

NAME (FIRST LAST): _____ SID: _____

TIME AND CONDITIONS: You have 2 hours and 30 minutes to complete the exam and 30 minutes to upload your submission to Gradescope **and** record explanations on the two questions assigned to you. The exam is open-book, open-Jupyter notebook, open-calculator, but no other materials are allowed.

QUESTIONS AND ANSWERS

- There are 9 questions. Honor code is the first one.
- **Give brief explanations or show calculations in each question** unless the question says this is not required. You may use, without proof, any result proved or used in lecture, the textbook, and homework, unless the question asks for a proof.
- Please leave answers as **unsimplified arithmetic or algebraic expressions including finite sums** unless the question asks for a simplification.

GRADING

- The exam is worth 96 points. Each question is worth 12 points.
- Please commit yourself to a single answer for each question. If you give multiple answers (such as both True and False) then please don't expect credit, even if the right answer is among those that you gave.
- It is your responsibility to complete and submit the exam on time. In the first part (7:00 PM - 9:30 PM), you will have **2 hours and 30 minutes** to complete the exam and you **must stop** working on your exam at or before 9:30 PM and write down the exact time that you stop working underneath your honor code. In the second part (9:30 PM - 10:00 PM), you will have 30 minutes to complete both of two tasks below:

(1) You scan your work and upload to Gradescope. If you have technical difficulties, you must email your work to your GSI or stat88exams@gmail.com immediately. We will ignore any email that comes after 10:00 PM. **A late submission will not be accepted under any circumstances.**

(2) You will receive an email at 9:30 PM assigning you two questions randomly selected from your exam. **You have to record an explanation of everything you write down on those two questions.** If you did not provide an explanation to either one, you will receive a 0 on that question. You will have to record with your camera on and **you must finish recording by 10:00 PM**, although you may submit the recording link up to 24 hours after the exam.

FORMAT

- **You must answer each question on a separate page for a total of 9 pages.**
- Writing solutions on an iPad or tablet is acceptable.
- You must select the correct page associated for each subpart of each question when submitting to Gradescope. If you do not, you will get a zero for that question on the exam.

1. Honor Code

Data Science and the entire academic enterprise are based on one quality – integrity. We are all part of a community that doesn't fabricate evidence, doesn't fudge data, doesn't present other people's work as our own, doesn't lie and cheat. You trust that we will treat you fairly and with respect. We trust that you will treat us and your fellow students fairly and with respect. Please read carefully the UC Berkeley's Honor Code below:

"I certify that all solutions will be entirely my own and that I will not consult or share information with other people during the exam. I promise I will act with honesty, integrity, and respect for others."

"I understand if I do not match my response to the corresponding questions on Gradescope, my response will run the risk of not being graded."

Please transcribe the statement below and sign your name next to it:

"I understand and will abide by the above statements."

2. A random variable X has the distribution shown in the table below. Let μ_X and σ_X denote the mean and standard deviation of X .

x	-1	0	1
$P(X = x)$	1/18	16/18	1/18

- (a) **[3 points]** Show that $\mu_X = 0$ and $\sigma_X = \frac{1}{3}$.
- (b) **[3 points]** Use the probability distribution of X to calculate $P(|X - \mu_X| \geq 3\sigma_X)$. Compare this exact probability with the upper bound provided by Chebyshev's inequality and show that the bound provided by Chebyshev's inequality is attained.
- (c) **[3 points]** Now we take n independent random samples, X_1, X_2, \dots, X_n , drawn from the same probability distribution as X . Let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ be the sample mean. Apply Chebyshev's inequality to the sample mean to determine the smallest sample size n so that $P(|\bar{X} - \mu_X| \geq 3\sigma_X)$ is less than 0.01. Here μ_X and σ_X are still the mean and standard deviation of X in previous parts.
- (d) **[3 points]** Do you think the same statement of part(b) holds for the sample mean \bar{X} , i.e. if we use the probability distribution of \bar{X} to calculate $P(|\bar{X} - E(\bar{X})| \geq 3 \cdot \text{SD}(\bar{X}))$, it attains the bound provided by Chebyshev's inequality? Show your work. (*Hint: consider the case when n is large enough and use the fact that $\Phi(-3) = 1 - \Phi(3) \approx 0.0044$.*)

3. An experiment was conducted to observe the effect of a new antibiotic to treat an infection. Suppose that 100 randomly selected patients are given a standard antibiotic, and 85 of them recovered from an infection. Of the 100 randomly selected patients given a new antibiotic, 90 patients recovered from an infection. The patients in the two groups are independent of each other.

