Additional Topics in Hypothesis Testing

- One-Sided vs. Two-Sided Tests
- Error Types
- Choice of Level of Significance
- Cautions About Hypothesis Testing
- Importance of Results



REVIEW: HYPOTHESIS TESTING TERMINOLOGY

- Null hypothesis: initial assumption made about a parameter
- Alternative hypothesis: conclusion that is made if the evidence suggests the null hypothesis is false
- Test statistic: a criterion used to decide if a hypothesized value is a plausible value for the parameter
- P-value: measure of how unusual the test statistic is
- Level of significance: the cutoff that signifies that a p-value is small enough to reject the null hypothesis
 - Denoted by α ("alpha") and usually set to either 0.01, 0.05, or 0.10

REVIEW: HYPOTHESIS TESTING PROCEDURE

• We will cover 8 hypothesis tests this semester. Each is customized to the variable situation, but all follow the same procedure:

Step	Description
	Determine the appropriate test to use
1	Determine the null and alternative hypotheses
	Collect the data
2	Ensure the conditions for performing the test hold
3	Calculate the test statistic
4	Calculate the p-value
5	Calculate a confidence interval that matches the level of significance
6	Write a conclusion using the above results





ERROR TYPES

- Type I error: occurs when the null hypothesis is true in reality, but the evidence tells us to reject H_0
 - Common Reason: A bad sample resulted in a small p-value
- Type II error: occurs when the null hypothesis is false in reality, but the evidence tells us to fail to reject H_0
 - **Common Reason:** Hypothesized value was wrong, but the sample did not pick up enough evidence to reject that value

Reality	of the	Situation
---------	--------	-----------

		H ₀ True	H ₀ False
Your	Reject H ₀	Type I error	Correct decision
Decision	Fail to reject H ₀	Correct decision	Type II error

EXAMPLE: IDENTIFYING ERROR TYPES

- Scenario: You are buying a two-person kayak for you and a friend to use on the ocean. The two of you together weigh 400 pounds, which is the claimed weight capacity of the kayak. You want to test the hypotheses $H_0: \mu = 400$ vs. $H_A: \mu < 400$ by randomly sampling a set of reviews from previous buyers.
- **Result:** Many people lied about their weight in their review. You find that the average weight the kayak can hold is less than 400 pounds when in fact the true weight capacity is 400 pounds.
- Question: What type of error was made?
- Answer: _____
 - Decision: _____
 - Reality: _____

EXAMPLE: IDENTIFYING ERROR TYPES
• Scenario: You are buying a two-person kayak for you and a friend to use on the ocean. The two of you together weigh 400 pounds, which is the claimed weight capacity of the kayak. You want to test the hypotheses $H_0: \mu = 400$ vs. $H_A: \mu < 400$ by randomly sampling a set of reviews from previous buyers.
• Result: Previous buyers were honest about their weights. You find that 400 is a plausible value for the weight based on the sample. After testing it out, you find that it cannot hold both of you.
Question: What type of error was made?
Answer:
• Reality:
EXAMPLE: RAMIFICATIONS OF ERRORS
• Scenario: You are buying a two-person kayak and want to test the weight capacity using the hypotheses $H_0: \mu = 400$ vs. $H_A: \mu < 400$.
Question: What are the ramifications of each error type?
• Answer: • Type I Frror: You buying a kayak that would have
and instead look for a
• Type II Error: You spend money on a kayak that Even if you can , you and your friend will
Question: Which type of error is worse?
• Answer:
EXAMPLE: DESCRIBING ERROR TYPES
• Scenario: Rubber pellets are often used on synthetic turf football fields. However, these pellets often contain lead. The federal limit for lead levels in these areas is 400 μ g/g. Suppose a new football field is built. We want to test the hypotheses $H_0: \mu = 400$ vs. $H_A: \mu < 400$ by randomly sampling pellets from the field.
Question: How could a Type I and Type II error occur?
Answer:
Type I Error: Conclude that the field is(i.e) when in reality the field is) Type II Error: Conclude that the field is)
when in reality the field is (i.e)

EXAMPLE: DESCRIBING ERROR TYPES

- Scenario: Rubber pellets are often used on synthetic turf football fields. However, these pellets often contain lead. The federal limit for lead levels in these areas is 400 μ g/g. Suppose a new football field is built. We want to test the hypotheses $H_0: \mu = 400$ vs. $H_A: \mu < 400$ by randomly sampling pellets from the field.
- Question: What are the ramifications of each error type?
- Answer:
 - Type I Error: Many athletes get _____.
 - Type II Error: The field is _____ and _____
- Question: Which error type is more serious?
- Answer: _____

ERROR TYPES AND LEVEL OF SIGNIFICANCE

- The level of significance and the errors types are related.
 - 5% is a commonly used level of significance that will suffice in a majority of situations
 - Your choice of a level of significance may also be dictated by if a Type I or Type II error is worse.
- If a Type I error is worse, we should _____ the level of significance.
- If a Type II error is worse, we should _____ the level of significance.
- If neither error type is clearly worse, then stick with _____.

