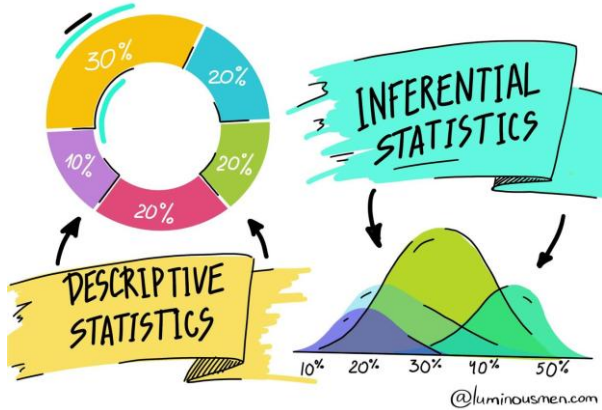
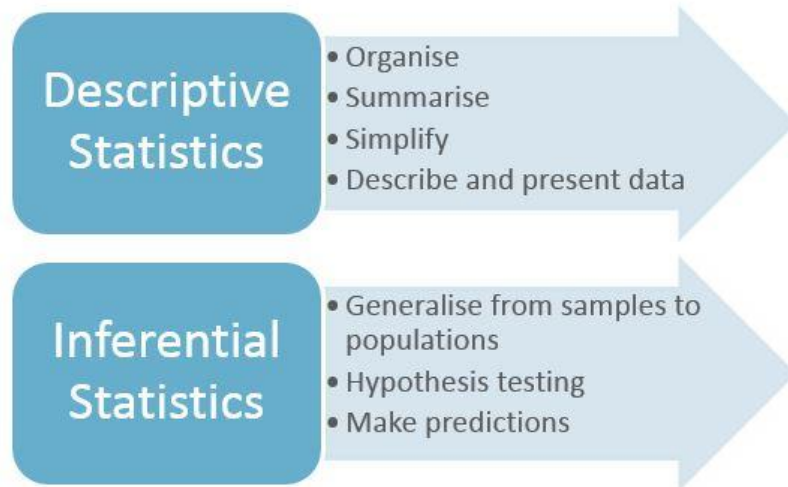


Hypothesis Testing

Why do we need inferential statistics?



Inferential statistics allow us to *infer* or generalize observations made with samples to the larger population from which they were selected.



What is a Hypothesis?

Research Question
(ideas)



A specific testable statement
(that guides an experiment and
statistical analysis)

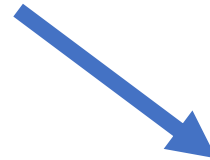
What Is a Real Hypothesis?

- A hypothesis is an educated guess, based on observation.
- Usually, a hypothesis can be supported or refuted through experimentation or more observation.
- A hypothesis can be disproven, but not proven to be true.



Research Question – Is online teaching effective?

Hypothesis Statement – Students taught offline perform better than students taught online



(ASSUMPTION) – based on previous studies, observations, experiences, etc.

Null Hypothesis and Alternative Hypothesis

$$H_0 \text{ vs } \begin{matrix} H_1 \\ \text{or} \\ H_a \end{matrix}$$

Students taught online vs offline
perform equally well on exams (no
difference/null)

Students taught offline perform
better than students taught online

Students taught online perform
better than students taught offline

You perform experiments to check if the H_0 holds true or not.
By disproving H_0 you accept the H_A

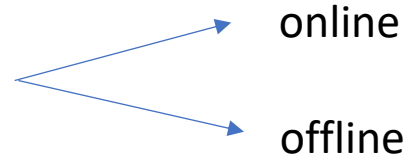
Variables in a hypothesis

Hypothesis Statement – Students taught offline perform better than students taught online

INDEPENDENT
VARIABLE

not changed by the other variables you are trying to measure

Teaching method



online

offline

DEPENDENT
VARIABLE

Value is changed or affected by the independent variable/s

Exam performance

Individuals with more years of education have higher income

Ho – No relationship between years of education and income

H₁ - Individuals with more years of education have higher income

Leopards are stronger than Tigers

H₀ – Leopards and Tigers are equally strong, no difference

H₁ – Tigers are stronger than Leopards

H₂ – Leopards are stronger than Tigers



Exercise effects on anxiety

H₀ - Exercise has no effect on anxiety

H₁ - Exercise lowers anxiety

H₂ – Exercise increases anxiety

IV – Exercise (exercising, not exercising)

