

FINAL EXAMINATION
ECON 103, STATISTICS FOR ECONOMISTS

MAY 1ST, 2018

YOU HAVE 120 MINUTES TO COMPLETE THIS EXAM. GRAPHING CALCULATORS, NOTES, AND TEXTBOOKS ARE NOT PERMITTED.

I pledge that, in taking and preparing for this exam, I have abided by the University of Pennsylvania's Code of Academic Integrity. I am aware that any violations of the code will result in a failing grade for this course.

Name: _____

Student ID #: _____

Signature: _____

Question:	1	2	3	4	5	6	Total
Points:	25	25	30	20	60	40	200
Score:							

Instructions: Answer all questions in the space provided, continuing on the back of the page if you run out of space. Show your work for full credit but be aware that writing down irrelevant information will not gain you points. Be sure to sign the academic integrity statement above and to write your name and student ID number on *each page* in the space provided. Make sure that you have all pages of the exam before starting.

Warning: If you continue writing after we call time, even if this is only to fill in your name, twenty-five points will be deducted from your final score. In addition, two points will be deducted for each page on which you do not write your name and student ID.

1. Someone in San Jose California, a city with 1 million residents, has stolen a car. The incident was captured on video, providing a detailed description of the thief. One out of every 10,000 people in San Jose fits this description. The police spot Amy walking down the street; she meets every detail of the description, so they arrest her. The only evidence against Amy is that she fits the description. The prosecutor argues as follows:

It is *highly* unlikely (far beyond a reasonable doubt) that an innocent person would fit this description. Hence, it is highly unlikely that Amy is innocent.

In parts (b)–(d) below, let I denote the event that someone from San Jose is innocent of the theft, and D be the event that she fits the description from the video.

- 5 (a) How many of the residents of San Jose fit the description from the video?

Solution: There are 1 million people in San Jose. One out of every 10,000 fits the description. Thus, $1000000/10000 = 100$ people fit the description.

- 5 (b) Calculate $P(I \cap D)$.

Solution: Of the 100 people in San Jose who fit the description, 99 are innocent and one is guilty. Hence, $P(I \cap D) = 99/1000000$.

- 5 (c) Calculate $P(D|I)$.

Solution: There are 1 million people in San Jose, and only one of them is guilty. Hence $P(I) = 999999/1000000$. By the definition of conditional probability,

$$P(D|I) = \frac{P(I \cap D)}{P(I)} = \frac{99/1000000}{999999/1000000} = \frac{99}{999999} \approx 1/10000$$

- 5 (d) Calculate $P(I|D)$.

Solution: From above, 100 out of people in San Jose fit the description. Thus, $P(D) = 100/1000000$. By the definition of conditional probability,

$$P(I|D) = \frac{P(I \cap D)}{P(D)} = \frac{99/1000000}{100/1000000} = 99/100$$

- 5 (e) In light of your responses to parts (c) and (d), evaluate the prosecutor's argument. Should the jury vote to convict Amy? Explain briefly.

Solution: The jury should *not* vote to convict Amy. The prosecutor is correct that, given someone is innocent, it is highly unlikely that she will fit the description: this probability is $P(D|I) \approx 1/10000$. The relevant probability, however, is the *reverse* conditional probability, namely the probability that someone is innocent given that she fits the description. This is $P(I|D) = 99/100$ so it is *highly likely* that Amy is innocent. The prosecutor's reasoning is fallacious. (FYI: this fallacy is commonly called the "prosecutor's fallacy.")

2. In this question you will determine the same probability two different ways: first by using the rules for calculating probabilities from class, and then by Monte Carlo simulation.

- 10 (a) There is an urn containing six balls, four of which are red. You make two random draws from the urn *without replacement*. What is the probability that both of the balls you draw are red?

Solution: There are $\binom{6}{2} = \frac{6!}{4!2!} = 15$ ways to choose two balls from a collection of six balls, provided that order doesn't matter. This is the denominator in our probability calculation. For the numerator, we need to count how many of these 15 possible draws contain exactly two red balls. Since there are 4 red balls in the urn, there are $\binom{4}{2} = \frac{4!}{2!2!} = 6$ ways to choose 2 of them, so the probability is $6/15 = 0.4$.

