# Decision Error

Grinnell College

Nov 15, 2024

# Rest of Semester

The rest of the tools we will learn about in class invovles testing for *differences* or *associations* between groups which may help inform your project goals

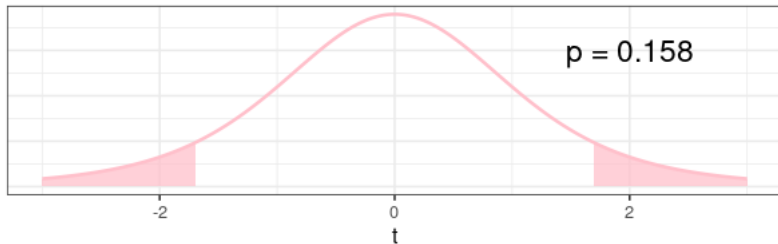| Type | Continuous | Categorical |
|---|---|---|
| Simple Test | $t$-test | Single Proportion |
| 2 Groups | Two-sample $t$-test, paired test | Difference in Proportion |
| Multiple Groups | ANOVA | $\chi^2$ Test |
| Mixed variables | Regression | Regression |

# Strength of Evidence

So far, our process has been as follows:

1. Being with a null hypothesis, $H_0 : \mu = \mu_0$
2. Collected data and comptute statistic, i.e, $\overline{x}$
3. Compare our statistic against the null distribution, i.e., $t = \frac{\overline{x} - \mu_0}{\hat{\sigma}/\sqrt{n}}$
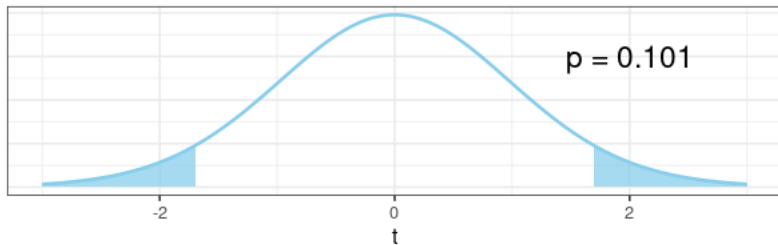4. Derive a $p$-value based on the statistic and the distribution

Today, we will address the question of whether or not our oberved data is consistent with our hypothesis

# Comparison

t = 1.69, df = 5



p = 0.158

t = 1.69, df = 30



p = 0.101

# Decision Making

Based on the evidence we have collected, we must ultimately decide between one of two decisions:

1. There is sufficient evidence to reject $H_0$
2. There is *not* sufficient evidence to reject $H_0$

# Decision Making

Just as our confidence intervals were correct or incorrect, so to may be our decision regarding $H_0$. In this case, however, there are two distinct ways in which our decision can be incorrect:

1. $H_0$ is *TRUE* (i.e., there is no effect), yet we reject anyway
2. $H_0$ is *FALSE* (i.e., there is an effect), yet we fail to reject it

# Decision Making

These two types of errors are known as Type I and Type II errors, respectively:

1. $H_0$ is *TRUE* (i.e., there is no effect), yet we reject anyway
   - Type I error
   - "False positive"
   - Evidence leads to wrong conclusion
2. $H_0$ is *FALSE* (i.e., there is an effect), yet we fail to reject it
   - Type II error
   - "False negative"
   - Not enough evidence to conclude

# Decision Making

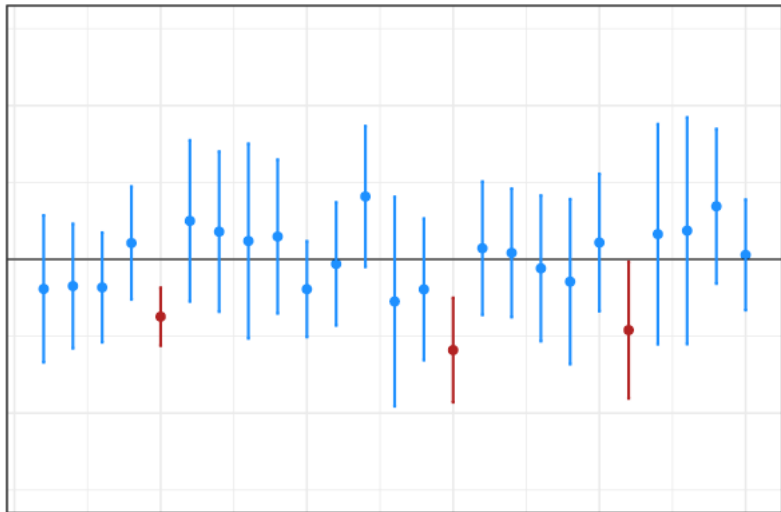|  | True State of Nature | |
|---|---|---|
| Test Result | $H_0$ True | $H_0$ False |
| Fail to reject $H_0$ | Correct | Type II Error |
| Reject $H_0$ | Type I Error | Correct |