DV – anxiety levels

Directionality in a hypothesis



```
graph TD; A[Directionality in a hypothesis] --> B[This prediction is typically based on past research, accepted theory, extensive experience, or literature on the topic.]; B --> C[Else your statistical outcome can be misleading, by ignoring other outcomes e.g. Does a technical degree impart technical skills?];
```

This prediction is typically based on past research, accepted theory, extensive experience, or literature on the topic.

Else your statistical outcome can be misleading, by ignoring other outcomes
e.g. Does a technical degree impart technical skills?

High quality of engineering education leads to higher technical skills

Ho – Quality of engineering education has no effect on technical skills

H1 - High quality of engineering education leads to higher technical skills



Directionality?

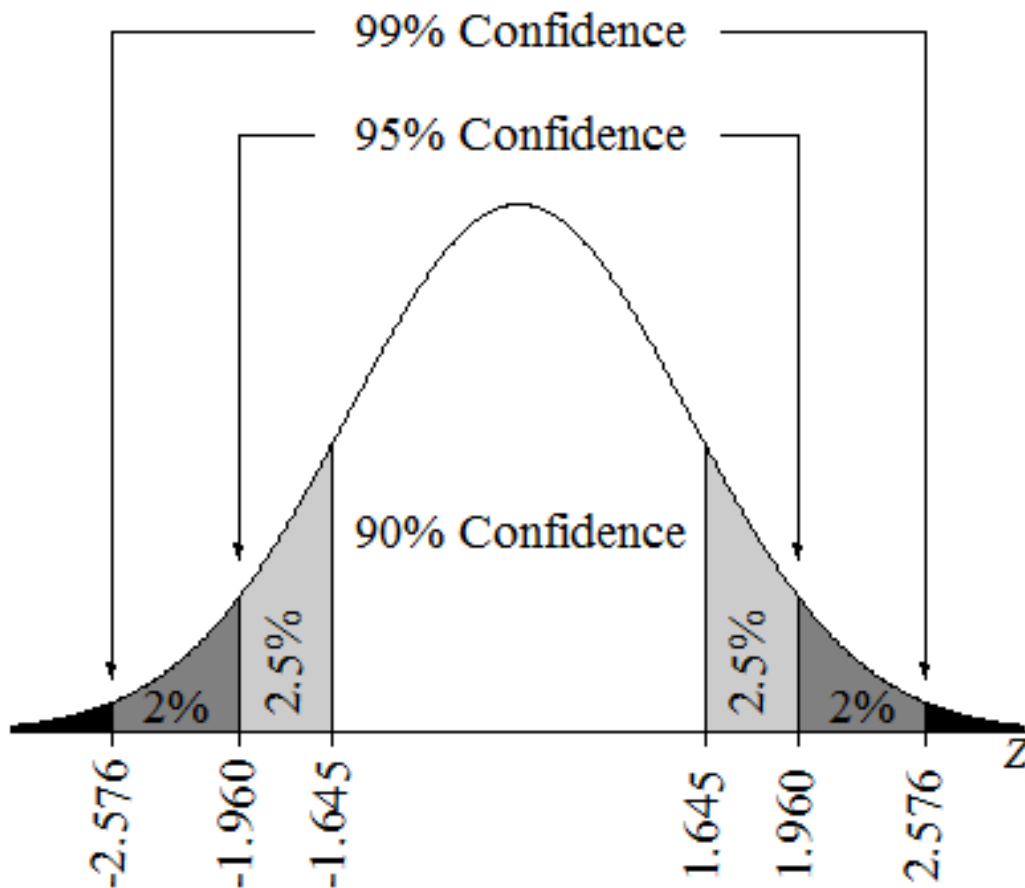
Exercise effects on anxiety

H_0 - Exercise has no effect on anxiety

H_1 - Exercise lowers anxiety

H_2 - Exercise increases anxiety

Confidence Intervals



Confidence Level	α (level of significance)	$Z_{\alpha/2}$
99%	1%	2.575
95%	5%	1.96
90%	10%	1.645

