

Write as clearly as you can and in the box:

CSCI 3022
Final Exam
Spring 2019

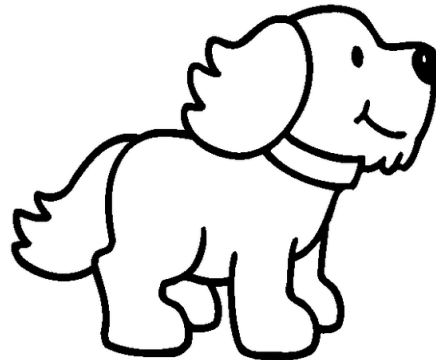
Name:

Student ID:

Read the following:

- **RIGHT NOW!** Write your name on the top of your exam.
- You are allowed **one** 3×5 in notecard of **handwritten** notes (both sides). No magnifying glasses!
- You may use a calculator provided that it cannot access the internet or store large amounts of data.
- You may **NOT** use a smartphone as a calculator.
- Clearly mark answers to multiple choice questions on the provided answer line.
- Mark only one answer for multiple choice questions. If you think two answers are correct, mark the answer that **best** answers the question. No justification is required for multiple choice questions.
- If you do not know the answer to a question, skip it and come back to it later.
- For free response questions you must clearly justify all conclusions to receive full credit. A correct answer with no supporting work will receive no credit.
- You have **150 minutes** for this exam.

Page	Points	Score
MC	32	
SA	8	
FR1	20	
FR2	20	
FR3	20	
Total	100	



Potentially Useful Values and Formulas

Standard Normal Distribution: Here $\Phi(z)$ is the cumulative distribution function for the standard normal distribution evaluated at z . Its equivalent form in Python is $\Phi(z) = \text{stats.norm.cdf}(z)$,


$\Phi(4.75) \approx 1.000$	$\Phi(3.00) = 0.999$	$\Phi(2.58) = 0.995$	$\Phi(2.32) = 0.990$	$\Phi(2.20) = 0.986$
$\Phi(2.00) = 0.977$	$\Phi(1.96) = 0.975$	$\Phi(1.80) = 0.964$	$\Phi(1.75) = 0.960$	$\Phi(1.64) = 0.950$
$\Phi(1.60) = 0.945$	$\Phi(1.50) = 0.933$	$\Phi(1.40) = 0.919$	$\Phi(1.28) = 0.900$	$\Phi(1.20) = 0.885$
$\Phi(1.04) = 0.850$	$\Phi(1.00) = 0.841$	$\Phi(0.93) = 0.825$	$\Phi(0.84) = 0.800$	$\Phi(0.76) = 0.775$
$\Phi(0.67) = 0.750$	$\Phi(0.60) = 0.725$	$\Phi(0.52) = 0.700$	$\Phi(0.45) = 0.675$	$\Phi(0.40) = 0.655$
$\Phi(0.32) = 0.625$	$\Phi(0.25) = 0.600$	$\Phi(0.19) = 0.575$	$\Phi(0.13) = 0.550$	$\Phi(0.00) = 0.500$

Student's t-Distribution: The following values of the form $t_{\alpha,v}$ are the critical values of the t -distribution with v degrees of freedom, such that the area under the pdf and to the right of $t_{\alpha,v}$ is α . Its equivalent form in Python is $t_{\alpha,v} = \text{stats.t.ppf}(1 - \alpha, v)$.

$t_{0.05,23} = 1.714$	$t_{0.025,23} = 2.069$
$t_{0.05,22} = 1.717$	$t_{0.025,22} = 2.074$
$t_{0.05,21} = 1.721$	$t_{0.025,21} = 2.080$
$t_{0.05,20} = 1.725$	$t_{0.025,20} = 2.086$

F-Distribution: The following values of the form F_{α,v_1,v_2} are the critical values of the F -distribution with v_1 and v_2 degrees of freedom, such that the area under the pdf and to the right of F_{α,v_1,v_2} is α . Its equivalent form in Python is $F_{\alpha,v_1,v_2} = \text{stats.f.ppf}(1 - \alpha, v_1, v_2)$.

