

Student Name (Print): \_\_\_\_\_  
(First Name) (Last Name)

**2020 Summer STAT 22000 Practice Final Exam**

*Important: read the following instructions carefully.*

1. During the exam, you may refer to the textbook, slides, homework, and the solutions and other materials posted on Canvas. You can use R during the exam. However, you cannot use Google or other search engines during the exam. You must do the exam all by yourself. You cannot get assistance from other people.
2. If a question asks you do some calculations, you must **show your work to receive full credit**.
3. Please check Canvas or email regularly during the exam. Yibi might send out corrections or clarifications about exam problems that you don't want to miss.
4. If you are unsure of what a question is asking for, **you may send questions to Yibi by email**.
5. Whenever appropriate, parts of a question will be graded conditionally on how you answered the preceding part(s). For example, even if you get part (a) of a question wrong, you will still get credit for the rest of the question provided your answers to parts (b), (c), etc. are consistent with how you answered part (a).

1. [Roulette] [16 points]

The game of roulette involves spinning a wheel with 38 slots: 18 red, 18 black, and 2 green. A ball is spun onto the wheel and will eventually land in a slot, where each slot has an equal chance of capturing the ball. One popular bet is that it will stop on a red slot; such a bet has an  $18/38$  chance of winning.

- (a) (2pts) Suppose a gambler bets on red in 6 different spins. What is the probability that he wins the first 3 spins but loses the next 3 spins?
- (b) (2pts) Let  $X$  be the number of times the gambler wins in the 6 spins. Explain why  $X$  has a binomial distribution.
- (c) (3pts) Suppose a gambler bets on red in 6 different spins. What is the probability that the gambler wins exactly 3 of the 6 spins, i.e.,  $P(X = 3)$ ?

For part (d-e) below, suppose the gambler bets on red in 50 different spins.

- (d) (4pts) How many times do you expect the gambler to win in the 50 spins? And with what standard deviation?
- (e) (5pts) Find an approximate value for the probability that the gambler wins at least 26 times in the 50 spins. Please calculate using normal approximation to Binomial WITH continuity correction. If you don't know how to do the continuity correction, you can also use normal approximation WITHOUT continuity correction and get 4 points if it is done correctly.



**3. [How often Read a Newspaper]** [9 points]

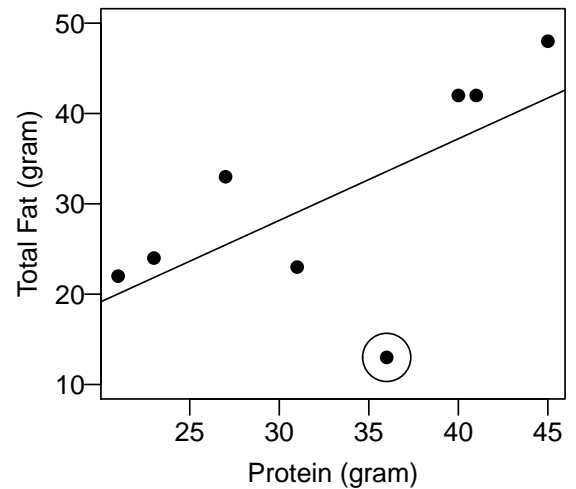
In a survey of a random sample of 50 students in a certain University, it is found that on average the subjects in the sample read a newspaper 4.1 times in the week prior to the survey, with a standard deviation of 3.0 times.

- (a) [4 points] Find a 95% confidence interval for the mean number of times students in this University read a newspaper the week prior to the survey.
- (b) [2 points] About the number of times students read a newspaper in the week prior to the survey, which of the following statement is true?
- (i) The population distribution is normal so it's legitimate to construct the confidence interval in (a).
  - (ii) The population distribution is not normal since this variable only take integer values. Hence we cannot construct a confidence interval based on a  $t$  distribution.
  - (iii) The population distribution is not normal, but we can construct the confidence interval in (a) as long as the sample contains no outlier and is not severely skewed, and the observations are independent.
- (c) [3 points] Explain to someone who knows no statistics what a 95% confidence interval in part (a) means. More specifically, what is the thing that has a 95% probability to happen?