$$CI = \bar{x} \pm z \frac{s}{\sqrt{n}}$$

CI = confidence interval

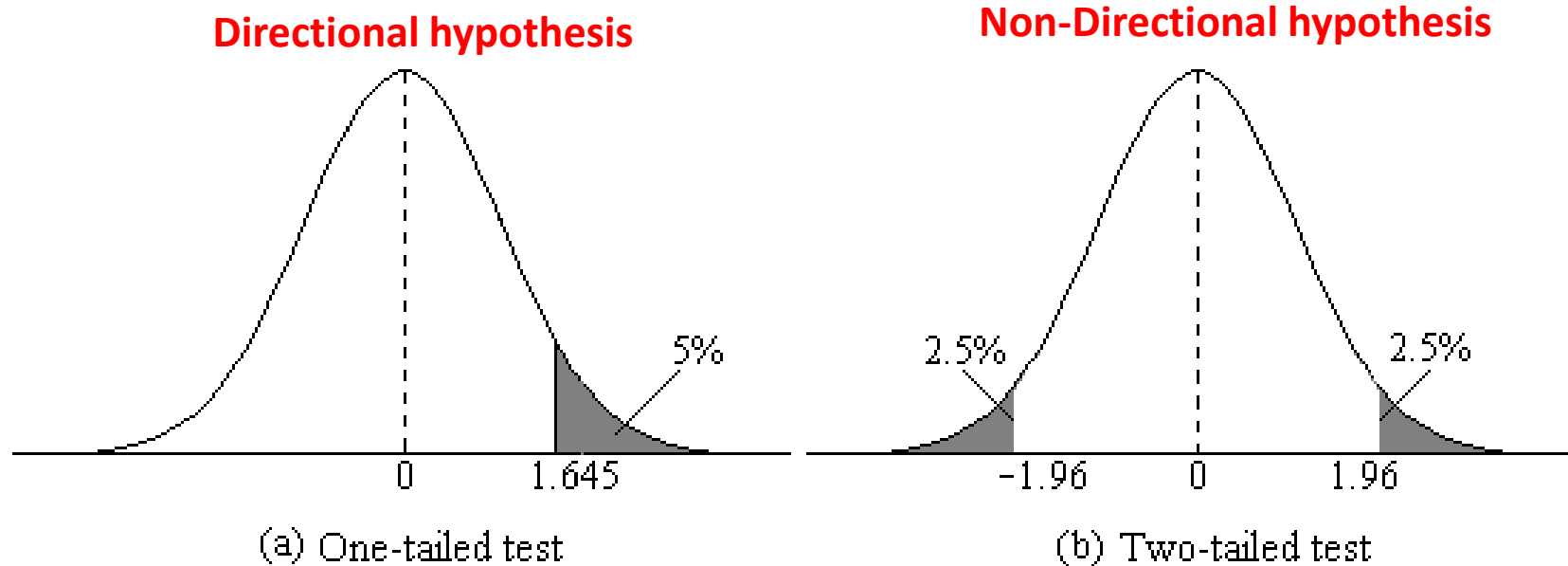
\bar{x} = sample mean

z = confidence level value

s = sample standard deviation

n = sample size

Hypothesis testing



General Rule: Use two-tailed test.

Only if direction is known from prior studies (justified reason), use one-tail test.

Criterion (α) for significance – 5% (0.05) for most behavioural studies (95 % CI)

If $p > 0.05 \rightarrow$ Accept the H_0

If $p \leq 0.05 \rightarrow$ Reject the H_0 & accept H_A

One-tailed vs two-tailed test

When is a one-tailed test NOT appropriate?

