

## Three More Inferences

Grinnell College

November 2025

# Big Picture

The last few slide decks we have been focusing on testing if a population mean is the same as BLANK or building confidence intervals

We are now going to *generalize* the tests to....

1. The difference between two means
2. One proportion
3. The difference between two proportions

Critically, they work almost the exact same as a z-test stuff....

## 2 Inferences: Hypothesis Test

Almost all hypothesis tests carry the same format...

- ▶ State our null hypothesis (default state of the world we want to disprove) and alternative (thing we are trying to show)
- ▶ After checking assumptions we have a sampling distribution.
  - ▶ For us it'll be normal so long as we know  $\sigma^2$

## Difference of Two Means

Set up: We have two distinct subpopulations and are interested in the difference of their means. Eg are Dream Island penguin's mean flipper length the same as Toregersen Island penguin's flipper length?

---

$$H_0: \mu_1 - \mu_2 = 0 \quad \text{OR} \quad \mu_1 - \mu_2 \geq 0 \quad \text{OR} \quad \mu_1 - \mu_2 \leq 0$$

$$H_A: \mu_1 - \mu_2 \neq 0 \quad \text{OR} \quad \mu_1 - \mu_2 < 0 \quad \text{OR} \quad \mu_1 - \mu_2 > 0$$

The alternative is below its corresponding null hypothesis.

It's possible to add a constant to one side (eg  $H_0: \mu_1 - \mu_2 = 3$ ) but rarely done in practice

---

$$H_0: \mu_{\text{Dream Island}} = \mu_{\text{Toregersen}}$$

$$H_A: \mu_{\text{Dream Island}} \neq \mu_{\text{Toregersen}}$$

# Difference in 2 Means: Assumptions

Assumptions:

- ▶ Randomly collected data
- ▶ Independent and identically distributed
  - ▶ We assume they have the same mean for the moment
  - ▶ We do NOT assume they have the same variance (ie this is the unequal variance t-test; also called the Welch t-test)
- ▶ Both sample sizes (from Dream Island and Toregeresen Island) are large or the populations are both normal
  - ▶ Have to check for both groups!
  - ▶ Often indicate relevant statistics with a subscript to denote the two groups
  - ▶ Eg  $n_D$  and  $n_T$  for the sample size from Dream Island and Toregeresen Island, respectively

## Test Stat

General Form: 
$$\frac{\text{Observed} - \text{Hypothesized}}{\text{Standard Error}}$$

$$\text{Z-test statistic} = \frac{\bar{x} - \mu}{\sqrt{\sigma^2/n}}$$

$$\text{t-test statistic} = \frac{\bar{x} - \mu}{\sqrt{s^2/n}}$$

$$\text{t-test stat for 2 means} = \frac{(\bar{x}_D - \bar{x}_T) - (\mu_D - \mu_T)}{\sqrt{\frac{s_D^2}{n_D} + \frac{s_T^2}{n_T}}} = \frac{\bar{x}_D - \bar{x}_T}{\sqrt{\frac{s_D^2}{n_D} + \frac{s_T^2}{n_T}}}$$

NOTE: We are assuming  $H_0$  is true so  $\mu_D - \mu_T = 0$

## Df and P-values

Degrees of freedom is actually super complicated to calculate in this situation (see [Satterthwaits approximation](#) for the gory details)

In practice computers calculate the degrees of freedom at the same time it calculates the p-value (and honestly the test statistic) which brings us to a new R command.... `t.test()`

From R...

- ▶ test statistic = 1.701
- ▶ degrees of freedom = 111.37
- ▶ p-value = 0.0916

So we have weak evidence the mean flipper length is different between the two islands.

## Difference in Means Confidence Interval

And we can give a range of what we think the difference in means might be.

$$\bar{x}_1 - \bar{x}_2 \pm t_{df} \sqrt{\frac{s_D^2}{n_D} + \frac{s_T^2}{n_T}}$$

Again, we just use R because finding the degrees of freedom (and by extension  $t_{df}$ ) is complex...

( -0.309, 4.062)

We are 95% confidence the true difference between flipper lengths from Dream Island Penguins to Toregersen Island Penguins is between -.309 and 4.062 mm.

## Difference in Means Summary

Has similar inference as that for 1 mean.

This is definitely the point where R needs to be used to run this test effectively.

Big picture for me is to think about the (difference of means) as

- ▶ A single variable
- ▶ with mean 0 (for hypothesis test)
- ▶ or mean  $\bar{x}_1 - \bar{x}_2$  (for confidence intervals)
- ▶ and a “pooled standard error”

## Ex: Wool Breaks

The number of breaks in yarn given the type of yarn it was made of (type A or type B). To do this, 54 batches of wool were formed into yarn either using type A method or type B method. 27 bundles of yarn were made in both type A and type B style. The resulting bundles were tested on a loom and the total number of breaks found.

Find the

1. Factor(s) and it's levels
2. Experimental unit (= observation unit for this example)
3. A t-test to see if there is a difference in the two means
4. A 90% confidence interval around the differences in means

# Switching Gears

We are now going to go **\*back\*** to the z-test.

What is the least normal distribution you can think of?

# Switching Gears

We are now going to go **\*back\*** to the z-test.