- 15 (b) Write R code to check your calculation from part (a) via Monte Carlo simulation using 10,000 simulation replications. For simplicity, I suggest that you create a vector called `urn` containing ones and zeros, where the ones represent the four red balls and the zeros represent the two other balls. You can then make random draws without replacement from `urn`.

Solution:

```
urn = c(1, 1, 1, 1, 0, 0)
urn_sim = function() {
  draw = sample(urn, 2, replace = FALSE)
  sum(draw) == 2
}
sims = replicate(10000, urn_sim())
mean(sims)
```

3. Let X be a continuous RV with support set $[0, 1]$ and pdf $f(x) = \alpha x^{\alpha-1}$ where $\alpha > 0$.

Name: _____

Student ID #: _____

- 10 (a) Calculate the CDF of X .

Solution:

$$F(x_0) = \int_{-\infty}^{x_0} f(x) dx = \alpha \int_0^{x_0} x^{\alpha-1} dx = \frac{\alpha x^\alpha}{\alpha} \Big|_0^{x_0} = x_0^\alpha$$

So the CDF is 0 for $x_0 < 0$, x_0^α for $x \in [0, 1]$ and 1 for $x_0 > 1$.

- 10 (b) Calculate $E[X]$.

Solution:

$$E[X] = \int_{-\infty}^{\infty} x f(x) dx = \alpha \int_0^1 x^\alpha dx = \frac{\alpha x^{\alpha+1}}{\alpha+1} \Big|_0^1 = \frac{\alpha}{\alpha+1}$$

- 10 (c) Calculate $\text{Var}(X)$. You do not have to simplify your answer.

Solution:

$$E[X^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \alpha \int_0^1 x^{\alpha+1} dx = \frac{\alpha x^{\alpha+2}}{\alpha+2} \Big|_0^1 = \frac{\alpha}{\alpha+2}$$

and by the shortcut rule,

$$\begin{aligned} \text{Var}(X) &= E[X^2] - (E[X])^2 = \frac{\alpha}{\alpha+2} - \left(\frac{\alpha}{\alpha+1}\right)^2 \\ &= \frac{\alpha(\alpha+1)^2 - \alpha^2(\alpha+2)}{(\alpha+2)(\alpha+1)^2} = \frac{\alpha}{(\alpha+2)(\alpha+1)^2} \end{aligned}$$

4. Let $X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ where σ^2 is known. In class we derived a $(1 - \alpha) \times 100\%$ confidence interval for μ that takes the form Estimator \pm ME. This is called a *two-sided* confidence interval, because it has two finite endpoints. It is also possible to construct *one-sided* confidence intervals, although we did not consider these in Econ 103. An *upper one-sided confidence interval* for a some parameter θ is defined as a range $[\text{LCL}, +\infty)$ constructed from the sample data such that $P(\text{LCL} \leq \theta) = 1 - \alpha$. In other words, the interval $[\text{LCL}, +\infty)$ covers θ with probability $1 - \alpha$. In this problem you will construct an upper one-sided confidence interval for μ based on X_1, \dots, X_n .

- 5 (a) What is the sampling distribution of $\sqrt{n}(\bar{X}_n - \mu)/\sigma$? Be sure to specify the values of any and all parameters of its distribution.

Name: _____

Student ID #: _____

Solution: $N(0, 1)$

- 5 (b) Continuing from the preceding part, write down the line of R code we would use to find the value of c such that $P(\sqrt{n}(\bar{X}_n - \mu)/\sigma \leq c) = 1 - \alpha$.

Solution: Let $Z = \sqrt{n}(\bar{X}_n - \mu)/\sigma$. We need to find c so that $P(Z \leq c) = 1 - \alpha$. Since $Z \sim N(0, 1)$, $c = \text{qnorm}(1 - \alpha)$.

- 10 (c) Using the expression from (b), derive the formula for LCL such that $[\text{LCL}, +\infty)$ is an upper one-sided confidence interval for μ with confidence level $(1 - \alpha)$.

Solution: Re-arranging the expression from (b),

$$P(\sqrt{n}(\bar{X}_n - \mu)/\sigma \leq c) = 1 - \alpha$$

$$P(-\mu \leq c\sigma/\sqrt{n} - \bar{X}_n) = 1 - \alpha$$

$$P(\mu \geq \bar{X}_n - c\sigma/\sqrt{n}) = 1 - \alpha$$

$$P(\bar{X}_n - c\sigma/\sqrt{n} \leq \mu) = 1 - \alpha$$

Hence, $\text{LCL} = \bar{X}_n - c\sigma/\sqrt{n}$ where $c = \text{qnorm}(1 - \alpha)$ by part (b).

