

Student name:

MATH 120: Elementary Statistics

Final exam

Section 1

May 6, 2019

Instructions:

- This is a regular “closed-book” test, and is to be taken without the use of notes, books, or any reference materials other than those provided with this test.
- Collaboration or group work is not permitted.
- Cell-phone usage, in any form, is prohibited for the entire duration of the test. This also applies to any restroom breaks taken during the test.
- The time limit for taking this test is 2 hours from the scheduled start time.
- This test adds up to 50 points.

Part I

Give short answers to each question as instructed.

- I.1. [4 pts.] Is there a relationship between the size of sand grains on a beach and its slope? Scatterplot of data from a sample of 20 beaches around the world suggests there may be a linear association. Based on these data, the following linear regression model was constructed to predict grain size (in mm.) from the beach slope (in degrees):

$$\widehat{\text{median sand diameter}} = 0.16 + 0.053 (\text{beach slope})$$

- (a) Identify the explanatory variable and the response variable, including their units.
(b) Interpret the meaning of the slope (with units) in this application context.

(a) Explanatory = Beach slope (in degrees)
Response = Median sand diameter (in mm.)

(b) Meaning of slope: For every 1° increase in beach slope, the predicted median sand diameter increases 0.053 mm.
The slope has units $\frac{\text{mm}}{\text{degree}}$

Grade: (a) = 2.5, (b) = 1.5

For (a): 1.5 pt = Explanatory + units; 1 pt = response + units

(b) 1 pt = correctly interpret slope; 0.5 pt = include units.

- I.2. [4 pts.] The following statistics summarize the age distribution (in years) of a 10-person team working on a project:

Median	Mean	IQR	SD
27.8	29	3.5	3.9

The oldest person, whose age is 34 years, leaves the team and is replaced by someone who is 36 years old. Compute, if possible, the new values of the above summary statistics. Show reasoning.

- * Median, Q1, Q3: Remain the same, since they are unaffected by increasing the largest value. Thus IQR remains the same.
- * Mean: Will increase since it depends on sum of all values.
old sum = $10 \times 29 = 290$. After the change, new sum = 292
New mean = $292/10 = 29.2$

- * SD: Should increase, but there is no way to find it without knowing all the data values.

Answer: Median = 27.8, Mean = 29.2, IQR = 3.5, SD = cannot determine

Grade: 1 pt each for correct median, mean, IQR, SD.
50/50 split between answer & reason or clarification

- I.3. [4 pts.] Explain the difference between each of the following pairs of statistical terms:

- (a) Sampling bias vs. sampling error.
- (b) Response bias vs. non-response bias.

(a) Sampling bias happens when a sample fails to accurately represent the underlying population, e.g., due to bad sampling methods, or undercoverage, or non-response.

Sampling error exists even when there is no bias, and it arises because every sample is necessarily an approximation of the underlying population.

(b) Response bias arises when a respondent gives an incorrect or misleading answer to a question because of how it is worded, or because of fear or embarrassment, etc.

Non-response bias arises when members of a sample fail to provide a response to a survey, or to question on it.

Grade: (a) = 2 points, (b) = 2 points
Roughly 1 point each for correctly addressing each of the 4 terms.

I.4. [4 pts.] A nationwide study on CEO compensation at publicly-traded companies surveyed 30 executives from Fortune 500 companies, and 54 from other smaller companies. They found the average annual compensation was \$9.2 million for the Fortune 500 sample, and \$7.7 million for the other group. We want to carry out a hypothesis test to determine whether this is evidence of a statistically significant difference.

(a) Write appropriate hypotheses, being sure to clarify what your parameters and subscripts denote.

Let μ_F = true mean compensation of ^{all} Fortune 500 CEO's,
 μ_0 = true mean compensation of all other company CEO's

Null hypothesis: $H_0: \mu_F - \mu_0 = 0$

Alt hypothesis: $H_A: \mu_F - \mu_0 \neq 0$

(b) Suppose the P -value for this test turns out to be 0.15. Will a 90% confidence interval for estimating the difference in mean compensation contain 0? Give reasons.

Yes, a 90% confidence interval will contain 0.

Since the P -value is 0.15, and the test was 2-tailed, the confidence level would need to be below 85%, (as that would match $\alpha = 0.15$) in order to find evidence of a significant difference.

I.5. [12 pts.] Each of the following questions requires a word, phrase, or numerical value as the answer. No reasoning or justification is needed.

i. A survey organization conducted telephone interviews in which 1,248 randomly selected adults in the United States were asked to respond to the question:

"At the present time do you think television commercials are an effective way to promote a new product?"

Identify the following as precisely as possible

* The population: All adults in the U.S.

* The parameter(s): The true % of all adults who think T.V. commercials are effective.

ii. Given the probabilities $P(A) = 0.4$, $P(A \text{ and } B) = 0.4$, $P(A \text{ or } B) = 0.4$, find

* $P(A | B)$: 1 [since $P(B) = 0.4$, and $P(A|B) = \frac{P(A \text{ and } B)}{P(B)}$]

* Are A and B disjoint events? (Yes, or No): No

2.5 points for (a):
 0.5 pt = clarify parameters
 2 pt = correct hypotheses

1.5 point for (b).
 1 pt = correct answer
 0.5 pt = reason

1 point each

1.5 pt →

0.5 pt →

- iii. A survey of employee job satisfaction at a large corporation reported the correlations shown in the table. The variables are: YS=years of service; SL=salary; PR=promotion rate; and JS=job satisfaction.

