

City University of New York (CUNY)

CUNY Academic Works

Open Educational Resources

Queensborough Community College

2020

Clear-Sighted Statistics: Module 13: Introduction to Null Hypothesis Significance Testing (NHST) (slides)

Edward Volchok

CUNY Queensborough Community College

[How does access to this work benefit you? Let us know!](#)

More information about this work at: https://academicworks.cuny.edu/qb_oers/156

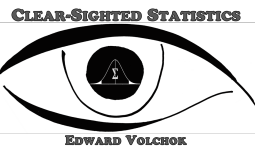
Discover additional works at: <https://academicworks.cuny.edu>

This work is made publicly available by the City University of New York (CUNY).

Contact: AcademicWorks@cuny.edu

Introduction to NHST (Null Hypothesis Significance Testing)

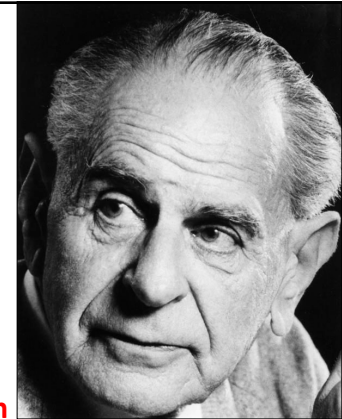
Module 13



© 0 1

0

“...I shall not require of a scientific system that it shall be capable of being singled out once and for all, in a positive sense; but I shall require that its logical form shall be such that it can be singled out, by means of empirical tests, in a negative sense: it must be possible for an empirical scientific system to be refuted by experience.”



Karl R. Popper
(1902 – 1994)

Falsification, nullification, refutation

© 0 1

Karl R. Popper, *The Logic of Scientific Discovery*, (Mansfield Centre, CT: Martino Publishing, 2014), pp. 40-41

1

Lecture Objectives

Define Null and Alternate Hypotheses

Define the significance level

Define Type I (α) errors

Define Type II (β) errors

Define statistical power

Describe the purpose of a decision rule

Define p-values

List the steps in the testing process

© 0 1

2

The Origins of NHST

© 0 1

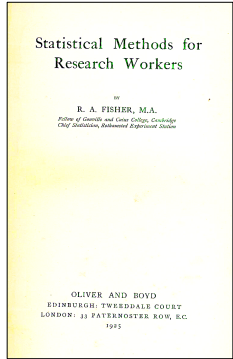
3

Clear-Sighted Statistics

1925: Fisher published the first edition of *Statistical Methods for Research Workers**

Laid the foundation for “significance tests”

Only one hypothesis
(Today NHST has 2 hypotheses)



*The 14th and final edition was published in 1970, 8 years after Fisher's death

4

Clear-Sighted Statistics

1928: Jerzy Neyman and Egon Pearson introduced Hypothesis Tests

Goal: Correct flaws in Fisher's approach

Added a 2nd hypothesis, Alternate Hypothesis (H_1)

Defined Type I or α Errors (false positive), and

Type II or β errors (false negative)

Jerzy Neyman, and Egon S. Pearson, E. S. "On the Use and Interpretation of Certain Test Criteria for Purposes of Statistical Inference: Part I." *Biometrika* 20A, 1928, pp. 175-240.

5

Clear-Sighted Statistics

Fisher and Neyman waged an acrimonious debate about the merits of their approaches until Fisher's death in 1962

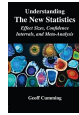
6

Clear-Sighted Statistics

Textbook authors combined Fisher's significance testing and Neyman-Pearson's hypothesis testing into a “unified” approach

7

Not everyone agrees that this amalgam is “unified”



Clear-Sighted Statistics

“It might be tempting to regard a mixture of the two approaches [Fisher and Neyman-Pearson] as possibly combing the best of both worlds, but the two frameworks are based on incompatible conceptions of probability. The mixture is indeed incoherent, and so it’s not surprising that misconceptions about NHST are so widespread.”
-- Geoff Cumming

© © © © Geoff Cumming, *Understanding the New Statistics: Effect Sizes, Confidence Intervals, and Meta-Analysis*. (New York: Routledge, 2012), p. 25.