4. [Fast Food] [9 points]

Data were obtained from the A&W Web site for the total fat in grams and the protein content in grams for various items on their menu. Some summary statistics and a scatter plot are also provided:

Item	Total fat (grams)	Protein (grams)
Kid' Cheeseburger	24	23
Kid' Hamburger	22	21
Original Bacon Cheeseburger	33	27
Original Bacon Double Cheeseburger	48	45
Original Double Cheeseburger	42	40
Crispy Chicken Sandwich	23	31
Grilled Chicken Sandwich	13	36
Papa Burger	42	41
Mean	30.875	33.000
SD	12.264	8.864
Correlation	$r = 0.653$	



- (a) [4 points] Find the equation of the least-squares regression line for predicting total fat from protein.
- (b) [5 points] For each of the following statements about the circled data point on the scatterplot, determine whether it is TRUE or FALSE. No explanation is required.
- (i) TRUE or FALSE: The circled data point on the scatterplot is for the Grilled Chicken Sandwich.
  - (ii) TRUE or FALSE: The residual associated with this data point will have a negative value.
  - (iii) TRUE or FALSE: This point would likely be considered an outlier.
  - (iv) TRUE or FALSE: This point has a high leverage.
  - (v) TRUE or FALSE: Without this point, the correlation between total fat and protein would be higher.

**5. [Hip Girth & Weight]** [14 points]

The scatterplot on the right shows the weights (in kg) and hip girths (in cm) of 46 physically active women age 35-44. The following regression output is for predicting women's body weight from their hip girth. Part of the output is blurred.

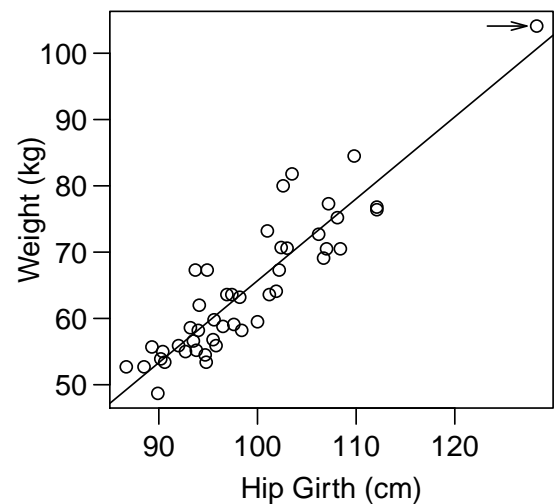
Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	XXXXXXXX	8.46196	XXXXXX	1.94e-08
hip.girth	XXXXXXXX	0.08523	XXXXXX	< 2e-16

Here is a summary of the data:

	Body weight (kg)	Hip girth (cm)
Mean	64.41	98.97
SD	10.80	7.95

Correlation  $r \approx 0.91$



- (a) (4pts) Write down the equation of the least square regression line for predicting women's body weights (in kg) from their hip girths (in cm).
- (b) (4pts) Calculate a 95% confidence interval for the slope of the regression line for predicting women's body weights (in kg) from their hip girths (in cm).

- (c) (2pts) If body weight is measured in pounds and hip girth is measured in inches, what will be the correlation between body weight and hip girth? (1 pound = 0.454 kg, 1 inch = 2.54 cm.)
- (d) (4pts) Predict the hip girth of a woman that weighs 60 kg using linear regression.

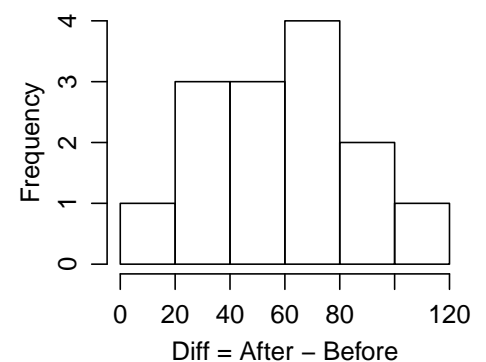
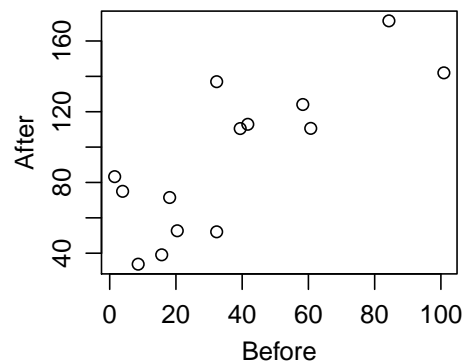
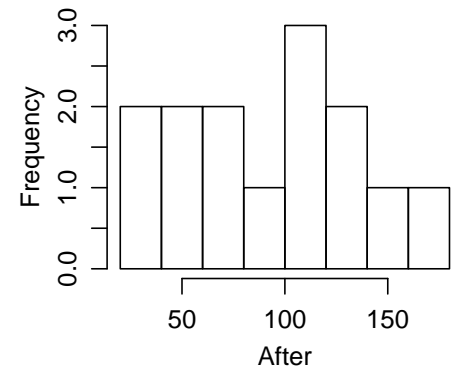
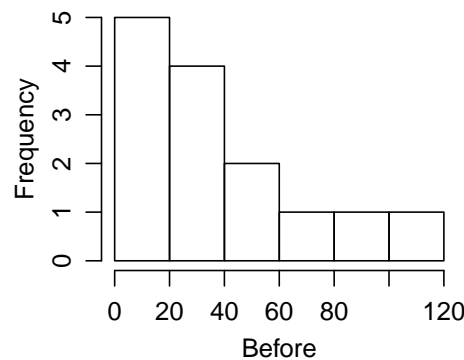