	YS	SL	PR	JS
YS	1			
SL	0.23	1		
PR	0.58	0.74	1	
JS	-0.79	0.82	0.88	1

Assuming the conditions necessary for interpreting correlations are met, are the following true or false:

* Higher promotion rates are associated with longer years of service: True

* Longer years of service are associated with greater job satisfaction: False

- iv. The monthly EPS (earnings per share) in dollars, over a 20-month period for a corporation are given below in ascending order:

Month	1	2	3	4	5	6	7	8	9	10	11	12
EPS	-2.4	-2.1	-1.7	-1.6	-0.2	-0.1	1.2	1.2	2.9	3.2	3.6	3.8

Month	13	14	15	16	17	18	19	20
EPS	4.0	4.1	4.1	4.2	4.4	4.7	4.8	5.0

Find the 5-number summary: $Q_1 = \text{between } 5^{\text{th}}/6^{\text{th}}$; Median = bet $10^{\text{th}}/11^{\text{th}}$
 $Q_3 = \text{between } 15^{\text{th}}/16^{\text{th}}$

5-number summary:

Min = -2.4, $Q_1 = -0.15$, Med = 3.4, $Q_3 = 4.15$, Max = 5.0

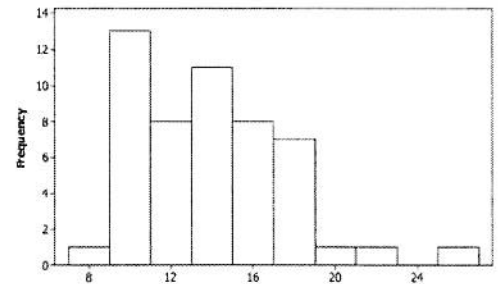
All values in dollars

- v. For the histogram shown, what summary statistics would best describe the center and spread?

Median and IQR

Will the mean be larger, or smaller, than the median?

Mean will be larger



- vi. A meteorology organization carried out a hypothesis test to investigate whether there is a difference in true mean rainfall between two cities ($H_0 : \mu_1 = \mu_2$; $H_A : \mu_1 \neq \mu_2$). They found evidence of a significant difference, using $\alpha = 0.05$.

* What does this tell about their P -value? (i.e., how large, or small, might its value be)

$P\text{-value} < 0.05$, since they found evidence of significance

* Suppose the significance level for the test is changed to $\alpha = 0.1$. Could that change the conclusion/inference from the test?

No, because

$P < 0.05$ guarantees that $P < 0.1$

Part II

Give complete and detailed solutions to each question. Grading will primarily be based on correct steps, reasons, relevant sketches, and clarity.

- II.1 [7 pts.] A team of scientists, researching the consequences of vitamin B₁₂ deficiency, tracked a group of 57 adults with B₁₂ deficiency for 7 years. At the end of this period they found 14 people in their sample exhibited symptoms of major depression.
- (a) What kind of study was this? (survey? observation? experiment? prospective? retrospective? etc.)
- (b) What is the nature and scope of conclusions this study can reach regarding B₁₂ deficiency and depression? Explain.
- (c) Use a confidence interval to estimate the true rate of depression among those with B₁₂ deficiency, based on data from this sample. Be sure to include all steps and state an appropriate conclusion.

(a) It was an observational, prospective study.

(b) Because it was not an experiment, the most the study can conclude is that there is an association between B₁₂ deficiency and symptoms of major depression among adults. However, even this claim would be difficult to make without knowing more details.

(c) Checking the conditions:

(I) Is the sample independent: not clear if it was randomly selected. But $n = 57$ is certainly less than 10% of adults with B₁₂ deficiency.

(II) Sufficiently large: successes = 14, failures = $57 - 14 = 43$
Both larger than 10. So the sample is large enough.

I will use a 90% confidence level. Thus $Z^* = 1.65$

$$\text{Margin of error} = 1.65 \times \sqrt{\frac{\frac{14}{57} (1 - \frac{14}{57})}{57}} = 0.0941$$

$$\text{Confidence Interval} = \left[\frac{14}{57} \pm 0.0941 \right] = [0.151, 0.3397]$$

Conclusion: We are 90% confident that the true rate of depression among those with B₁₂ deficiency lies between 15.1% and 34%. Note, however, the sample did not fully meet the required conditions for inference.

Grade: (a) = (b) = 1.5 points. (c) = 4 points

For (a): 1 pt = observational; 0.5 pt = prospective

For (b): Full credit for any reasonable statement that emphasizes association & why.