$F_{0.1,2,9} = 3.006$	$F_{0.1,3,12} = 2.606$
$F_{0.05,2,9} = 4.256$	$F_{0.05,3,12} = 3.490$
$F_{0.025,2,9} = 5.271$	$F_{0.025,3,12} = 5.715$
$F_{0.01,2,9} = 8.022$	$F_{0.01,3,12} = 5.953$

Bayes' theorem	$p(A B) = \frac{p(B A)p(A)}{p(B)}$	Law of total probability	$p(E) = \sum_{i=1}^N p(E F_i)p(F_i)$
Union of sets	$p(A \cup B) = p(A) + p(B) - p(A \cap B)$	Conditional probability	$p(A B) = \frac{p(A \cap B)}{p(B)}$
Sigmoid function	$\text{sigm}(z) = \frac{1}{1 + e^{-z}}$	Regression	$\hat{\sigma}^2 = \frac{SSE}{n-2}, \quad SE(\hat{\beta}) = \frac{\hat{\sigma}}{\sqrt{\sum (x_i - \bar{x})^2}}$
	Get that bread!	Binomial coefficients:	$C(n, k) = \binom{n}{k} = \frac{n!}{k! (n-k)!}$
Three types of F you might use:	$F = \frac{(SST - SSE)/p}{SSE/(n-p-1)}$	$F = \frac{SSB/df_{SSB}}{SSW/df_{SSW}}$	$F = \frac{(SSE_{red} - SSE_{full})/(p-k)}{SSE_{full}/(n-p-1)}$
Some confidence intervals:	$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$	$\bar{x} \pm t_{\alpha/2, n-1} \frac{\sigma}{\sqrt{n}}$	$\left[\frac{(n-1)S^2}{\chi_{\alpha/2, n-1}^2}, \frac{(n-1)S^2}{\chi_{1-\alpha/2, n-1}^2} \right]$

Multiple choice problems: Write your answers in the boxes, or they will not be graded!



Figure 1. Tasty or snuggly?

- (2 points) You are writing a program to determine whether an object is a puppy or a blueberry muffin (see Figure 1), so you know which objects you can eat and which objects you can snuggle. The program tests the null hypothesis that the object is a puppy (i.e., that it is *not* edible) against the alternate hypothesis that the object is a blueberry muffin (i.e., that it *is* edible). If your program determines that an object is a blueberry muffin, you will eat it without question. What type of error is preferable?

- Type I error (false positive)
- Type II error (false negative)
- Type III error
- Type IV error
- Your program will classify an object perfectly every time, so there is no need to worry about errors

- (2 points) You are playing your favorite video game when an announcement appears on the screen: “ENTERING THE DATA SCIENCE LEVEL.” Your character appears in front of a graveyard. Each headstone has a birth date, a death date, and an animal mascot of Data Science semesters past. An instruction appears on your screen: “Press F-statistic to pay respects.”

Curious as to whether or not the features of (i) animal type, (ii) weight, (iii) hair majesty, and (iv) number of appendages affects the life span of each animal in a statistically significant manner, you perform a multiple linear regression F test on the data. You compute an F-statistic of $F = 2.5$, and you calculate that the F-critical value (F_{crit}) for your significance level is at 3.2. Assuming your F test is valid, which **one** of the following corresponds to an appropriate conclusion?

- All of the features are statistically significantly different from 0
- At least one of the features is statistically significantly different from 0
- None of the features is statistically significantly different from 0
- At least one of the features is **not** statistically significantly different from 0
- All of the features’ slope coefficients are statistically significantly different from 0
- At least one of the features’ slope coefficients is statistically significantly different from 0
- None of the features’ slope coefficients is statistically significantly different from 0
- At least one of the features’ slope coefficients is **not** statistically significantly different from 0

3. (2 points) You are at a psychic getting a reading done. You look into his crystal ball and see a multiple linear regression model. The psychic holds out two cards face down. One has an R^2 value on the other side, and the other has an adjusted R^2 value (R^2_{adj}). Both values correspond to the same multiple linear regression model, but one of the values is much greater than the other one. The psychic is not very good at his job, however, and asks you to tell him which card has the greater value. Which do you pick?

