

## Stat 218 - Exam 2 (Practice Exam)

## Formula sheet

	Single Variable	
	Proportion	Quantitative
Statistic	$\hat{p} = \frac{\text{Num. successes}}{n}$	$\bar{x} = \frac{\sum x}{n}$ , where $\sum x$ is the sum of all observations in the sample
Standard Error	$SE(\hat{p}) = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$	$SE(\bar{x}) = \frac{s}{\sqrt{n}}$
Standardized Statistic	$z = \frac{\hat{p} - \pi_0}{\sqrt{\pi_0(1 - \pi_0)/n}}$	$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$
Theory-Based CI	statistic $\pm 2 \times SE(\text{statistic})$	

Difference in proportions	
Overall proportion of successes $\hat{p} = \frac{n_1\hat{p}_1 + n_2\hat{p}_2}{n_1 + n_2}$	Standardized Statistic $z = \frac{\hat{p}_1 - \hat{p}_2 - 0}{\sqrt{\hat{p}(1 - \hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$
Theory Based CI $\hat{p}_1 - \hat{p}_2 \pm 2 \times \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}$	

**Concepts**

\_\_\_\_\_/15

Circle T (true) or F (false) for each of the following statements[1pt each]:

- T F** If a statistic A has a smaller standard error than statistic B, the width of the 95% interval for A will be smaller than the 95% interval for B
- T F** 2 *imes* SE is called the margin of estimation
- T F** The midpoint of a confidence interval is the parameter value
- T F** A confidence interval is the set of all possible values for the parameter
- T F** Using a higher confidence level produces a wider interval of plausible values
- T F** A representative sample is required to make unbiased inference about the population
- T F** Sampling bias occurs when there is not random assignment
- T F** An experiment is a study in which subjects are not randomly assigned to treatment groups or conditions
- T F** A quasi-experiment is one that manipulates the explanatory variable and uses random assignment of treatments to groups
- T F** A confounding variable is a variable that is related to only the explanatory variable
- T F** The relative risk is one way to summarize the association between two quantitative variables
- T F** We examine the difference in proportions across treatment groups because we want to know what proportion of people in the population respond to treatment.
- T F** In order to use theory-based inference when estimating the difference in two proportions, we must have a total of 40 observations with at least 10 observations in each group
- T F** The usual null hypothesis when comparing two proportions is  $H_0 : \mu_1 - \mu_2 = 0$
- T F** There are validity conditions for using simulation-based inference for the difference between two proportions

## Plague, Inc.

You are playing a game of Plague, Inc, where you are attempting to mutate your bacterium, ‘BubonicEbola’ so that it can kill everyone on earth. You decide to use statistics to investigate how the game mechanics work so that you can determine the best strategy.

### Infecting Madagascar \_\_\_\_\_/19

Early in the game, the goal is to evolve enough mutations that increase the infectivity of the virus, but are unlikely to be noticed (symptoms aren’t too visibly obvious). If the disease becomes noticed, countries will start working on a cure and will implement quarantine measures that make it difficult for the disease to spread across the globe. The country which is seen as hardest to infect is Madagascar, because it is only accessible by boat (no airports) and will close its ports as soon as an even moderately serious disease is identified.

You play 30 games using the same starting parameters and minor mutations (coughing, sneezing, cysts) that increase infectivity but are unlikely to be noticed, and in each game, you record whether or not Madigascar is infected at the end of 3 in-game years. In total, you manage to infect Madagascar 12 times.

1. [1pt] What is the observational unit?
  
  
  
  
  
2. [1pt each] For each of the following, fill in the value. If a quantity does not apply to this problem, write NA. If a quantity is unknown, write ‘unknown’. Show your work for any calculations.

$$\mu = \qquad \qquad \qquad \pi = \qquad \qquad \qquad \bar{x} =$$

$$\hat{p} = \qquad \qquad \qquad n = \qquad \qquad \qquad \text{successes} =$$

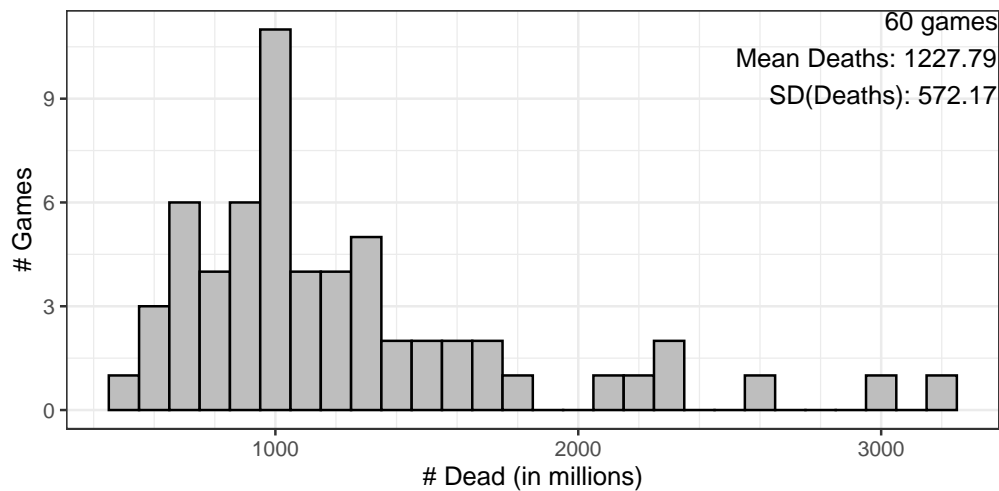
3. [1pt] What type of variable is this?
  
  
  
  
  
4. [2pt] What are the validity conditions for theory-based inference on this type of data?
  