8

What NHST does

Clear-Sighted Statistics

Uses sample data and probability theory to determine whether a proposition about population data should be rejected

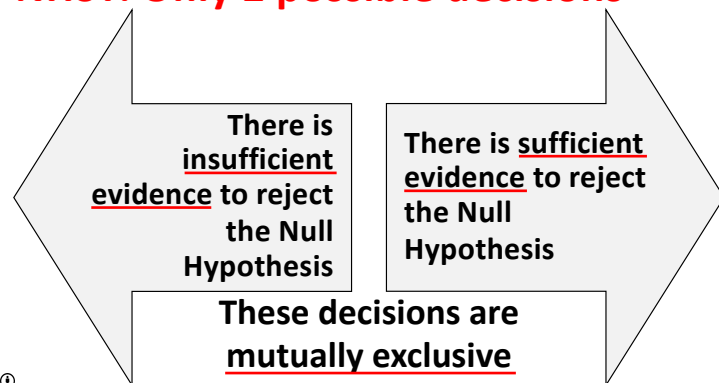
Does not verify a hypothesis

© © © ©

9

NHST: Only 2 possible decisions

Clear-Sighted Statistics



© © © ©

10

What an Hypothesis is not

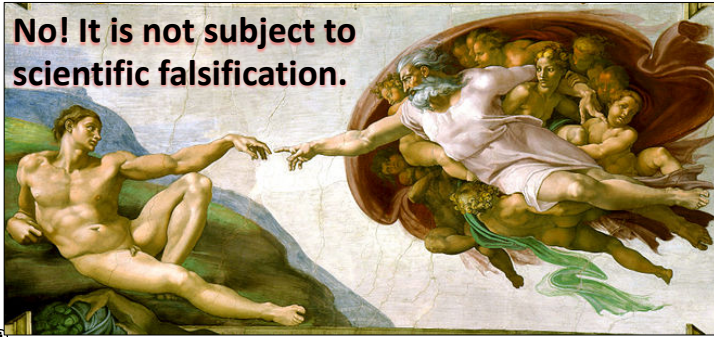
Clear-Sighted Statistics

© © © ©

11

Is Michelangelo's painting of God creating Adam an hypothesis?

No! It is not subject to scientific falsification.



12

A Hypothesis is not a Theory

13

What is an Hypothesis?

Preliminary proposition that provides an explanation of a phenomenon

Can be tested, refuted, falsified

14

What is a Theory?

Unified explanation of phenomena

A theory was initially only a hypothesis

Withstood repeated attempts at falsification

Yet, theories can also be falsified

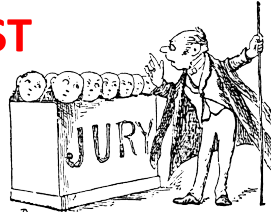
Paradigm Shift:
Scientific revolution occurs when a theory is falsified

15

Clear-Sighted Statistics

Non-Mathematical NHST

Criminal Jury Trial



16

16

Clear-Sighted Statistics

Trials & NHST have a shared premise

Defendant is assumed not guilty until the prosecution can convince a jury of the defendant's guilt beyond a reasonable doubt

Null Hypothesis is presumed an acceptable explanation until the data proves otherwise

17

17

Clear-Sighted Statistics

Null Hypothesis, H_0

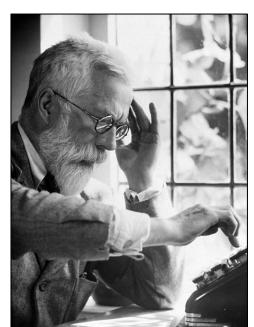
<p>With parametric techniques: H_0 always pertains to the population</p>	<p>Sample statistic (\bar{X}, p, s^2, or r) = population parameter (μ, π, σ^2, or ρ)</p>
<p>With a nonparametric technique like Chi-square: Observed Frequencies = Expected Frequencies</p>	<p>Implies "no difference" or "no effect"</p>