- 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, no matter how "close" to significant the two-tailed test was.
- Using statistical tests inappropriately can lead to invalid results that are not replicable and highly questionable—a steep price to pay to show significance in your results

Exercise effects on anxiety

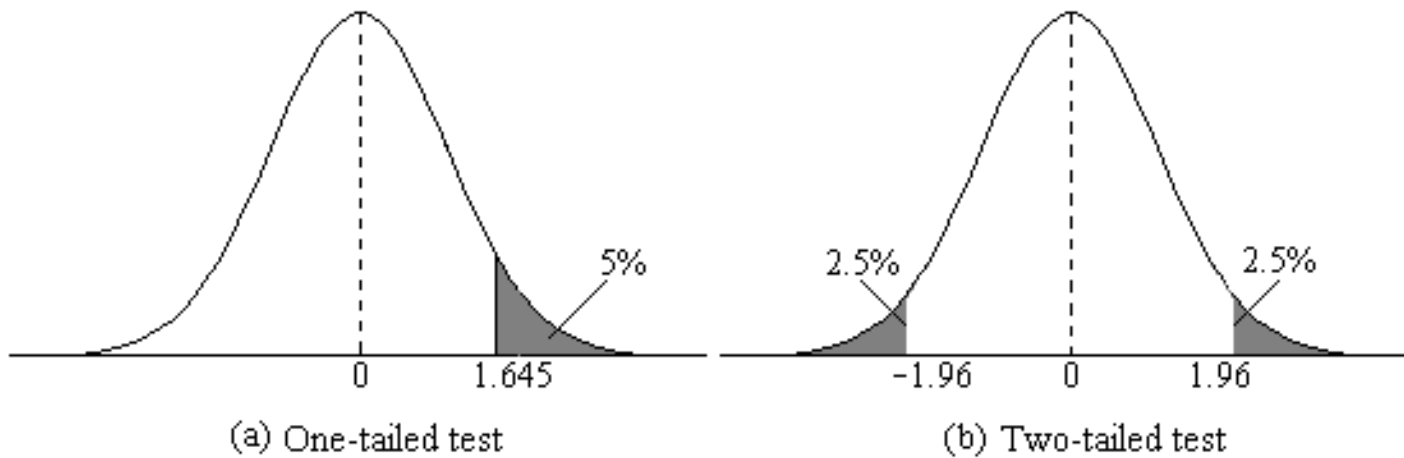
H_0 - Exercise has no effect on anxiety

H_1 - Exercise lowers anxiety

H_2 - Exercise increases anxiety?

EXERCISE AND **ANXIETY**

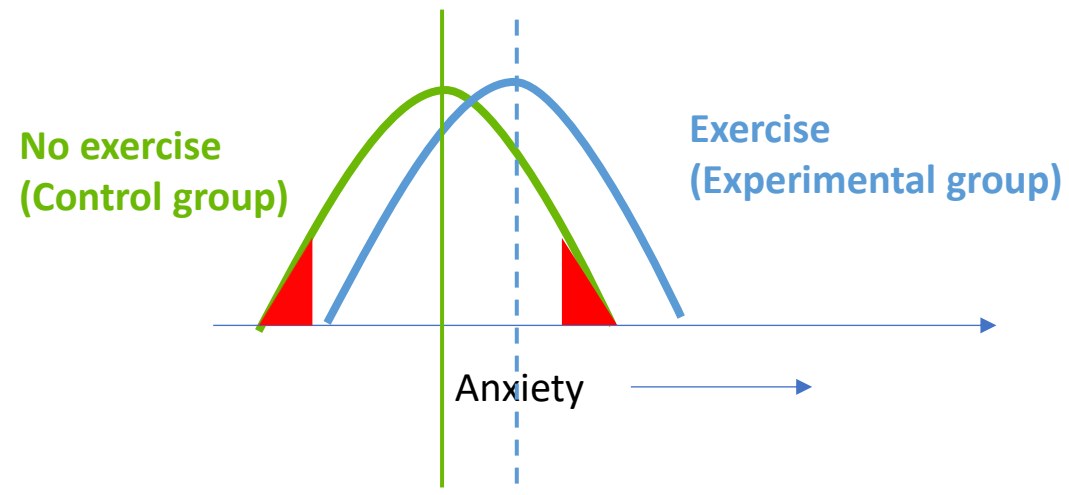
Studies show that it is very effective at enhancing overall cognitive function.