  
  
  
  
5. [1pt] Are these conditions met? Explain.

6. [1pt] What formula should you use to calculate the standard deviation of the sample statistic?
7. [1pt] What is the value of the standard deviation of the sample statistic?
8. [1pt] What formula should you use to calculate an approximately 95% confidence interval for the long-run proportion of times you can infect Madagascar using your current strategy?
9. [1pt] Calculate an approximately 95% confidence interval for the long-run proportion of times you can infect Madagascar using your current strategy. Show your work. If you do not have a value for a component of the formula, define a variable for that value and show as much work as you can.
10. [2pt] Interpret your interval from 9 in the context of the problem. If you did not get an interval from 7, use  $(a, b)$  as your values and interpret as normal.
11. [2pt] If you wanted to decrease the width of the interval you calculated in 9, while still having 95% confidence, what approach would you use? Explain why your solution will decrease the width of the interval.

## Number of Deaths

\_\_\_\_\_/21

Once your disease has spread across the globe, the next part of your strategy is to significantly increase the lethality of BubonicEbola by evolving symptoms such as Dysentery, Hemorrhagic Shock, and Total Organ Failure. This time, you play 60 games, using this strategy starting at the beginning of year 3. For each game, you record how many people are dead worldwide (in millions) at the end of year 4. You would like to assess the long run performance of this strategy, as measured by body count.



1. [3pt] Describe the shape, spread, and center of the distribution of the number of deaths in the 60 games in the sample.
2. [2pt] Is this an experiment or an observational study? Explain why.
3. [1pt] What type of variable is this?
4. [1pt] Describe the population (in words)

Name:

ID:

---

5. [1pt] Is it possible to collect data on every item in the population? If yes, how would you do it? If no, why not?
6. [1pt] Can you use simulation-based inference for this data? Why or why not?
7. [2pt] What are the validity conditions for theory-based inference on the type of data you selected in question 3?
8. [2pt] Are these validity conditions met? Why or why not?
9. [1pt] What formula would you use to calculate the standard deviation of the sample mean?
10. [1pt] Calculate the standard deviation of the sample mean for your sample.

11. [2pt] What formula would you use to calculate an approximately 95% confidence interval for the population mean number of fatalities due to BubonicEbola using your strategy? Write it using the appropriate mathematical symbols for this type of variable.
12. [1pt] Calculate the approximately 95% confidence interval for the population mean number of BubonicEbola fatalities at the end of year 4. Show your work. If you do not have a value for a component of the formula, define a variable for that value and show as much work as you can.
13. [3pt] Interpret your interval from 12 in the context of the problem. If you did not get an interval from 13, use  $(a, b)$  as your values and interpret as normal.

## Comparing Infection Strategies

\_\_\_\_\_/27

When evolving a disease in Plague, Inc., you have multiple different upgrades you can focus on - Symptoms, Transmission, and Abilities. Symptoms increase the likelihood of infection, while Transmission upgrades focus on how the disease is spread (e.g. through animal contact, blood transmission, air, or in humid environments). You decide to compare strategies, playing 60 games in total. You will randomly sample 30 numbers from 1 to 60; games with sampled numbers will use only transmission-focused upgrades to increase the spread of the disease, while games with numbers not in the random sample will use only symptom upgrades to increase infectivity. For each game, you will play 3 in-game years, and at the end of that period, you will assess whether the disease has spread so that at least 50% of the world is infected (a success). Is focusing on one type of upgrades more likely to infect at least 50% of the world in 3 years than focusing on the other type of upgrades?

Strategy	Killed $\geq 50\%$	Killed $< 50\%$	Total
Symptom	12	18	30
Transmission	17	13	30

1. [1pt] Is this an observational study or an experiment?
2. [1pt] What is the research question? (Please state in the form of a question, with appropriate punctuation)
3. [1pt] Are the variables categorical or quantitative?
4. [4pt] If you were using a hypothesis test to answer the research question,
  - What would your null hypothesis be in words?
  - What is the null hypothesis, in symbols? Be sure to use the appropriate symbols for the variable type you identified in 2!
  - What is the alternative hypothesis, in words?
  - What is the alternative hypothesis, in symbols?
5. [1pt] What is the relevant statistic for this problem? What is the population equivalent? Be sure to use the correct symbols.



6. [2pt] What are the validity conditions for theory-based inference which are applicable to this problem? Are they met? Why or why not?
7. [2pt] Calculate the sample proportion for each group; using these values, calculate the value you identified in 5.
8. [2pt] What formula should you use to calculate the standard error of the quantity in 5, if you were to conduct a theory-based hypothesis test? If the formula uses symbols not defined in previous questions, please define those symbols as well. (Hint: Be sure you use the standard error formula for hypothesis tests, not confidence intervals! You may have to find this formula within a formula that is given to you.)
9. [2pt] What formula should you use to calculate the standard error of the quantity in 5, if you want to construct a confidence interval? (Hint: This is not the same formula as you used in question 7, and you may have to find it inside a different formula.)
10. [3pt] Calculate a 2SD confidence interval for the difference in shiny encounter rates between the two species; be sure to use correct notation and show your work.

11. [4pt] Interpret your interval in the context of the problem. If you did not get an answer for 10, use the (wrong) interval (0.1503, 0.3201).
12. [2pt] Given your answer in 10, what can you conclude about the value of the standardized statistic? Explain your conclusion.  
Note: Do not calculate the standardized statistic; instead, base your explanation on the confidence interval you interpreted.
13. [2pt] Given your answers in 10 and 11, would you reject or fail to reject the hypothesis  $H_0$  : the shiny encounter rates between the two species are the same?