For (c): 1 point each for: ① conditions; ② M.E.; ③ C.I., ④ conclusion

II.2 [7 pts.] The 2-way table shows the distribution of rank of university faculty at all U.S. medical schools by sex in 2016 (Source: Association of American Medical Colleges):

	Faculty Rank			Total
	Assistant professor	Associate professor	Professor	
Female	36,498	12,589	8,708	57,795
Male	43,561	21,991	28,634	94,186
Total	80,059	34,580	37,342	151,981

- Find the probability that a random person from this sample is not an assistant professor.
- What is the probability that of 3 randomly selected persons none is an assistant professor?
- Find the probability that someone who is a professor is female.
- Using probabilities, determine whether faculty rank and sex are independent.

[Be sure to show all steps and justify their use.]

(a) This question is asking for: $P(\sim \text{assistant prof})$

$$P(\sim \text{assistant prof}) = 1 - P(\text{assistant Prof}) = 1 - \frac{80,059}{151,981} = \boxed{0.4732}$$

(b) This is asking for: $P(\sim \text{AP and } \sim \text{AP and } \sim \text{AP})$, AP = Assistant Prof.

we may assume independence, since the 3 are selected randomly. Also, since the sample size is large, I'll assume a constant probability for all 3 persons.

$$P(\sim \text{AP and } \sim \text{AP and } \sim \text{AP}) = (0.4732)^3 = \boxed{0.1060}$$

(c) This is asking for: $P(\text{Female} | \text{Professor}) = \frac{8708}{37,342} = \boxed{0.2332}$

(d) There are several options of probabilities we could look at to check whether $P(A|B) = P(A)$. I'll pick one of the easiest ones:

$$P(\text{Female}) = \frac{57,795}{151,981} = 0.3803 \quad \left. \vphantom{P(\text{Female})} \right\} \text{Not equal} \Rightarrow \text{Not independent}$$

$$\text{From (c), } P(\text{Female} | \text{Professor}) = 0.2332$$

Thus, faculty rank and sex are not independent.

Grade: (a) = 1 point. (b) = (c) = (d) = 2 points

(a) 50/50 split between step + answer

(b) 1.5 pt = show calculation step & reason; 0.5 pt = answer

(c) 1 pt = correctly interpret question; 1 pt = correctly compute

(d) 1 pt = know/show probability meaning of independence; ie. $P(A|B) = P(A)$

1 pt = correct computation + answer

II.3 [8 pts.] Does race matter when applying for National Institutes of Health grants? A study (reported in *Science*, August 2011) found that of 58,148 applications submitted by white researchers, 15,700 were funded by the NIH. Additionally, 198 of 1164 applications submitted by black researchers were funded. Is this evidence that the chance of funding is different for white and black researchers? Carry out a hypothesis test and state your conclusion. Show all steps, including: hypotheses; conditions check; sampling distribution model, with sketch; all calculations; and inference, with clear indication of what significance level you're using.

Preliminary inventory: This problem is about proportions; it involves 2 samples, and we want a hypothesis test.

I will use a significance level of 0.1 (i.e., $\alpha = 10\%$)

Let P_W, P_B denote the true proportion of White and Black researchers, respectively, who are funded by the NIH.

* Hypotheses: $H_0: P_W = P_B$
 $H_A: P_W \neq P_B$

* Conditions check:

① Is each sample independent? This is not clear from the information provided. There is certainly no indication of random selection. Possibly the samples are representative of their respective populations.

② Are the samples independent of each other?

Again, no indication of random selection, or other strategy to ensure independence.

③ Is each sample large enough? Yes, there are more than 10 successes and failures in each sample.

Overall verdict: Conditions may not be met. Interpret results with caution

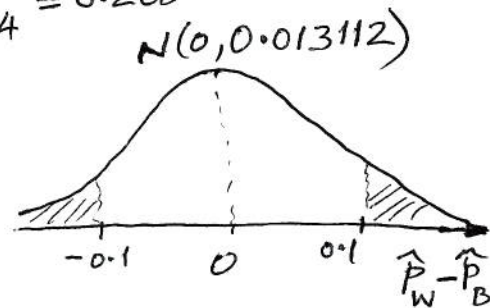
* Sampling distn. model:

Pooling is required here. $\hat{P}_{\text{pool}} = \frac{15,700 + 198}{58,148 + 1164} = 0.268$

The model is $N\left(0, \sqrt{\frac{0.268(1-0.268)}{58,148} + \frac{\text{same}}{1164}}\right)$

$$\hat{P}_W = \frac{15,700}{58,148} = 0.27, n_W = 58,148$$

$$\hat{P}_B = \frac{198}{1164} = 0.17, n_B = 1164$$



* P-value calculation:

$$Z = \frac{0.27 - 0.17}{0.013112} = 7.627$$

For this large Z-score, the P-value is 0 for all practical purposes

Conclusion: Since the P-value is below α , we reject the null hypothesis and infer that there is strong evidence of difference in the chance of NIH funding for White and Black researchers.

However, we note these conclusions are questionable, as the samples may not have met the conditions necessary for inference

Grade: 0.5 pt = clear significance level choice;

0.5 pt = clear statement of what P_W, P_B represent

1 pt each for the following 7 items: ① conditions check;

② Hypotheses; ③ sampling dist model; ④ Pooling;

⑤ P-value calculation; ⑥ Sketch; ⑦ conclusion