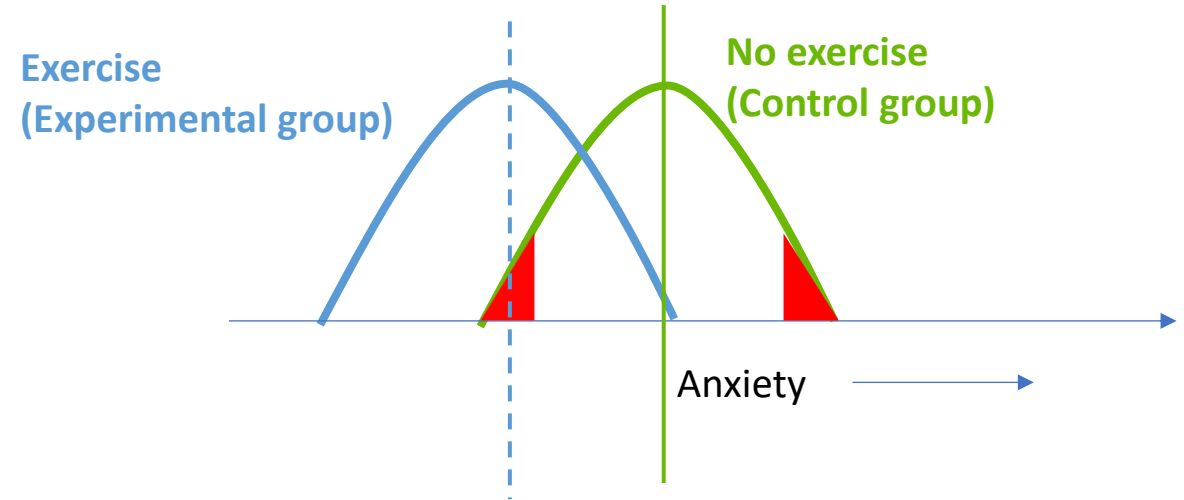
When $p \leq .05$, we reject the null hypothesis - there is a '**significant**' difference between the two groups.

When $p > .05$, we retain the null hypothesis - there is less difference between the groups.

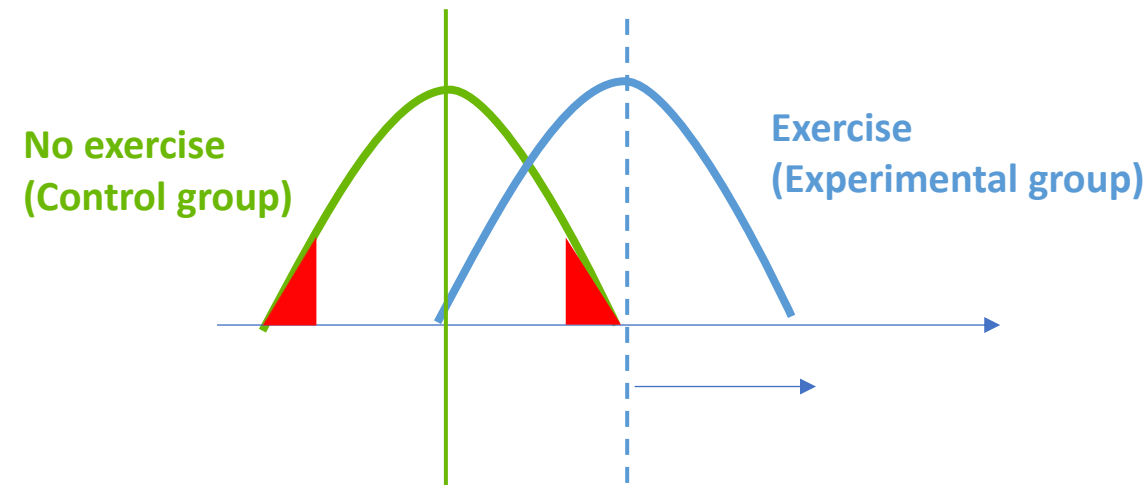
H0 - Exercise has no effect on anxiety



H1 - Exercise lowers anxiety



H2 - Exercise increases anxiety



Another Directional Hypothesis

You have a new drug to treat pain that is cheaper than the existing drug and you only want to confirm if the new drug is less effective than the existing drug

Whether the new drug is better than the existing drug does not matter.