EXAMPLE: APPROPRIATE LEVEL OF SIGNIFICANCE

- Scenario: You are buying a two-person kayak and want to test the weight capacity using the hypotheses $H_0: \mu = 400$ vs. $H_A: \mu < 400$.
- Question: What level of significance would be appropriate to use?

-2

 $\alpha =$

0

- Answer: _____
 - _____ error is more serious want to be confident that _____



quite as much ______ for significance,

- giving us a better chance of _____ H_0
 - Makes Type II error ______

14



- Researchers have a decent amount of control over how significant their results are as they can often choose:
 - The population being studied
 - Form of the alternative hypothesis \rightarrow Discussed earlier
 - Value of the hypothesized mean
 - Sample size
- It is also possible to run multiple analysis on different samples of the data or using different hypotheses, but only report the desired results.
 - This is called **p-hacking** and is extremely unethical

EXAMPLE: CHOICE OF POPULATION

- Scenario: A market researcher is doing a study on the largest amount that females would pay for a pair of jeans and believes the mean is greater than \$100. She surveys 30 random females who are shopping at a large, suburban mall. Assume $\sigma = 25 .
 - **Results:** $\overline{X} = 116$, Z = 3.51, p = 0.0002, C.I.: (107.05, 124.95)
 - **Conclusion:** Researcher writes an article claiming the average amount a female would be willing to pay for jeans exceeds \$100.
- Question: What is wrong with the conclusion here?
- Answer:
 - Population that was studied was _____ not all females ______
 - Could be able to generalize it to all females who shop at ______
 - Sample is not _____ of the _____

• Scenario: An SAT prep course claims to improve scores. The mean improvement between attempts for 40 students was 85 points. Assume the standard deviation of improvements is 80 points. The prep course needs to decide how to present the results.					
Results: Claim Hypotheses Z P-Value					
0+ Pt. Improvement $H_0: \mu = 0$ vs. $H_A: \mu > 0$ 6.72 9.08 x 10 ⁻¹²					
50+ Pt. Improvement $H_0: \mu = 50$ vs. $H_A: \mu > 50$ 2.77 0.0028					
 Question: Which result makes a more compelling case for the effectiveness of the course? 					
Answer: Improving scores by					
 Test was still – more convinced the course actually as opposed to scores increasing 	/				
	- 20				
Example: Choice of Sample Size	20				
 Scenario: Two analyses are performed using α = 0.05: Test if mean fastball speed in baseball exceeds 96 mph Test if mean PPM of a chemical released from a car engine is less than 8 	3				
• Results: Analysis Hypotheses \overline{x} n Z P-Value					
Fastball $H_0: \mu = 96$ vs. $H_A: \mu > 96$ 96.06 7059 1.68 0.0465					
PPM $H_0: \mu = 8 \text{ vs. } H_A: \mu < 8$ 6 5 1.63 0.0515					
 Question: Which result is more meaningful? 					
 Answer:					
Importance of Results	21				
 Scenario: Study of fastballs in baseball found an average speed of 96.06 mph. This result yielded a p-value of 0.0465 when testing <i>H</i>₀: μ = 96 vs. <i>H_A</i>: μ > 96 from a random sample of 7059 pitches. Question: Why does the p-value not indicate whether a result is important? 					
Answer: We can by simply taking a enough sample	-				
• If we the standard error $\left(\frac{\sigma}{\sqrt{n}}\right)$ by <i>n</i> , then the test statistic $\left(Z = \frac{\bar{x} - \mu}{\sigma_{\perp}}\right)$ will					
 Smaller p-values result from test statistics Is the difference between 96 mph and 96.06 mph 	?				

IMPORTANCE OF RESULTS

•	Scenario: Study of a new car engine found that it releases only 6
	ppm of a certain chemical into the air. This result yielded a p-value
	of 0.0515 when testing $H_0: \mu = 8$ vs. $H_A: \mu < 8$ from a random
	sample of 5 engines.

- Question: Why might a study with a non-significant result still be important?
- Answer: Studies with small sample sizes have ______ achieving significance, but are often used as ______ to show ______

• Almost attained significance with only 5 engines

DANGER OF RUNNING MULTIPLE TESTS

Could convince an investor to ______ that could afford a ______ sample size

• Scenario: A hospital studied the relationship between cell phone usage and 20 types of brain cancers (called gliomas). There was no relationship found between cell phones and 19 of the gliomas, but one glioma was found to be related to cell phone usage at the 5% level of significance. Only this one result was reported.

• Question: Why is this unethical?

- Answer:
 - _____ the 19 gliomas that were _____ reporting only one of the results can cause _____ in the general population
 - Would expect _____ of the 20 tests to be significant ______
 - $\alpha = 0.05$ means extreme results occur 5% of the time even when H_0 is _____
 - Likely a ____
 - <u>Relevant xkcd</u> comic depicting a humorous play on p-hacking

23