

Practice Exam #2 - Sketch Solution

Ryan Miller

The following information will appear verbatim on the first page of Exam 2. You do not need to memorize this information, but you should be familiar with it.

Directions

- Answer each question using *no more than specified number of sentences* and not attempt to avoid these guidelines by using run-on sentences. Answers that are unnecessarily verbose may result in point loss.
- Do not include superfluous information in your answers, you may be penalized if you make an inaccurate statement even if you go on to provide a correct answer. Your answers should be clear, concise, and include only what is needed to answer the question that was asked.

Formula Sheet

Definitions:

- **Risk:** relative frequency of an event/outcome
- **Relative Risk:** ratio of the risks across two groups
- **Odds:** ratio of how often an event/outcome is observed relative to how often it is not observed
- **Odds ratio:** ratio of odds across two groups

Formulas:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$r = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

Statistic	Standard Error	Conditions
\hat{p}	$\sqrt{\frac{p(1-p)}{n}}$	$np \geq 10$ and $n(1-p) \geq 10$
\bar{x}	$\frac{\sigma}{\sqrt{n}}$	normal population or $n \geq 30$
$\hat{p}_1 - \hat{p}_2$	$\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$	$n_i p_i \geq 10$ and $n_i(1-p_i) \geq 10$ for $i \in \{1, 2\}$
$\bar{x}_1 - \bar{x}_2$	$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$	normal populations or $n_1 \geq 30$ and $n_2 \geq 30$

Question #1 (conceptual questions)

Part A: In 1-2 sentences, explain the meaning of “95% confidence” in the statistical term “95% confidence interval”. Avoid using the words “confident” or “confidence” as core components of your explanation.

- The term “confidence” describes the long-run success rate of the method used to come up with the interval estimate. Thus, if a procedure for constructing a confidence interval has a 95% confidence level we’d expect that procedure to produce an interval that contains the population parameter of interest 95% of the time (ie: for 95% of random samples that could arise)

Part B: Suppose your friend wants to use a statistical significance threshold of $\alpha = 0.01$ in a hypothesis test instead of the “conventional” threshold of $\alpha = 0.05$. In 1-2 sentences, briefly explain the trade-offs involved using your friend’s proposed threshold (relative to the conventional threshold) in terms of Type 1 and Type 2 errors.

- This change would make it harder to reject H_0 by requiring more evidence in the sample data. Consequently, the Type 1 error rate would decrease, but the Type 2 error rate would increase.

Part C: Consider a genomics experiment that involves 12,000 hypothesis tests each corresponding to a different genetic marker. Briefly explain why using a statistical significance threshold of $\alpha = 0.05$ could be problematic in this experiment.

- Using $\alpha = 0.05$ will give each hypothesis test where the null hypothesis is true a 5% chance of producing a Type 1 error. Thus, if none of these genetic markers are related with the outcome you’d expect 600 “statistically significant” findings, but all of them would be Type 1 errors.

Part D: Briefly explain the difference between a *sample* and a *population*. Why is this distinction important when describing the relationships found when analyzing sample data? Limit your response to at most 3 sentences.

- A population is the entire group of cases that you’d like to make a statement about. While a sample is a subset of that group that you have data on. This difference is important because the relationships in a sample are unlikely to perfectly reflect those in a population due to sampling error (sampling bias and/or sampling variability).

Question #2 (confidence intervals)

Earlier in the semester you encountered the American Community Survey data, a random sample of $n = 1287$ United States residents that is administered by the US Census Bureau on a rolling basis.