5. This question is taken from your homework. It is based on a dataset containing the results of the tae kwon do event in the 2004 Athens Olympics. The competition is a tournament consisting of a number of bouts. In each bout, a pair of competitors fight each other, points are awarded, and a winner is declared by the judges. In accordance with Olympic regulations, one of the competitors in each bout is *randomly chosen* to wear blue body protectors. The other wears red body protectors. This question investigates whether wearing one color or the other gives an advantage in the competition. The data are stored in a data table called `taekwondo`. Each row corresponds to a *single bout*:

<code>red.id</code>	competitor id number for the fighter who wore red
<code>blue.id</code>	competitor id number for the fighter who wore blue
<code>round</code>	round of the tournament (i.e. semifinals, finals, etc.)
<code>winner</code>	color worn by the fighter who won the bout
<code>method</code>	method of win (i.e. points, knockout, etc.)
<code>red.points</code>	number of points awarded to the fighter who wore red
<code>blue.points</code>	number of points awarded to the fighter who wore blue

There are no missing values in the dataset. Here are the first few rows:

```
red.id blue.id  round winner  method  red.points blue.points
```

Name: _____

Student ID #: _____

1:	5816	5818 last 16	Blue Points	9	5
2:	5817	5824 last 16	Blue Points	3	5
3:	5819	5825 last 16	Red Points	15	16
4:	5820	5822 last 16	Red Points	14	15
5:	5821	5827 last 16	Red Points	13	12
6:	5828	5823 last 16	Red Knockout	7	3

- 4 (a) We'll restrict attention to the "last 16" round of the competition to ensure that each row contains a *unique* pair of fighters. Write R code to extract only those rows of `taekwondo` for which the value in the column `round` is "last 16" and store the result in a data table called `last16`.

Solution:

```
last16 = taekwondo[round == "last 16"]
```

- 6 (b) To begin, we'll analyze the *proportion* of bouts won by the blue fighter. Write R code to: (i) count the number of elements in the column `winner` of `last16` and store the result in a variable called `n`, and (ii) count the number of bouts won by the blue fighter and store the result in a variable called `n.blue`.

Solution:

```
n = last16[, length(winner)]
n.blue = last16[, sum(winner == 'Blue')]
```

- 10 (c) There are 32 bouts in `last16` of which 19 were won by the blue fighter. Using this information, calculate an approximate 95% confidence interval for the population proportion of bouts won by fighters wearing blue based on the approximation provided by the CLT. Do your results suggest that wearing one color versus the other conveys a competitive advantage? Explain briefly.

Solution:

$$\hat{p} = 19/32 \approx 0.59$$

$$\tilde{SE}(\hat{p}) = \sqrt{\hat{p}(1-\hat{p})/n} = \sqrt{\left(\frac{19}{32} \times \frac{13}{32}\right)/32} \approx 0.09$$

Hence, the CI is approximately $0.59 \pm 2 \times 0.09$ or roughly (0.41, 0.77). We do not find convincing evidence that either color conveys an advantage. If we absolutely had to guess, we would say that blue might convey a slight advantage but our results are perfectly consistent with the reverse as well: the difference

between the estimated proportion and 0.5 could easily be nothing more than sampling variability.

- 10 (d) Now suppose that you wanted to test the null hypothesis that the population proportion of bouts won by fighters wearing blue equals 0.5 against the two-sided alternative. Approximately what is your p-value for this test? Explain your results.

Solution: The test statistic is:

$$T = \frac{\hat{p} - 0.5}{\sqrt{0.5^2/n}} = \frac{19/32 - 0.5}{\sqrt{0.25/32}} \approx 1.06$$

If the test statistic were *exactly* one, the p-value for a two-sided test would be $2 * (1 - \text{pnorm}(1)) \approx 2 \times 0.16 = 0.32$. The test statistic here is slightly larger than one, so the p-value should be slightly smaller than 0.32. This is a very large p-value: we would *fail* to reject the null at any of the standard significance levels (i.e. 10%, 5%, 1%). We have not found convincing evidence that wearing either color conveys a competitive advantage.