- (a) **[4 points]** Let p_1 denote the probability of recovery under the standard treatment and let p_2 denote the probability of recovery under the new treatment. Construct an approximate 95% confidence interval for $p_1 - p_2$, the difference of percentages of recovery under the standard treatment and the new treatment.
- (b) **[4 points]** Perform a test at a 5% level whether the percentage of recovery under the standard antibiotic differs from the percentage of recovery under the new antibiotic. You should follow the 5 steps for hypothesis testing discussed in class. (*Hint: use the fact that $\Phi(-1) \approx 0.1587$.*)
- (c) **[4 points]** Now suppose that 200 patients participated in the experiment and out of 200, a simple random sample of 100 participants received the standard treatment and the remaining participants received the new treatment. Again, 85 patients recovered in the first group, and 90 patients recovered in the second group. Perform a test whether the standard and the new antibiotic have different effects on the infection. State your null hypothesis, alternative hypothesis, test statistics, and p-value. You can write the p-value as an expression of finite sum and don't have to make a conclusion about rejecting or accepting the null.

4. Let X_1, X_2, \dots, X_n be independent exponentially distributed random variables with mean $1/\lambda$.

(a) [4 points] Show that $S = \min\{X_1, X_2, \dots, X_n\}$ has an exponential distribution with parameter λn .

(b) [4 points] Find an unbiased estimator of $1/\lambda$ that can be written in terms of S .

(c) [4 points] Let T be the unbiased estimator found in part(b). Find the mean squared error (MSE) of T and compare it with the MSE of the sample mean $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. (Hint: use the result of part(a).)

5. Suppose that X_1, X_2, \dots, X_n are i.i.d. random draws from $\text{Uniform}(a, b)$ and our goal is to estimate the unknown parameters a and b based on n samples. Let $T_1 = \min\{X_1, X_2, \dots, X_n\}$ and $T_2 = \max\{X_1, X_2, \dots, X_n\}$.

(a) [4 points] Show that $E(T_1 - a) = \frac{b-a}{n+1}$ and $E(b - T_2) = \frac{b-a}{n+1}$.

(b) [4 points] Find the expectations of $T_1 + T_2$ and $T_1 - T_2$. (*Hint: combine the two equations given in part(a).*)

(c) [4 points] Find unbiased estimators of a and b in terms of both T_1 and T_2 . (*Hint: combine the expressions for $E(T_1 + T_2)$ and $\frac{n+1}{n-1}E(T_1 - T_2)$ using the result of part(b).*)

6. Consider rolling a die for which the probability of rolling a six is $p \in (0, 1)$. David wants to test if the die is fair, i.e. there is no bias towards a six. David conducted 10 trials and obtained 1 six which was from the last trial. Then David left the lab.

- (a) [**6 points**] Danny continues David's work and performs a hypothesis testing for whether the die is biased towards a six or not. Formally, the null hypothesis is $H_0 : p = \frac{1}{6}$ and the alternative hypothesis is $H_A : p < \frac{1}{6}$. Danny uses the number of sixes out of David's 10 trials as the test statistic. Find the exact p -value. You can leave your answer as unsimplified expressions.

- (b) [**6 points**] In fact, David actually stopped his experiment immediately after 1 six, because his boss had instructed him to do so. David looked at the p -value that Danny obtained in part(a) and explains to him that the p -value should be changed. Explain why Danny's p -value should be changed. Which test statistic should be used and what is the distribution of the test statistic under the null distribution? Find the correct p -value. (*Hint: the fact that the experiment was stopped immediately after the first six should remind you of the distribution familiar to you.*)

7. Let (X, Y) be a random pair and let $\hat{Y} = \hat{a}X + \hat{b}$ be the linear regression of Y based on X . For notation, we use $r = r(X, Y)$ to denote the correlation between X and Y and use $D = Y - \hat{Y}$ to denote the residual. We write $E(X) = \mu_X$, $\text{Var}(X) = \sigma_X$, $E(Y) = \mu_Y$, and $\text{Var}(Y) = \sigma_Y$.

(a) [**6 points**] Show that $D = Y - \hat{Y}$ and \hat{Y} are uncorrelated, i.e. $r(D, \hat{Y}) = 0$.

(b) [**6 points**] Show that $\text{Var}(Y)$ can be decomposed into $\text{Var}(Y) = \text{Var}(\hat{Y}) + \text{Var}(D)$.

8. For each part (a) - (e) of this question, write down the numbers of **ALL** statements that are true. No credit will be given for partially correct answers (i.e. if a false statement is selected, or if a true statement is omitted.)