H₀ - Null hypothesis – No difference between new drug and existing drug to treat pain

H₁ - Alternate hypothesis – Is the new drug less effective than the existing drug

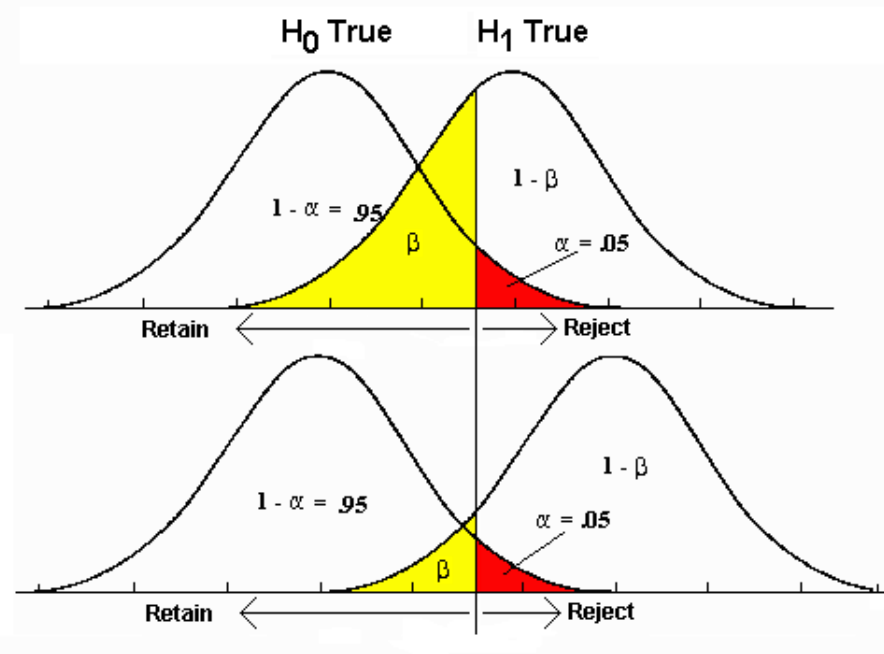
EXAMPLES

- H_0 , H_A , IV, DV, one or two tailed test?
- Smoking is injurious to the lungs
- Videogaming can lower attention span
- Does repetition in advertising improve sales?
- Air pollution is more fatal than COVID19
- Is there a difference in leadership style between men and women?

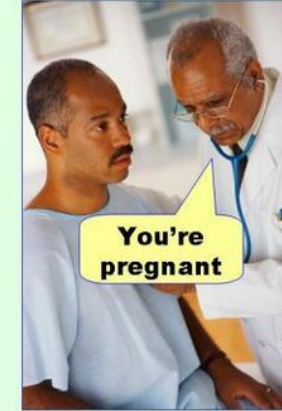
Ever wondered how and why statistics came into picture?

<https://nautil.us/how-eugenics-shaped-statistics-238014/>

Types of Errors in hypothesis testing



Type I error
(false positive)



Type II error
(false negative)



[Reality:]