- 6 (e) For the remainder of the question, we will examine the relative difference in the number of *points* scored by the blue and red fighters in each bout. Write R code accomplish the following: (i) select only those rows of `last16` for which the value in the column `method` is `Points` and store the result in a data table called `last16.points`, (ii) create a vector called `D` whose entries contain the *difference* in the number of points scored by blue versus red (Blue - Red) in each bout with `method` equal to `Points`.

Solution:

```
last16.points = last16[method == 'Points']
D = last16.points[, blue.points - red.points]
```

- 4 (f) I calculated the mean of the column `red.points` in `last16.points` and got 10.1. Similarly, I calculated the mean of the column `blue.points` and got 11.7. If I were to run the command `mean(D)` at the R console what result would I get?

Solution: $11.7 - 10.1 = 1.6$

- 10 (g) I entered the command `var(D)` at the R console and got 25. Next I entered `var(last16.points$red.points)` and `var(last16.points$blue.points)` and got

17 and 31, respectively. Calculate the sample correlation between the columns `red.points` and `blue.points` of the data table `last16.points`.

Solution: Rearranging the formula from class and substituting values from the question statement:

$$\begin{aligned} s_d^2 &= s_x^2 + s_y^2 - 2s_x s_y r_{xy} \\ 2s_x s_y r_{xy} &= s_x^2 + s_y^2 - s_d^2 \\ r_{xy} &= \frac{s_x^2 + s_y^2 - s_d^2}{2s_x s_y} \\ &= \frac{17 + 31 - 25}{2\sqrt{17} \times 31} = \frac{23}{2 \times \sqrt{527}} \approx 0.5 \end{aligned}$$

10

- (h) To test the null hypothesis that red and blue fighters are awarded, on average, the same number of points against the two-sided alternative, should we use a test for independent samples or matched pairs data? Explain briefly and then carry out the appropriate test at the 5% level based on the CLT. To answer, you will need the fact that there are 29 rows in the data table `last16.points`. Be sure to report: (i) the test statistic, (ii) the decision rule, and (iii) the result of the test.

Solution: This is matched pairs data: the score earned by the red fighter in a given bout cannot possibly be independent of the score earned by the blue fighter *in the same bout*. The test statistic is

$$T = \frac{\bar{D}}{s_d/\sqrt{n}} = 1.6/(5/\sqrt{29}) \approx 1.7$$

For a 5% test, the decision rule is: Reject H_0 if $|T| > 2$. In this case we fail to reject the null hypothesis.

6. This question is based on a data table called `face` containing data from a paper investigating whether “inferences of competence based solely on facial appearance predicted the outcomes of U.S. congressional outcomes.” Here are the first few rows of the dataset:

	year	state	Dcomp	Dshare	blueState
1:	2002	AK	0.3335918	0.1166021	0
2:	2004	AK	0.4193548	0.4777151	0
3:	2002	AL	0.5510950	0.4048368	0
4:	2004	AL	0.3223141	0.3236140	0

Name: _____

Student ID #: _____

5: 2002	AR	0.2738834	0.5386246	0
6: 2004	AR	0.7360000	0.5582051	0

Each row corresponds to a particular U.S. congressional election: **year** gives the year of the election, **state** gives the state, and **Dshare** gives the Democratic vote share. For example, the value of 0.1166021 for **Dshare** in the first row indicates that the Democratic candidate in the 2002 congressional race in Alaska got just under 12% of the vote. The column **blueState** is a dummy variable taking the value 1 if the corresponding state is a “blue state,” defined as a state in which most voters have favored the Democratic presidential candidate in past elections. We see from the first six rows of this dataset that Alaska (AK), Alabama (AL), and Arkansas (AR), are *red states*. The column **Dcomp** is a measure of the “facial competence” of the Democratic congressional candidate, constructed as follows. For each congressional race, a group of students were shown *unlabeled* photographs of the Democratic and Republican candidates side-by-side:



Which person is the more competent?