## 6. [Fortified Orange Juice] [9 points]

V. Tangpricha et al. conducted a study to determine whether fortifying orange juice with vitamin D would increase serum 25-hydroxyvitamin D [25(OH)D] concentration in the blood<sup>1</sup>. In medicine, the blood concentration of 25(OH)D is considered the best indicator of how much vitamin D is in the body. Most experts consider a serum 25(OH)D level of less than 30 nmo/L as indicative of vitamin D deficiency.

In the study, 14 subjects drank 240 mL per day of orange juice fortified with 1000 IU of vitamin D. Concentration levels were recorded at the beginning of the experiment and again at the end of 12 weeks. The before and after serum 25(OH)D concentrations in the blood, in nanomoles per liter (nmo/L), of the 14 subjects and their differences. The plots below are the histograms of the before and after serum 25(OH)D level and their scatter plot, as well as the histogram of their difference.

Subject	Before	After	Difference = After - Before
1	8.6	33.8	25.2
2	3.9	75.0	71.1
3	32.3	137.0	104.7
4	1.5	83.3	81.8
5	60.7	110.6	49.9
6	18.1	71.5	53.4
7	20.4	52.7	32.3
8	100.9	142.0	41.1
9	39.4	110.5	71.1
10	84.3	171.4	87.1
11	15.7	39.1	23.4
12	32.3	52.1	19.8
13	58.3	124.1	65.8
14	41.7	112.9	71.2
Mean	37.01	94.00	56.99
SD	29.94	42.10	26.20



- (a) [5 points] Construct a 99% confidence interval for the mean increase of the serum 25(OH)D concentration after 12 weeks of drinking fortified orange juice.

<sup>1</sup>V. Tangpricha et al. (2003) Fortification of Orange Juice with Vitamin D: A Novel Approach for Enhancing Vitamin D Nutritional Health. *American Journal of Clinical Nutrition*, Vol. 77, pp. 1478-1483

- (b) [2 points] Which of the following statement is true? No explanation is required.
- (i) The normality assumption for constructing the confidence interval in part (a) is violated because the distribution of the “Before” blood 25(OH)D level appears to be skewed.
  - (ii) The independence assumption for constructing the confidence interval in part (a) is violated because from the scatter plot, the “Before” and “After” blood 25(OH)D level appear to be highly correlated.
  - (iii) Both the above are true.
  - (iv) All the above are false.
- (c) [2 points] From the scatter plot, the correlation coefficient between the “Before” and “After” blood 25(OH)D level is closest to which of the following?

- (i)  $-0.7$       (ii)  $-0.2$       (iii)  $0.3$       (v)  $0.8$

No explanation is required.