Shown below is some R output involving data. You may need to use some, all, or none of this output in the questions that follow.

```
acs = read.csv("https://remiller1450.github.io/data/EmployedACS.csv")

acs %>% group_by(Race) %>% summarize(Mean_Income = mean(Income),
                                     SD_Income = sd(Income),
                                     Mean_Age = mean(Age),
                                     SD_Age = sd(Age),
                                     n_HealthInsurance = sum(HealthInsurance == 1),
                                     n_Race = n())

## # A tibble: 4 x 7
##   Race Mean_Income SD_Income Mean_Age SD_Age n_HealthInsurance n_Race
##   <chr>      <dbl>    <dbl>  <dbl>  <dbl>         <int>  <int>
## 1 asian         61.2      77.5   41.2   14.1             86    92
## 2 black         34.5      29.9   43.5   14.2            106   116
## 3 other         31.2      42.4   38.9   13.6             80   102
```

```
## 4 white          45.5      55.5      43.7      15.4          907      977
```

```
cor(acs$Age, acs$Income)
```

```
## [1] 0.1649631
```

```
cor_results = cor.test(acs$Age, acs$Income, conf.level = 0.99)
```

```
cor_results$conf.int
```

```
## [1] 0.09431882 0.23395441
```

```
## attr(,"conf.level")
```

```
## [1] 0.99
```

Part A: Calculate a 95% confidence interval estimate for the difference in the proportions of white individuals and black individuals with health insurance in the United States. Show all of your work. You should use the value $c = 1.96$ to calibrate your interval.

- Here the 95% CI is given by point estimate $\pm c * SE$
 - The point estimate is $\hat{p}_1 - \hat{p}_2 = 907/977 - 106/116$
 - The calibration constant, c , is given as 1.96
 - Using CLT, $SE = \sqrt{\frac{907/977*(1-907/977)}{977} + \frac{106/116*(1-106/116)}{116}} = 0.02733455$
 - * Note that a pooled proportion is *not* used here because we are not interested in being consistent with a null hypothesis.
- The final answer is: (-0.039, 0.068)

Part B: Does the interval you calculated in Part A support the claim that white individuals are more likely to have health insurance than black individuals in the United States? Briefly explain, limiting your response to no more than 2 sentences.

- No, the interval suggests a difference in proportions of zero is plausible, so we are not confident that either group is more likely to have health insurance.

Part C: Suppose your friend decides that a 99% confidence level is more appropriate for the interval estimate you calculated in Part A. Would this interval suggest a *wider range* or a *narrower range* of plausible differences in the population? Briefly explain your answer, limiting your response to no more than 2 sentences.

- To increase the confidence level we'd need a procedure that generates wider intervals. Systematically wider intervals are necessary to produce a higher long-run success rate.

Part D: Consider the task of estimating the mean income of all individuals in United States belonging to each of the four racial categories used by the ACS researchers. This task would involve 4 different confidence intervals, one for each racial category. Briefly explain which of these interval estimates would have the smallest margin of error. You should assume all of the intervals use the same confidence level.

- For a single mean $SE = s/\sqrt{n}$, which is smallest for the sample of white individuals ($SE = 1.777$). Since each interval will have the same confidence level and the margin of error only depends upon the confidence level and the SE it will be this interval that has the smallest margin of error.

Part E: Without considering any other factors, interpret the relationship between age and income in the United States population based upon an analysis of the ACS data.

- As evidenced by $r = 0.165$, there is a weak positive correlation. We are confident that there is a positive correlation in the population and not just that sample by noticing the 99% CI doesn't contain zero.

Part F: Now consider the third variable "race" in regard to the interpretation you provided in Part E. Do you believe it is necessary to perform a stratified analysis by race to have a reasonable understanding of how age and income are related in the United States population? Briefly explain.

- Race is clearly associated with income, but it is not associated with age. Thus it doesn't confound the association between age and income and a stratified analysis is unnecessary.

Question #3 (hypothesis testing)

An experiment conducted at the University of Sydney in Australia investigated whether electrical stimulation to the brain could help participants successfully solve problems that required non-routine approaches.

The experiment trained 40 participants to solve problems in a particular way and then asked them solve an unfamiliar problem that required a creative solution, 20 of the participants were randomly assigned to receive electrical stimulation to the brain, and the other 20 received a placebo condition (the same apparatus without any electricity).