18

18

Clear-Sighted Statistics

H_0 never proven, but may be disproven



“...the null hypothesis is never proved or established, but is possibly disproved, in the course of experimentation. Every experiment may be said to exist only in order to give the facts a chance of disproving the null hypothesis.”*

Ronald A. Fisher
(1890 – 1962) *Ronald Aylmer Fisher, *The Design of Experiments*, (Edinburgh, UK: Oliver and Boyd, 1935), p. 19.

19

19

Hypotheses & Skepticism

The truth of any hypothesis is always subject to doubt

When the H_0 is not rejected, we never say it is proven

We say, "we failed to reject it"

When the H_0 is rejected, we do not consider the H_1 true



Alternate Hypothesis, H_1

This is the proposition the prosecutor hopes to prove: Defendant is guilty

Falsification of the H_0

Sample statistic (\bar{X} , p , s^2 , or r) \neq population parameter (μ , π , σ^2 , or ρ)

Observed Frequencies do not fit the Expected Frequencies



H_0 & H_1 are...

Mutually Exclusive

Collectively Exhaustive

Either/Or proposition



Significance Level, α

Jury trial verdicts are based on the "reasonable doubt" standard

NHST decisions are based on the significance level, α

α is the inverse of the confidence level ($1 - CL = \alpha$)

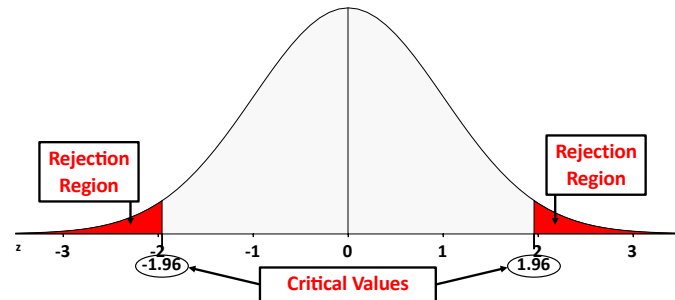
α delineates the rejection region from the region where the H_0 is not rejected

α is set at the start of the NHST

Neyman-Pearson & Fisher used a 5% α most frequently



Rejection regions & critical values



Two possible outcomes, two possible errors

Two Outcomes:
 1) Reject H_0 , 2) Fail to reject the H_0

Two Possible Errors:
 1) False Positive (Type I or α)
 2) False Negative (Type II or β)

Jury Trial: Correct Decision #1

- Innocent man declared not guilty
- H_0 is not rejected
- Any difference between the statistic and parameter is merely sampling error

Jury Trial: Correct Decision #2

- Guilty man is convicted
- H_0 is rejected
- Difference between the statistic and parameter is statistically significant (not sampling error)

Type I Error: False Positive

- Innocent man is convicted
- H_0 is wrongfully rejected
- α sets the acceptable risk of a Type I error

Type II Error: False Negative

- Guilty man is acquitted
- Wrongfully failing to reject the H_0

2 x 2 NHST Outcomes Matrix

	Accept Null Hypothesis	Reject Null Hypothesis
Correct Decision	Jury correctly acquits the defendant. The analyst correctly fails to reject the null hypothesis.	Jury correctly convicts the defendant. The analyst correctly rejects the null hypothesis.
Incorrect Decision	Type I Error, α error, or false positive. An innocent man is convicted. The analyst wrongly rejects the null hypothesis.	Type II Error, β error, or false negative. A guilty man is acquitted. The analyst wrongly fails to reject the null hypothesis.