(a) [**2 points**] Let $(5, 10)$ be an instance of 95% confidence interval for a parameter θ . Which of the following statement is true:

- i The chance that the interval $(5, 10)$ contains θ is 95%.
- ii The chance that the interval $(5, 10)$ contains θ is 0% or 100%.
- iii None of the above.

(b) [**3 points**] Let T_1 and T_2 be two different estimators for a parameter θ . Which of the following statement is true:

- i If T_1 and T_2 are both unbiased estimators of θ , then $\frac{T_1+T_2}{2}$ is also an unbiased estimator.
- ii If T_1 is an unbiased estimator but T_2 is biased, then $\text{MSE}(T_1)$ is smaller than $\text{MSE}(T_2)$.
- iii If T_1 and T_2 are both unbiased estimators of θ and $\text{Var}(T_1) < \text{Var}(T_2)$, then $\text{MSE}(T_1)$ is smaller than $\text{MSE}(T_2)$.
- iv None of the above.

(c) [**3 points**] Let X be a continuous random variable with density f and CDF F . Which of the following statement is true:

- i f can be derived from F .
- ii F can be derived from f .
- iii For any x , $P(X = x) = 0$.
- iv $\int_{-\infty}^{\infty} F(x)dx = 1$.

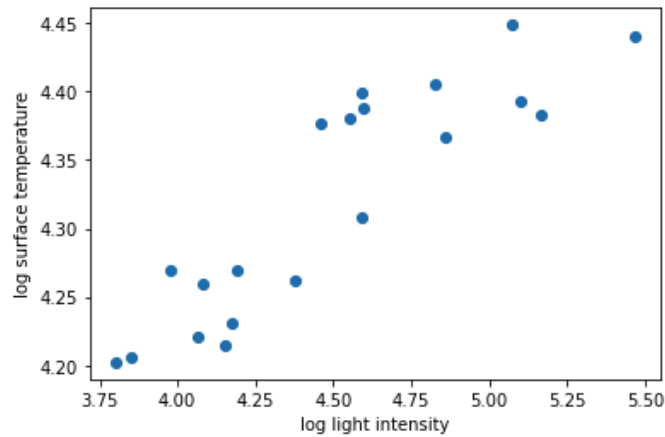
(d) [**2 points**] Let X_1, X_2, \dots, X_n be i.i.d. with mean μ and SD σ . Which of the following statement is true:

- i If n is large enough, $S = X_1 + X_2 + \dots + X_n$ approximately follows a normal distribution.
- ii If X_1, X_2, \dots, X_n follow normal distribution, $S = X_1 + X_2 + \dots + X_n$ exactly follows a normal distribution.
- iii As n increases, $S = X_1 + X_2 + \dots + X_n$ is more concentrated around its mean $n\mu$, i.e. for any $\epsilon > 0$, $P(|S - n\mu| < \epsilon) \rightarrow 1$ as $n \rightarrow \infty$.
- iv None of the above.

(e) [**2 points**] Let X have a $\text{Uniform}(0, 1)$. What is the distribution of $Y = -2 \log X$:

- i Exponential distribution.
- ii Normal distribution.
- iii Uniform distribution.
- iv None of the above.

9. A study was conducted to determine whether a linear relationship exists between log surface temperature (response) and log light intensity (predictor variable) for some nearby stars. Below is a scatter plot based on a random sample of 20 stars.



- (a) [2 points] From the scatter plot, discuss if the assumptions of the Simple Linear Regression Model are satisfied.

Now we run the linear regression on this data and the summary of the Python regression output is given below.

	coef	std err	t	P> t	[0.025	0.975]
const	3.6032	0.087	41.427	0.000	3.420	3.786
x	0.1597			0.000		

We additionally know that the sample mean and SD of log light intensity are 4.50 and 0.46, the sample mean and SD of log surface temperature are 4.32 and 0.082, and the SD of the residuals is 0.0039.

- (b) [2 points] What is the predicted log surface temperature if the log light intensity is 5?

(c) [**3 points**] Find the standard error of $\hat{\beta}_1$, $SE(\hat{\beta}_1)$.

(d) [**3 points**] Find the 95% confidence interval for β_1 . You may use one of the following percentile outputs from python:

```
stats.t.ppf(0.975, df=18) = 2.1  
stats.t.ppf(0.950, df=18) = 1.73  
stats.t.ppf(0.975, df=20) = 2.09  
stats.t.ppf(0.950, df=20) = 1.72
```

(e) [**2 points**] Find the sample correlation between log light intensity and log surface temperature. The sample correlation between (x_1, x_2, \dots, x_n) and (Y_1, Y_2, \dots, Y_n) is defined as

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