After viewing the photographs for one second, the students were asked which person appeared more competent, based on the photographs. The variable **Dcomp** records the *fraction* of students who thought that the Democratic candidate appeared more competent. For example, the value of 0.3335918 for **Dcomp** in the first row of **face** indicates that just over a third of the students thought that the Democratic candidate from the 2002 Alaska congressional race appeared more competent than his Republican rival. Any student who recognized either photograph was excluded from the calculation, so that **Dcomp** is constructed purely “based on judgements derived from facial appearance in the absence of prior knowledge about the person.” To be clear: the students did not know that the people in the photos were congressional candidates, or which was a Democrat.

Name: _____

Student ID #: _____

To answer the following parts you will need to refer to the regression results on the final page of this exam. You may want to tear out that page for ease of reference.

- 5 (a) What is the sample correlation between **Dcomp** and **Dshare**?

Solution: The R-squared from Regression #1 equals 0.19 and gives the *square* of this correlation. The correlation itself equals $\sqrt{0.19} \approx 0.44$.

- 5 (b) Construct an approximate 95% confidence for the slope in a linear regression in which **Dcomp** *alone* is used to predict **Dshare**. Briefly interpret your results.

Solution: From the results of Regression #1, our estimate of the slope is 0.33 with a standard error of 0.06. Thus, the approximate 95% confidence interval is 0.33 ± 0.12 or (0.21, 0.33). There is a *strong* positive relationship between **Dcomp** and **Dshare**.

- 5 (c) Continuing from the previous part, how accurately does **Dcomp** predict **Dshare**?

Solution: The residual standard deviation from Regression #1 is 0.13, indicating that **Dcomp** predicts **Dshare** to an accuracy of about 13 percentage points.

- 5 (d) Suppose you wanted to test the null hypothesis that Democratic congressional candidates do equally well, in terms of vote share, in blue and red states against the two-sided alternative. What is the value of your test statistic? Approximately that is the p-value of the test? Briefly interpret your results.

Solution: Using the results of Regression #2, the test statistic is $0.12/0.03 = 4$ yielding a p-value of less than 0.01. There is very strong evidence that Democratic congressional candidates do better (i.e. have a higher vote share) in blue states compared to red states.

- 5 (e) Is there evidence that the relationship between **Dcomp** and **Dshare** differs in red and blue states? Justify your answer.

Solution: For full points, an answer should point out that we need to look at estimate of the interaction term **blueState:Dcomp** in Regression #4. It should explain the meaning of the point estimate of -0.18, namely that the estimated slope is *smaller* in blue states, but also point out that the standard error is large. This difference could easily be due to sampling variability, a point that could be made either with a confidence interval or hypothesis test.

- 15 (f) Based on the full set of regression results, do you agree with the claim that “inferences of competence from faces predict election outcomes?” Justify your answer using the tools you have learned in Econ 103. Clear, concise answers will be treated more favorably than long-winded ones.

Solution: For full points, students should recognize that they need to examine the results of Regressions #1, #3, and #4. They should explain what each of these regressions tells us and point out that even after we control for `blueState`, either with or without an interaction, there is a strong positive relationship between `Dcomp` and `Dshare`. Facial competence clearly does predict results in US congressional elections. A complete answer should rely on statistical inferences to support this claim: either confidence intervals, hypothesis tests, or both.

Regression #1

```
lm(formula = Dshare ~ Dcomp, data = face)
      coef.est coef.se
(Intercept) 0.34    0.03
Dcomp        0.33    0.06
---
n = 118, k = 2
residual sd = 0.13, R-Squared = 0.19
```

Regression #2

```
lm(formula = Dshare ~ blueState, data = face)
      coef.est coef.se
(Intercept) 0.46    0.02
blueState    0.12    0.03
---
n = 118, k = 2
residual sd = 0.14, R-Squared = 0.15
```

Regression #3

```
lm(formula = Dshare ~ blueState + Dcomp, data = face)
      coef.est coef.se
(Intercept) 0.33    0.03
blueState    0.09    0.02
Dcomp        0.28    0.06
---
n = 118, k = 3
residual sd = 0.13, R-Squared = 0.27
```

Regression #4

```
lm(formula = Dshare ~ blueState + Dcomp + blueState:Dcomp, data = face)
      coef.est coef.se
(Intercept)  0.30    0.04
blueState    0.18    0.07
Dcomp        0.35    0.08
blueState:Dcomp -0.18    0.13
---
n = 118, k = 4
residual sd = 0.13, R-Squared = 0.29
```

Name: _____

Student ID #: _____