# Type I Errors

A Type I error describes a situation in which we incorrectly identify a null effect:

▶ Conclude that an intervention works when it does not
▶ Conclude that there is a relationship between two variables when there are not

A Type I error will occur, for example, when our constructed confidence does not contain $\mu_0$ when $\mu_0 = \mu$

# Type I Errors

# Type I Error Rate

We can control the rate at which we commit Type I errors with adjusting the *level of significance*, denoted $\alpha$.

This is also called the *Type I error rate*

The Type I error rate has a *one-to-one* correspondence with our confidence intervals: a 95% confidence interval will permit a Type I error 5% of the time, corresponding to $\alpha = 0.05$

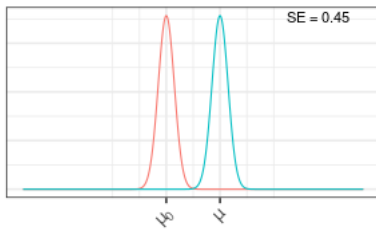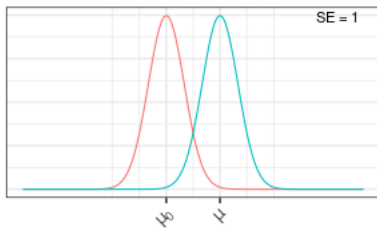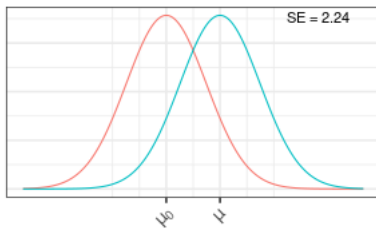We *reject* our null hypothesis when *p*-value $< \alpha$

# Type II Errors

A Type II error describes a situation in which the null hypothesis is false, yet based on the evidence gathered we fail to reject it:

- ▶ An intervention has a clinical effect, but it is not detected
- ▶ An email is considered spam, but the filter does not detect it

Typically, a Type II error is the result of one or more factors:

- ▶ Too few observations in our sample
- ▶ The population has large variability
- ▶ The effect size is small

Line — Null — True

# Type II Error Rate

The Type II error rate is typically denoted $\beta$

More frequently, we consider the rate at which Type II errors do not occur $(1 - \beta)$, a term we refer to as *power*

A study that is unable to detect a true effect is said to be *underpowered*

# Power

Consider the following analogy[1]: you send a child into the basement to find an object

- ▶ What is the probability that she actually finds it?
- ▶ This will depend on three things:
  - ▶ How long does she spend looking?
  - ▶ How big is the object she is looking for?
  - ▶ How messy is the basement?

---

# Power

If the child spends a long time looking for a large object in a clean, organized basement, she will most likely find what she's looking for

If a child spend a short amount of time looking for a small object in a messy, chaotic basement, it's probably that she won't find it

Each of these has a statistical analog:

- ▶ How long she spends looking? = How big is the sample size?
- ▶ How big is the object? = How large is the effect size?
- ▶ How messy is the basement? = How noisy/variable is the data?

# Drawing Conclusions

As we never truly know whether $H_0$ is correct or not, we must simultaneously be prepared to combat both types of error

| Test Result | True State of Nature | |
|---|---|---|
| | $H_0$ True | $H_0$ False |
| Fail to reject $H_0$ | Correct $(1 - \alpha)$ | Type II Error $(\beta)$ |
| Reject $H_0$ | Type I Error $(\alpha)$ | Correct $(1 - \beta)$ |

- Type I error $= P(\text{Reject } H_0 | H_0 \text{ true}) = \text{false alarm}$
- Type II error $= P(\text{Fail to reject } H_0 | H_A \text{ true}) = \text{missed opportunity}$

# *p*-values

Although the $\alpha = 0.05$ is customary for Type I error rate and a cut-off for "statistical significance", this is no substitute for correctly evaluating context

For example, a highly publicized study in 2009 involving a vaccine protecting against HIV found that, analyzed one way, the data suggested a *p*-value of 0.08. Computed a different way, it resulted in a *p*-value of 0.04

Debate and controversy ensued, primarily because the consequence of using a particular method was the difference between a result being on other side of the $p < \alpha$ threshold

But is there really that much a difference between $p = 0.04$ and $p = 0.08$?