A philosopher on Type I & Type II Errors



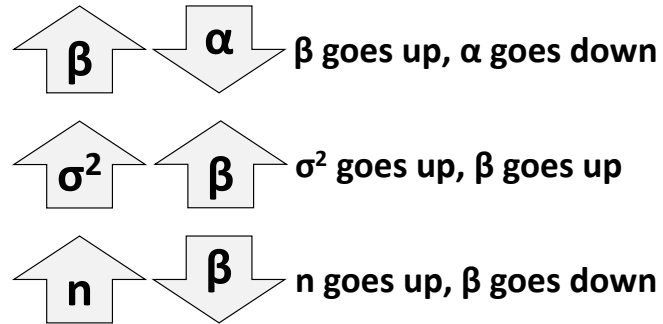
Type I Error
Reject "True" H_0

"...one can be deceived in believing what is untrue, but on the other hand, one is also deceived in not believing what is true."

Type II Error
Retain "False" H_0

*Kierkegaard, Søren 1962. Works of Love: Some Christian Reflections in the Form of Discourse. Translated by Howard V. and Edda H. Hong. New York, NY: Harper Torchbooks, p. 23.

What affects the P(α) & P(β)?



Calculating Type I & Type II Errors

The allowable risk of a Type I error is set when the significance level is chosen

The probability of a Type II error is a calculation that will be explored

Goal for many analysts: $P(\alpha) \leq 0.05$ and the $P(\beta) \leq 0.20$
(Type I Errors considered more serious than Type II Errors)

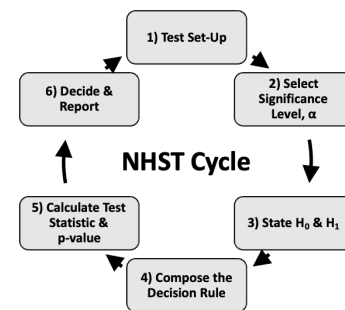
Statistical Power: (1 - β)

Power of a test to find a significant difference between the parameter and the statistic when one exists

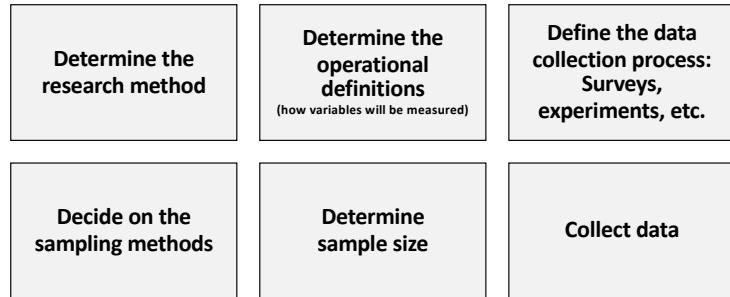
Minimum level of statistical power is 0.80 [$P(\beta) = 0.20$]

Low powered tests are unreliable: Miss important effect and confuse sampling error for an effect