**7. [T or F & Multiple Choice] [14 points]**

- (a) (6pts) Determine whether the following statements are true or false. No explanation is required.
- (i) TRUE or FALSE: The significance level of a test is the probability of making a Type 1 error when the null hypothesis is true.
  - (ii) TRUE or FALSE: Increasing the significance level will increase the probability of making a Type 2 error.
  - (iii) TRUE or FALSE: A Type 2 error is made if a correct null hypothesis is rejected.
  - (iv) TRUE or FALSE: We usually try to avoid making a Type 2 error more than to a Type 1 error
  - (v) TRUE or FALSE: The  $t$ -distribution has heavier tails than the normal distribution.
  - (vi) TRUE or FALSE: Confidence intervals based on a  $t$ -distribution will be shorter than confidence intervals based on the Normal distribution.
- (b) (2pts) A certain brand of cigarettes advertises that the mean nicotine content of their cigarettes is  $\mu = 1.5$  milligrams (mg). To test this, a random sample of 100 cigarettes of this brand were examined and the  $p$ -value for testing  $H_0 : \mu = 1.5$  mg versus  $H_a : \mu \neq 1.5$  mg was found to be  $= 3.2\%$ . Determine whether the following statements are true or false. No explanation is required.
- (i) TRUE or FALSE: A  $p$ -value of  $3.2\%$  means the probability that  $H_0$  is true is  $3.2\%$ . So the evidence supporting  $H_0$  is weak.
  - (ii) TRUE or FALSE: The value  $1.5$  mg is in the  $95\%$  confidence interval for the actual mean nicotine content of cigarettes of this brand.

(c) (4pts) A hospital administrator hoping to improve wait times decides to estimate the average emergency room waiting time at her hospital. She collects a simple random sample of 64 patients and determines the time (in minutes) between when they checked in to the ER until they were first seen by a doctor. A 95% confidence interval based on this sample is (128 minutes, 147 minutes), which is based on the normal model for the mean. Determine whether the following statements are true or false. No explanation is required.

- (i) TRUE or FALSE: About 95% patients at this hospital's emergency wait between 128 and 147 minutes
- (ii) TRUE or FALSE: The margin of error is 9.5 minutes and the sample mean is 137.5 minutes
- (iii) TRUE or FALSE: Doubling the sample size would cut the margin of error by half.
- (iv) TRUE or FALSE: If we used a different confidence level, the interval would be symmetric about the sample mean

(d) (2pts) The World Bank reports that 1.7% of the US population lives on less than \$2 per day. A policy maker claims that this number is misleading because of variation from state to state and rural to urban. To investigate this, she takes a random sample of 100 households in Atlanta to compare with the national average and finds that 2.1% of the Atlanta population live on less than \$2/day. Select the null and alternative hypothesis to test whether Atlanta differs significantly from the national percentage.

- |   |   |
|---|---|
| (i) $H_0: p = 2.1, H_a: p \neq 2.1$     | (ii) $H_0: \mu = \$2 \text{ per day}, H_a: \mu > \$2 \text{ per day}$ |
| (iii) $H_0: p = 0.017, H_a: p = 0.021$  | (iv) $H_0: p = 0.021, H_a: p \neq 0.021$                              |
| (v) $H_0: p = 0.017, H_a: p \neq 0.017$ |   |

No explanation is required.

(e) (2pts) The alumni association for U of Chicago has gathered a large dataset on graduating seniors. Some of the variables in the data set are gender, major, GPA, and starting salary. They are interested in looking at relationships among these variables. For which of the following pairs of variables would a two sample  $t$ -test be appropriate to examine whether there is a relationship between the two variables?

- |                        |                     |
|------------------------|---------------------|
| (i) gender and major   | (ii) GPA and gender |
| (iii) major and salary | (iv) GPA and salary |
| (v) none of the above  |                     |

No explanation is required.

8. [Offshore drilling] [20 points]

A 2010 survey asked 827 randomly sampled registered voters in California “Do you support? Or do you oppose? Drilling for oil and natural gas off the Coast of California? Or do you not know enough to say?” Below is the distribution of responses, separated based on whether or not the respondent is a college graduate.

	<i>College Grad</i>	
	Yes	No
Support	154	132
Oppose	180	126
Do not know	104	131
Total	438	389

- (a) (4 points) What proportion of respondents **opposed** offshore drilling among those with a college degree? What is the corresponding proportion among those without a college degree?
- (b) (8 points) Test whether or not those with a college degree had a different tendency to oppose offshore drilling from those without a college degree. Please state the null and alternative hypotheses, give an appropriate test statistic, report the  $p$ -value, and state the conclusion in the context using a 0.05 significance level.
- (c) (5 points) Construct a 95% confidence level for the difference  $p_{c,s} - p_{n,s}$  where  $p_{c,s}$  the proportion **supported** offshore drilling among those with a college degree, and  $p_{n,s}$  is the corresponding proportion among those without a college degree.
- (d) (3 points) Does the confidence interval in part (c) agree or contradict with the conclusion of the test in part (b)? Is it surprising or possible or is there anything wrong? Explain.