What is the least normal distribution you can think of?

Probably an indicator variable that is 0/1?

## Z-test Review: Sampling Distribution

Previously, given our data was...

1. randomly selected
2. independent and identically distributed
3. pop is normal or the sample is large

our sampling distribution would be.....

$$\bar{x} \sim N(\mu, \sigma^2/n)$$

# Pop Question

What's the mean for an indicator variable called?

## Pop Question

What's the mean for an indicator variable called?

The mean of an indicator variable is called a **proportion**

Parameter (population proportion):  $p$

Statistic (sample proportion):  $\hat{p}$

## Pop Question

What's the mean for an indicator variable called?

The mean of an indicator variable is called a **proportion**

Why on earth am I defining what a proportion is?

## Pop Question

What's the mean for an indicator variable called?

The mean of an indicator variable is called a **proportion**

Why on earth am I defining what a proportion is?

Because you test a proportion the same way you test a mean with known variance, the z-test

## One last background piece....

Earlier I said the normal model can be used to approximate several other distributions.

One of those distributions is the *binomial distribution* which is (number of success) our of (number of trials).

There is some weird things about the binomial distribution....

- ▶ Defined by the probability of success  $p$  and number of trials  $n$
- ▶ It doesn't mention it's variance (unlike Normal)
- ▶ The variance is instead a function of  $p$  and  $n$ 
  - ▶  $np(1 - p)$
- ▶ Short hand is  $\text{Bin}(p, n)$  or  $\text{Binom}(p, n)$  or  $\text{Binomial}(p, n)$

## Inference for a proportion

First the sampling distribution and then we will talk about it's (not different) assumptions

---

$$\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$$

---

$$\hat{p} \sim N(p, \frac{p(1-p)}{n})$$

---

So the mean of a sample (be it normal or an indicator) will follow a normal distribution centered at the population's mean and with a (known) variance

# 1 Prop: Assumptions Unchanged

For a sample's mean...

1. The sample was randomly collected
2. The observations are independent and identically distributed
3. The population is normal or  $n$  is large

For a sample's proportion...

1. The sample was randomly collected
2. The observations are independent and identically distributed
3. ~~The population is normal or  $n$  is large~~
  - ▶ My variable is an indicator so clearly not normal
  - ▶ How large  $n$  should be now follows different guidelines

# 1 Prop: How large of $n$ do we need \*now\*??

The needed size for  $n$  is dictated by  $p$  actually....

- ▶  $np \geq 10$  AND
- ▶  $n(1 - p) \geq 10$

They must both be passed which will happen, for any  $p$  not 0 or 1, given  $n$  is sufficiently large.

---

Alternative way to think about this condition: Each category needs 10 observations (and this extends to more complicated situations, like the multinomial)

# 1 Prop: Hypothesis Test Statistic

For sample means, we assumed we knew  $\mu$ ....

$$\frac{\bar{x} - \mu}{\sqrt{\frac{\sigma^2}{n}}} \sim N(0, 1)$$

---

For sample proportions, we assume we know  $p$ ....

$$\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim N(0, 1)$$

# 1 Prop: P-value and Decision

P-value calculation remains the same as for the sample mean

- ▶ Calculate the probability in the tails of a normal distribution
- ▶ Which tail (or both tails) is dependent on  $H_A$

Decisions need to be stated in terms of a proportion

Eg: "We have strong evidence to suggest the proportion of general townsfolk of Grinnell are more religious than the student body"

# 1 Prop: Confidence Intervals

General Formula:

$$(\text{estimate}) \pm (\text{distribution value})(\text{st. error})$$

---

Sample Mean:

$$\bar{x} \pm z_{1-\alpha/2} \left( \sqrt{\frac{\sigma^2}{n}} \right)$$

---

Sample Proportion:

$$\hat{p} \pm z_{1-\alpha/2} \left( \sqrt{\frac{p(1-p)}{n}} \right)$$

## Caveat on Assumptions for Confidence Intervals

There is one thing that is annoying when it comes to the “ $n$  is large”

- ▶ For means we just gave guidelines on how not-normal the sample looked
- ▶ For 1 proportions hypothesis test we said  $np$  and  $n(1 - p) > 10$ 
  - ▶ In HT we assume we know  $p$

Confidence intervals don't assume we know the true  $p$  so what can we do?

## Caveat on Assumptions for Confidence Intervals

There is one thing that is annoying when it comes to the “ $n$  is large”

- ▶ For means we just gave guidelines on how not-normal the sample looked
- ▶ For 1 proportions hypothesis test we said  $np$  and  $n(1 - p) > 10$ 
  - ▶ In HT we assume we know  $p$

Confidence intervals don't assume we know the true  $p$  so what can we do?

Let's just plug in  $\hat{p}$ , our best guess at what  $p$  is

## Large $n$ Summary

For a sample mean...

Shape	needed $n$
Symmetric, not outliers	20-ish
Skewed	50-ish
Most anything	100-ish

---

For a sample proportion HT with a hypothesized proportion  $p$ ....

Formula	size
$np$	$\geq 10$
$n(1 - p)$	$\geq 10$

---

For a sample proportion CI....

Formula	size
$n\hat{p}$	$\geq 10$
$n(1 - \hat{p})$	$\geq 10$