In the electrical stimulation group 60% successfully solved the problem, while only 20% of the placebo group solved the problem.

Part A: Perform the initial steps of an appropriate z or t test to evaluate the hypothesis that electrical stimulation improves the ability to solve non-routine problems, as summarized by a difference in proportions. Your answer should clearly state your hypotheses, and properly calculate a test statistic. *You do not need to find a p -value or make a conclusion.*

- $H_0 : p_1 - p_2 = 0$ vs. $H_A : p_1 - p_2 \neq 0$
- Assuming H_0 , the pooled proportion is the overall solving rate of $16/40 = 0.4$, so $Z = \frac{(0.6-0.2)-0}{\sqrt{\frac{0.4*(1-0.4)}{20} + \frac{0.4*(1-0.4)}{20}}} = 2.581$

Part B: Considering the assumptions of the hypothesis test you began in Part A, would you trust that the p -value that would arise from this test would be accurate? Briefly explain.

- No, in the placebo group there were only 4 observed successes, which violates the sample size conditions of a two-sample Z-test. An exact approach should be used to ensure a reasonably accurate p -value.

Part C: The figure below shows a randomization distribution for these data (generated under the hypothesis that an equal proportion of each group successfully solve the problem). Use this randomization distribution to *estimate* the p -value. Your estimate does not need to be exact; any estimate within reason will be scored as a correct answer.

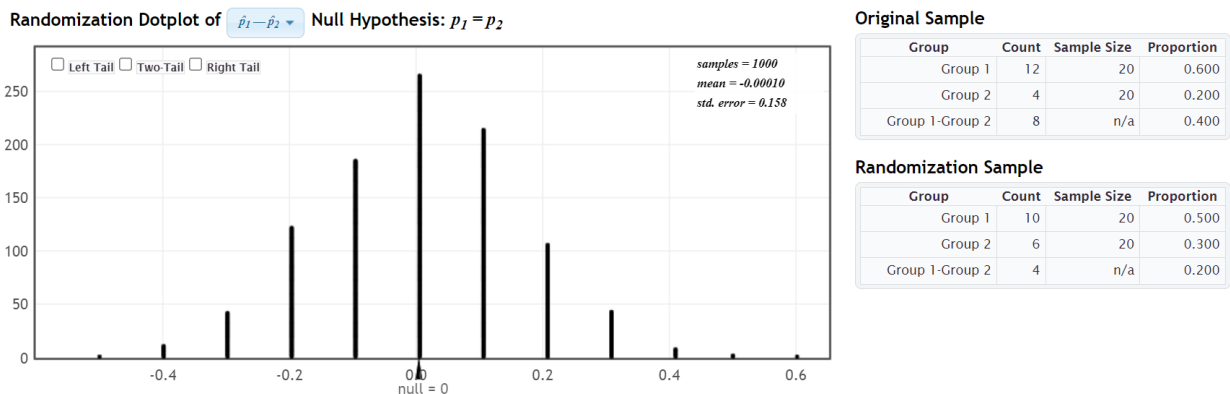


Figure 1: StatKey output used in Question Q3-C

- The observed difference in proportions was 0.4, and there are approximately 40 of 1000 simulated samples showing a difference of 0.4 or greater, so the two-sided p -value is approximately 0.04.

Part D: Provide a 1-2 sentence conclusion using the p -value you estimated in Part C and the context of the experiment.

- This p -value indicates our observed difference in proportions is large compared to what we'd expect under H_0 , thus we can conclude that the electrical stimulation is associated with an increase in problem solving ability relative to the control.

Part E: Suppose you used bootstrapping to construct a 95% confidence interval estimate for the difference in the proportion of problems solved. Would you expect this confidence interval estimate to suggest that a difference of zero is plausible? Briefly explain.

- No. The two-sided p -value is less than 0.05 (1 minus the confidence level), so in order for the interval and hypothesis test to yield a consistent conclusion, we know that CI should not contain the hypothesized difference of zero as specified in H_0 .