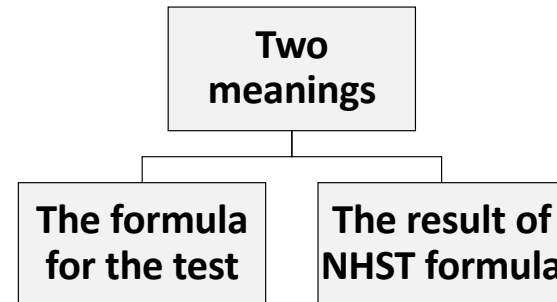
6 Steps of the NHST Cycle



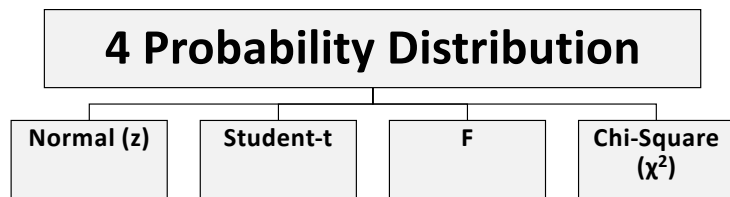
Step 1: Test Set-Up



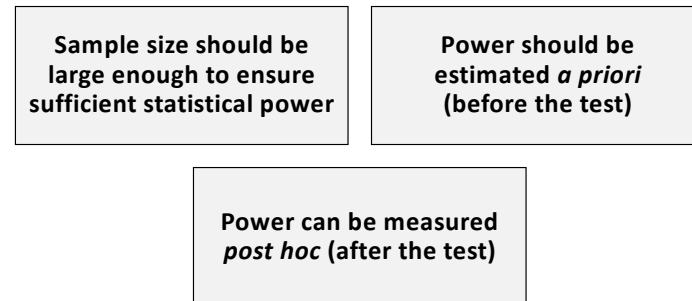
What's a test statistic?



Test statistics: 4 probability distributions



Sample size is a big issue



Significance Levels, α

Acceptable risk of committing a Type I Error

Chosen before data is analyzed

Typically set at 0.05, but sometimes 0.01 or 0.10

The lower the α , the harder to reject H_0

The lower the α , the less statistical power

Reject H_0 when the probability of obtaining a test statistic is \leq to α

α & Critical Value(s)

The critical value or values of the test statistic that mark the boundary between the region where the H_0 is rejected or not rejected

Critical values are based on the α and the test statistic

Interpreting Statistical Significance

When the probability of the test statistic is $\leq \alpha$, the results are “statistically significant”

Reject the H_0

Results unlikely due to random sampling error

Statistical Significance \neq Practical Significance

Practical significance refers to the magnitude of the effect

Having statistical significance does not mean results have any real application (practical significance)

With a large enough n , all results have statistical significance, but may lack practical significance

Practical significance without statistical significance is possible

Null Hypothesis (H_0)

Hypothesis that researchers seek to “nullify”

No difference between the statistic and parameter

No effect

Takes an equal sign: =, ≤, or ≥

Alternate Hypothesis (H_1)

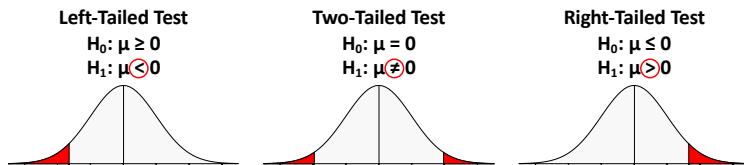
Opposite of the H_0

Something is happening:
Statistic \neq Parameter or there is an effect

Takes a not equal sign: \neq , $<$, or $>$

Does not imply practical significance

z and t tests are directional



Note the direction of the sign in the H_1

Decision Rules

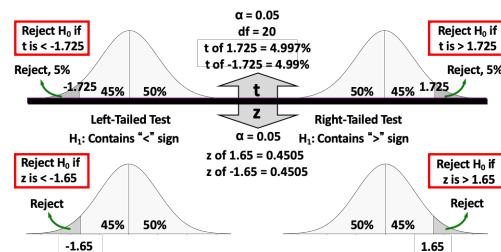
States the criterion for rejecting H_0

Based on critical value(s)

Decision Rule Syntax

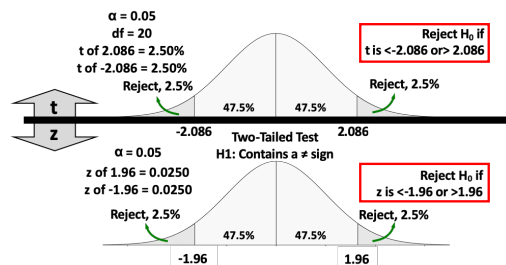
Reject the Null Hypothesis if
(z, t, F, χ^2) is (<, >, or < and >)
the critical value(s)

