
Wild-Type cells-GPP5 cell line -28.71111 -32.81497 -24.60725   0

  `$Treatment:Cell`
    diff      lwr      upr p adj
Serum Starved:GPP5 cell line-Normal Culture:GPP5 cell line 63.4 55.7614904 71.03851 0.0000000
Normal Culture:Wild-Type cells-Normal Culture:GPP5 cell line -5.6 -13.2385096 2.03851 0.1936359
Serum Starved:Wild-Type cells-Normal Culture:GPP5 cell line 8.5 0.3981371 16.60186 0.0382505
Normal Culture:Wild-Type cells-Serum Starved:GPP5 cell line -69.0 -76.6385096 -61.36149 0.0000000
Serum Starved:Wild-Type cells-Serum Starved:GPP5 cell line -54.9 -63.0018629 -46.79814 0.0000000
Serum Starved:Wild-Type cells-Normal Culture:Wild-Type cells 14.1 5.9981371 22.20186 0.0007910
```



NIAD

## Try It!

~ 5 mins

- Try the same process of Ordinary Two-way ANOVA (with multiple comparisons)
- Try to run the R codes
- Compare the results of Prism and R



NIAD

## Two interesting facts about Multiple Comparisons

**1. If the overall ANOVA finds a significant difference among groups, am I certain to find a significant post multiple comparison test?**

This is not necessarily true

Because It is possible that the overall mean of group A and group B combined differs significantly from the combined mean of groups C, D and E.

**2. If the overall ANOVA finds no significant difference among groups, are the multiple comparisons test results valid?**

All the multiple comparisons tests offered by Prism are valid even if the overall ANOVA did not find a significant difference among means.



NIAID

## Outline

### 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

### 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis**

### 3. Application in Prism and R



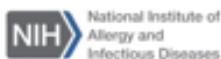
NIAID

## 2.6.Survival Analysis

In many clinical and animal studies, the outcome is survival time. The goal of the study is to determine whether a treatment changes survival. Prism creates survival curves, using the product limit method of Kaplan and Meier, and compares survival curves using both the **log-rank test** and the **Gehan-Wilcoxon test**.

Examples of **events** and **censor**:

Events	Censor
1. Death	1. Still alive at the end of the study
2. Targeted Events	2. Drop out of the study



NIAID

### The comparison between log-rank test and Gehan-Breslow-Wilcoxon test

Log-rank Test	Gehan-Breslow-Wilcoxon test
<ul style="list-style-type: none"> <li>• Gives equal weights to all time points</li> <li>• More standard</li> <li>• More powerful if the assumptions of proportional hazards** is true.</li> </ul>	<ul style="list-style-type: none"> <li>• Gives more weight to deaths at early time points</li> <li>• The results can be misleading when a large fraction of patients is censored at early time points.</li> <li>• Doesn't require a consistent hazard ratio, but does require that one group consistently have a higher risk than the other.</li> </ul>

\*\* Proportional hazards mean that the ratio of hazard functions (deaths per time) is the same at all time points. One example of proportional hazards would be if the control group died at twice the rate as treated group at all time points.



NIAID

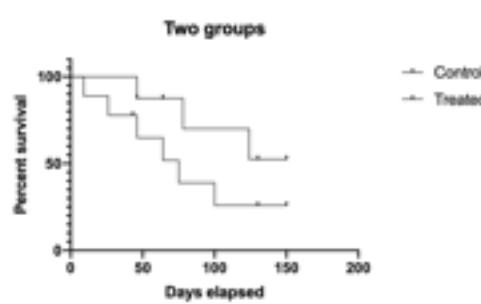
## Survival Analysis – Comparing Two Groups

Each row represents one subject. The X values are time. The Y values are entered into two columns that define the treatment groups

The Y value is "1" when the subject died at the specified time, and "0" when the subject's data was censored at that time.



X	Group A	Group B
Days elapsed	Control	Treated
X	Y	Y
46	1	
46	0	
64	0	
78	1	
124	1	
130	0	
150	0	
150	0	
9		1
26		1
43		0
46		1
64		1
75		1
100		1
130		0
150		0



Survival Curve comparison		
1	Comparison of Survival Curves	
2		
3	Log-rank (Mantel-Cox) test	
4	Chi square	2.010
5	df	1
6	P value	0.1563
7	P value summary	ns
8	Are the survival curves sig different?	No
9		
10	Gehan-Breslow-Wilcoxon test	
11	Chi square	2.532
12	df	1
13	P value	0.1115
14	P value summary	ns
15	Are the survival curves sig different?	No
16		
17	Median survival	
18	Control	Undefined
19	Treated	75.00
20		
21	Hazard Ratio (Mantel-Haenszel)	A/B
22	Ratio (and its reciprocal)	B/A
23	95% CI of ratio	0.6907 to 1.448
24		0.6907 to 10.02
25	Hazard Ratio (logrank)	A/B
26	Ratio (and its reciprocal)	B/A
27	95% CI of ratio	0.1034 to 1.442
28		0.6933 to 9.668

NIAID

# Outline

## 1. Statistical Testing Process

- 1.1. Formulate Null and Alternative hypotheses
- 1.2. Choose the appropriate test statistics
- 1.3. Calculate Sample Size
- 1.4. Compute the probability of observing the test statistic under the null hypothesis
- 1.5. Make statistical decision and biological conclusion

## 2. Common Statistical Tests

- 2.1. Choose appropriate tests
- 2.2. Data Preprocessing
- 2.3. Correlation
- 2.4. Two-sample T-test
- 2.5. ANOVA and Multiple Comparisons
- 2.6. Survival Analysis

## 3. Application in Prism and R



NIAID

### How to learn R?

- [Datacamp - Introduction to R](#): Interactive course for R beginners

### How to learn Prism?

- [Prism 8 User Guide](#)
- [Prism 8 Statistics Guide](#)
- [Prism 8 Curve Fitting Guide](#)
- [Hands-on Introduction to Statistics and Visualization using Prism 8](#) by BCBB

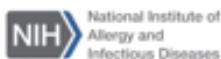


NIAID

## Conclusions

- Statistical Testing methods are much more than these!
- The materials of this seminar will be updated on: [Github - Statistical Testing](#)
- If you have any specific statistical problem, please send email to  
[bioinformatics@niaid.nih.gov](mailto:bioinformatics@niaid.nih.gov)
- Check our training schedule: [http://www.eventzilla.net/user/NIAID\\_OCICB\\_BCBB](http://www.eventzilla.net/user/NIAID_OCICB_BCBB)
- Take the survey and tell us what statistical topics you want to learn about next:  
[BCBB Statistical Training – Suggest a class!](#)
- Ask question on slack group!

*The invitation to the slack group will be sent after the workshop*



NIAID

**Thank You!**



NIAID