# Multiple Testing

Consider conducting 2 hypothesis tests, each with a Type I error rate of 5%

For any given test, the probability of *not* making an error is

$$P(\text{No type I error}) = 0.95$$

1. What is the probability that neither test has a Type I error?
2. What is the probability that *at least* one test has a Type I error?

# Example

Suppose that I am interested in testing if there is a non-zero correlation between cost and average faculty salary in each of the 8 regions of our college dataset

Suppose further we are testing for significance at the level $\alpha = 0.05$

|   | Region | $p$-value |
|---|--------|-----------|
| 1 | Far West | 0.7667 |
| 2 | Great Lakes | 0.0085 |
| 3 | Mid East | 0.0001 |
| 4 | New England | 0.0061 |
| 5 | Plains | 0.9487 |
| 6 | Rocky Mountains | 0.7394 |
| 7 | South East | 0.0143 |
| 8 | South West | 0.0344 |

# Example

Suppose that I am interested in testing if there is a non-zero correlation between cost and average faculty salary in each of the 8 regions of our college dataset

If my Type I error rate for each test is 5%, what is the probability that I make at least one Type I error?

$$P(\text{At least one Type I error}) = 1 - P(\text{Probability of no Type I errors})$$
$$= 1 - (1 - 0.05)^8$$
$$= 33.6\%$$

That is, instead of making a Type I error 1 in 20 times, we are now making it 1 in 3 times

# Family-wise error rates (FWER)

For a collection of independent hypothesis tests, the **family-wise error rate (FWER)** describes the probability of making one or more Type I errors

For $m$ independent tests with a Type I error rate of $\alpha$, the FWER is defined as

$$\text{FWER} = 1 - (1 - \alpha)^m$$

# FWER Correction

Just as we control the Type I error rate of a single hypothesis test with $\alpha$, we also have an interest in controlling the FWER

For $m$ hypothesis tests controlled at level $\alpha$, the correction $\alpha^* = \alpha/m$ is known as the **Bonferonni Adjustment**

If instead for a series of $m$ tests we reject the null hypothesis when $p < \alpha^*$, we will control the FWER at level $\alpha$

Assuming the 8 regions of our hypothesis test are independent, our Bonferonni adjustment for $\alpha = 0.05$ should be

$$\alpha^* = 0.05/8 = 0.00625$$

| | Testing $p < \alpha$ | |
|---|---|---|
| | Region | $p$-value |
| 1 | Far West | 0.7667 |
| 2 | Great Lakes | 0.0085 |
| 3 | Mid East | 0.0001 |
| 4 | New England | 0.0061 |
| 5 | Plains | 0.9487 |
| 6 | Rocky Mountains | 0.7394 |
| 7 | South East | 0.0143 |
| 8 | South West | 0.0344 |

| | Testing $p < \alpha^*$ | |
|---|---|---|
| | Region | $p$-value |
| 1 | Far West | 0.7667 |
| 2 | Great Lakes | 0.0085 |
| 3 | Mid East | 0.0001 |
| 4 | New England | 0.0061 |
| 5 | Plains | 0.9487 |
| 6 | Rocky Mountains | 0.7394 |
| 7 | South East | 0.0143 |
| 8 | South West | 0.0344 |

# Review

Based on the evidence observed, we will ultimately make one of two decisions:

1. Reject $H_0$
2. Fail to reject $H_0$

Depending on the true state of $H_0$, we can be incorrect in two ways:

1. Type I Error ($\alpha$): $H_0$ is true, yet we reject anyway
2. Type II Error ($\beta$): $H_0$ is false, yet we fail to reject it

Finally, there is the issue of *multiple testing*

1. Family-wise error rate
2. Bonferonni correction