Ho False

Ho True

Decision from
statistical tests

Reject Ho

Accept Ho

Correct Decision
Sensitivity/Power
 $1 - \beta$

Type 1 Error
"False Positive"
 α

Type 2 Error
"False Negative"
 β

Correct Decision
Specificity
 $1 - \alpha$

Observe difference
when none exists

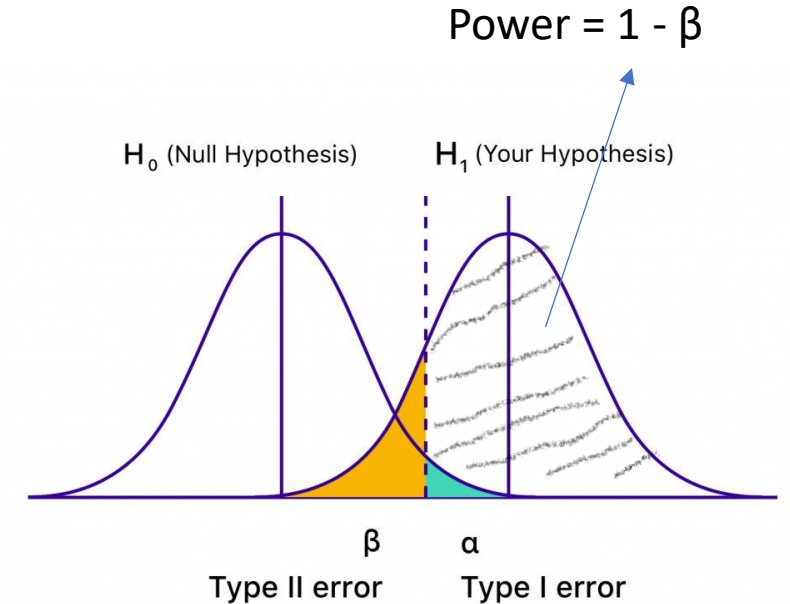
Overreacting!

Fail to find a difference
when there is one
Underreacting!

- Sample size it too small (high variability)
- Choosing one-tailed instead of two-tailed test
- Wrong statistical test

Power

- Power - the probability that your test will find a statistically significant difference when such a difference actually exists.
- In other words, power is the probability that you will reject the null hypothesis when you should (and thus avoid a Type II error).
- It is generally accepted that power should be .8 or greater; that is, you should have an 80% or greater chance of finding a statistically significant difference when there is one.



Power

Power is calculated using statistical software. You need to know –

- What type of test you plan to use (e.g., independent t-test, paired t-test, ANOVA, correlation, regression, etc.)
- The alpha value or significance level you are using (usually 0.05 or 0.01)
- The expected effect size
- The sample size you are planning to use

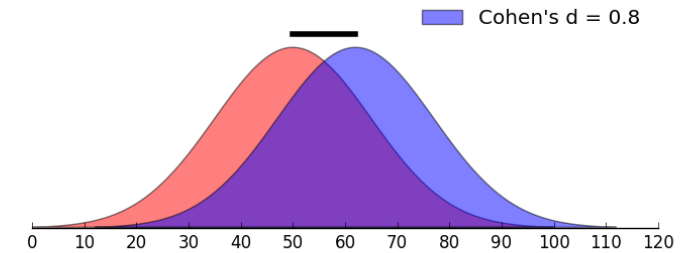
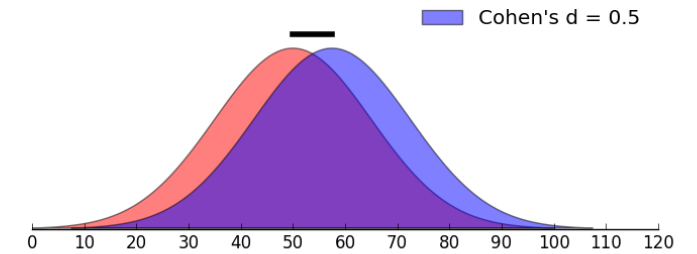
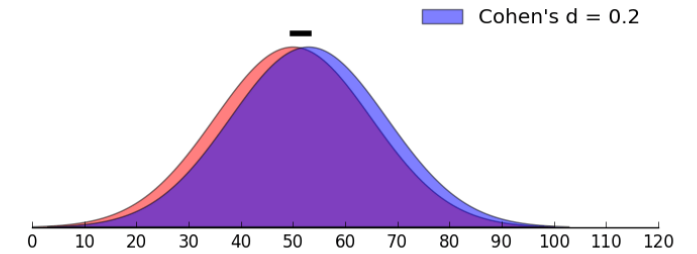
As your sample size increases, so does the power of your test.

