

FIRST MIDTERM EXAMINATION
ECON 103, STATISTICS FOR ECONOMISTS

FEBRUARY 14TH, 2017

You will have 70 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: _____

Signature: _____

Student ID #: _____ Recitation #: _____

Question:	1	2	3	4	5	6	Total
Points:	20	20	20	20	20	40	140
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, a point will be deducted for each page on which you do not write your name and student ID.

The following question appeared on the homework assignment to accompany Lecture #4:

- 20 1. What value of a minimizes $\sum_{i=1}^n (y_i - a)^2$? Prove your answer.
- 20 2. I have three boxes. The first contains two red balls, the second contains two blue balls, and the third contains one red ball and one blue ball. I choose a box at random so that each box is equally likely to be selected. I then reach inside and pull out a ball: it's red. Calculate the probability that the other ball in the box is also red. To help us grade your answer more easily, let B_1 be the event that I chose the first box, B_2 the second box, and B_3 the third box. Let R be the event that I drew a red ball.

3. Suppose we have a dataset $(x_1, y_1), \dots, (x_n, y_n)$ of scores for n students on two exams: x denotes exam #1 and y denotes exam #2. Let s_{xy} denote the sample covariance between exam scores, r denote the sample correlation, s_x and s_y the respective sample standard deviations, and \bar{x} and \bar{y} the respective sample means. Let $\hat{y} = a + bx$ where a is the intercept and b the slope of a linear regression calculated from this dataset.

- 2 (a) Write down the formula for a in terms of b , \bar{x} , and \bar{y} .
- 2 (b) Write down the formula for b in terms of s_{xy} and s_x .
- 2 (c) Write down the formula for r in terms of s_{xy} , s_x and s_y .
- 6 (d) Suppose that $b > r$. Then which is larger: s_x or s_y ? Explain.
- 8 (e) Suppose the sample mean score on the first exam, \bar{x} , was 75% while the sample mean score on the second exam, \bar{y} , was 85%. Moreover, our regression model predicts that someone who got 50 on the first exam will get 70 on the second. Calculate b .

- 20 4. Veronica wants to know the fraction of undergraduates who drink underage. Because she worries that her subjects may be unwilling to admit to breaking the law, Veronica conducts her interviews as follows:

Hi, I'm carrying out a survey on underage drinking. In a moment I'm going to ask if you drink underage, but before I do I'd like you to flip this fair coin in secret. If you get heads, answer my question truthfully. If you get tails, answer YES regardless of whether you drink underage. Because I don't know the outcome of your coin flip, if you answer YES I'll never know whether you really drink underage or just flipped a tails.

Veronica interviews a large number of subjects who all follow her instructions exactly: 85% answer YES. Calculate the fraction of undergraduates who drink underage.

5. Let X be a random variable with support set $\{-1, 1\}$, $p(1) = q$, and $p(-1) = 1 - q$.

5

(a) Write down the CDF of X .

3

(b) Calculate $E[X]$.

3

(c) Calculate $E[X^2]$.

3

(d) Calculate $Var(X)$.

6

(e) Let X_1 and X_2 be independent RVs both of which have the same pmf as X , defined in the problem statement. Write out the support set and pmf of $Y = X_1 + X_2$.

6. This question concerns a `data.table` called `star`. Here are its first and last five rows:

```
      race classtype yearssmall hsgrad g4math g4reading
1:      1          3           0    NA     NA         NA
2:      2          3           0    NA    706         661
3:      1          3           0     1    711         750
4:      2          1           4    NA    672         659
5:      1          2           0    NA     NA         NA
---
6321:    2          2           0    NA     NA         NA
6322:    1          2           0    NA     NA         NA
6323:    2          2           0     0     NA         NA
6324:    1          1           1    NA     NA         NA
6325:    1          1           1    NA     NA         NA
```

The data come from a randomized controlled experiment carried out in Tennessee in the 1980s. Students were randomly assigned to one of three class types, as indicated by the column `classtype`: 1 = small class, 2 = regular-sized class, 3 = regular-sized class with a teacher's aid. In the years that followed, researchers measured several outcomes: `hsgrad` is an indicator that takes the value 1 if a student graduated high school and 0 if not, `g4math` and `g4reading` are the student's scores on standardized 4th grade reading and math tests. In this question we will not work with the columns `race` or `yearssmall`.

- 5 (a) Say that the data are stored at `http://data.com/star.csv`. Give the full set of commands needed to read `star.csv` into R as a `data.table` called `star`.
- 5 (b) Give the R command needed to create a new `data.table` called `small` that only contains students who were randomly assigned to a small class.

- 5 (c) Write R code to carry out a linear regression that predicts a student's reading score from her math score. Use the data for all students in the experiment.

- 5 (d) Here are the results of running the command from the previous part:

Coefficients:

(Intercept)	g4math
317.5765	0.5696

If Beatrice scored 10 points higher in math than Alejandro, how many points higher would we predict that Beatrice will score in reading?

- 5 (e) Write R code to calculate the average high school graduation rate separately for students assigned to each different kind of class. Be sure to properly account for missing values.

- 5 (f) The results of the preceding part were as follows:

```
      classtype grad_rate
1:           1 0.8359202
2:           2 0.8251619
3:           3 0.8392857
```

Based on these results, do smaller class sizes improve high school graduation rates relative to regular-sized classes? Explain in no more than three sentences.

- 10 (g) Write an R function called `mycor` that takes two input arguments, a numeric vector `x` and another numeric vector `y`, and returns the sample correlation between them. In your answer you may use any R commands you like *except* the built-in correlation command `cor`. For simplicity you may assume that neither `x` nor `y` contains any missing values.