One-Tailed Tests



Note the decision rules

Two-Tailed Tests



Note the decision rules

Critical Values for z-tests

α	Left-Tailed Test	Two-Tailed Test	Right-Tailed Test
0.01	-2.33 (-2.326)	-2.58 & 2.58 (-2.576 (2.576)	2.33 (2.326)
0.05	-1.65 (-1.645)	-1.96 & 1.96 (-1.960 & 1.960)	1.65 (1.645)
0.10	-1.28 (-1.283)	-1.65 & 1.65 (-1.645 & 1.645)	1.28 (1.283)

Values within the parentheses are used by Excel

Calculating the Test Statistic

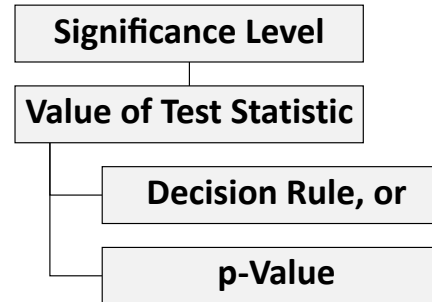
Formula varies with the type of test

Typically a fraction

Numerator: Sampling error

Denominator: Standard Error of the Mean or Proportion

Making a Decision



What are p-Values?

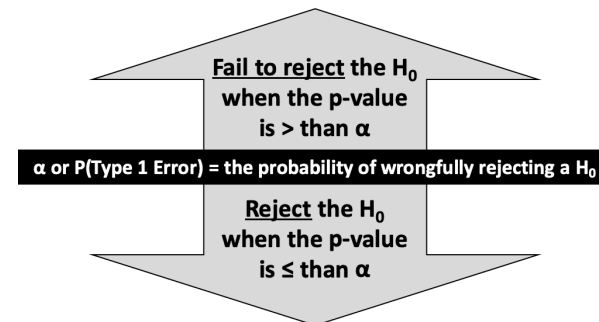
Slippery concept



Measures compatibility of the data with the H_0

p-value: Probability of getting a test statistics that is as extreme or more extreme than the one you just calculated

Interpreting p-Values



p-Values & the decision to reject the H_0

Possible scenarios: $\alpha = 0.05$ and the p-value is 0.05 or <0.001

H_0 is rejected in both scenarios

We would have more confidence that the results are statistically significant when $p < 0.001$ than $p = 0.05$

**P-values are widely
misunderstood**

**March 2016: American Statistical
Association's statement on p-values**

ASA on p-Values

"The contents of the ASA statement and the reasoning behind it are not new—statisticians and other scientists have been writing on the topic for decades. But this is the first time that the community of statisticians, as represented by the ASA Board of Directors, has issued a statement to address these issues."

-- Jessica Utts, President, ASA

6 Principles

#1: P-values can indicate how incompatible the [sample] data are with a specified statistical model [or the null hypothesis].

#2: P-values do not measure the probability that the studied [null] hypothesis is true, or the probability that the data were produced by random chance alone.

#3: Scientific conclusions and business or policy decisions should not be based only on whether a p-value passes a specific threshold [significance level].

6 Principles (Continued)

#4: Proper inference requires full reporting and transparency.

#5: A p-value, or statistical significance, does not measure the size of an effect or the importance of a result.

#6: By itself, a p-value does not provide a good measure of evidence regarding a model or hypothesis.

What's Next?

Test of Means, μ	Tests of Proportions, π	Test of Variance, σ^2	Tests of Relationships
<ul style="list-style-type: none"> • z-Tests • t-Tests • ANOVA Tests 	<ul style="list-style-type: none"> • z-Tests 	<ul style="list-style-type: none"> • F-test for Equality σ^2 	<ul style="list-style-type: none"> • Chi-Square • Correlation & Regression



Except where otherwise noted *Clear-Sighted Statistics* is licensed under a Creative Commons License. You are free to share derivatives of this work for non-commercial purposes only. Please attribute this work to Edward Volchok.