- optimal sample means that you have collected more information -- which makes it easier to correctly reject the null hypothesis when you should.
- A power value is between 0 and 1.
- If the power is less than 0.8, you typically need to increase your sample size.

Effect Size

E.g. you evaluate the effect of a group discussion on student knowledge using pre and post tests on 1000 students. The mean score on the pre test was 83 out of 100 while the mean score on the post test was 84.

- What if you simply found a statistical difference by virtue of a large sample size (> 1000 or 10000)?
- If you calculate the effect size – you get a standard method to defining the importance of the statistical difference



Cohen's d effect size interpretation

< 0.1 = trivial effect

$0.1 - 0.3$ = small effect

$0.3 - 0.5$ = moderate effect

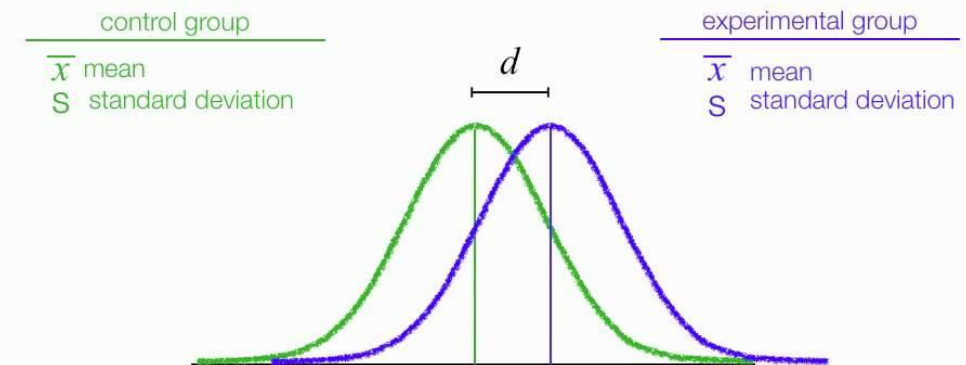
> 0.5 = large difference effect

Effect Size

- Effect size is a quantitative measure of the *strength of a phenomenon*.
- Effect size emphasizes the **size** of the difference or relationship
- Examples:
 - the correlation between two variables (specifically r^2)
 - $r=.1$ weak, $r=.5$ moderate, $r=.7$ strong, $r=.9$ very strong
 - the regression coefficient in a regression (B_0, B_1, B_2)
 - Relative to model and field
 - the mean differences in t tests (use Cohen's D)
 - $d = .2$ is small; $r = .5$ is medium; $r = .8$ is large
 - The mean differences in ANOVA (use eta)
 - .01 is small, .06 medium, .14 large

$$\text{Cohen's Effect size} = \frac{(\text{Mean}_{\text{treatment}} - \text{Mean}_{\text{control}})}{\text{Standard deviation}_{\text{pooled}}}$$

$$d = \frac{\bar{x} - \bar{x}}{S}$$



Basic formula for sample size - Continuous data

$$\text{Number of samples per group (n)} = \frac{2 \times (Z_{(1-\alpha/2)} + Z_{\beta})^2 \times \sigma^2}{\Delta^2}$$

Where Δ = size of difference, minimal effect of interest

α = significance level (eg 0.05)

β = power, probability of detecting a significant result (typically 80%, 90%)

σ = SD of data

Z_p = points on normal distribution to give required power and significance

DV: Anxiety level

Do people who exercise have lower levels of anxiety?

Does exercise lower anxiety?

IV: Exercise

Experimental group

Exercise

Anxiety level

Control group

No Exercise

Anxiety level

Between groups
(this does not allow you to measure change)

Experimental condition

Anxiety level

Exercise

Anxiety level

Within group/Repeated measures
(crossover design)

- Participant fatigue
- Longer experimental duration
- Carry over effects

Experimental group

Anxiety level

Exercise

Anxiety level

Anxiety level

No Exercise

Anxiety level

Control group

Mixed design

Between groups & Within groups