- A. The R^2 card because the model does not fit the data well
- B. The R^2 card because the model might not have enough features
- C. The R^2 card because the model might have extraneous (too many) features
- D. The R^2_{adj} card because the model does not fit the data well
- E. The R^2_{adj} card because the model might not have extraneous features
- F. Neither - they are always the same value
- G. You go home because you just wanted to see your future data science grade and this guy is terrible at his job

4. (2 points) You are playing a card game with a hedgehog named Alpaca. Alpaca is known to be a bit of a cheater. The probability that Alpaca is cheating, without knowing anything else, is 50%. Alpaca is a very good cheater, and when Alpaca cheats, they have an 80% chance of winning, as opposed to a 50% chance of winning when not cheating. Alpaca just won. What is the probability that Alpaca was cheating (to 3 decimal places)?

- | | |
|----------|----------|
| A. 0.1 | E. 0.5 |
| B. 0.25 | F. 0.615 |
| C. 0.315 | G. 0.815 |
| D. 0.4 | H. 0.9 |

5. (2 points) Samuel, your hedgehog companion, is working in a film class, and has recently posted a trailer for a movie he's planning. You show this to 100 of your internet friends and they each rate how much they like the trailer on a scale from 0 to 100. Then, you compute a 90% confidence interval for the mean rating of the trailer. The confidence interval you obtain is $[40, 60]$. You want to check this result using a larger sample so you share the video with 300 **more** people, for a new total sample size of 400. The mean rating from your new, larger sample is lower than the smaller sample by 2, but the standard deviation is the same. Using the sample size of 400, what is the 90% confidence interval for the mean rating of Samuel's movie trailer?

- A. $[39, 61]$
- B. $[43, 53]$
- C. $[40, 60]$
- D. $[45, 55]$
- E. $[38, 58]$

6. (2 points) You are walking down the street living your best life when a man dressed entirely in yellow and green attire yells at you from across the street, clutching Joseph the hedgehog. He says, "I have a challenge for you, Data Scientist. If you pass this challenge, I shall release Joseph. If not, then I will take Joseph to a farm upstate."

You bravely accept the challenge and begin walking toward this bizarre man. He sneers and asks:

"Oh, so you are approaching me? Well then, suppose you take an average of 4 steps per second. Which of these distributions is most appropriate if you wanted to find the probability that you take 7 steps in some amount of time?"

- | | |
|----------------------|--------------|
| A. Geometric | E. Uniform |
| B. Exponential | F. Poisson |
| C. Normal | G. Binomial |
| D. Negative binomial | H. Bernoulli |

☐

7. (2 points) Consider the following function. The function output constitutes a sample from which one of the following distributions?

```
def neat_o_complete_o(p):  
    x = 0  
    y = 5  
    while y > 0:  
        draw = np.random.choice([0,1], p=[1-p, p])  
        x += 1  
        if draw == 0:  
            y -= 1  
    return x
```

- | | |
|----------------------|--------------|
| A. Geometric | E. Uniform |
| B. Exponential | F. Poisson |
| C. Normal | G. Binomial |
| D. Negative binomial | H. Bernoulli |

☐

8. (2 points) Consider the following function. The function output constitutes a sample from which one of the following distributions?

```
def what_the_function():  
    x = 0  
    y = 5  
    draw = stats.uniform.rvs(loc=-5, scale=20) # uniform random draw between -5 and 15  
    if (draw >= x) and (draw <= y):  
        return 1  
    else:  
        return 0
```

- | | |
|----------------------|--------------|
| A. Geometric | E. Uniform |
| B. Exponential | F. Poisson |
| C. Normal | G. Binomial |
| D. Negative binomial | H. Bernoulli |

☐

9. (2 points) Suppose that you are performing a binary logistic regression classification to assign a class label $y \in \{0, 1\}$ to each data point and you model the probability that data point x belongs to Class 1 by

$$p(y = 1 \mid x) = \text{sigm}(\hat{\beta}_0 + \hat{\beta}_1 x)$$

where $\hat{\beta}_0 = 3$ and $\hat{\beta}_1 = 2$. If $p(y = 1 \mid x) \geq 0.5$, you classify as Class 1, otherwise, you classify as Class 0. How would your model classify a data point with $x = 2$?

Recall: $\text{sigm}(z) = \frac{1}{1 + e^{-z}}$

- A. inconclusive
- B. class 0
- C. class 1
- D. $\hat{y} = 0.5$
- E. $\hat{y} = 7$
- F. the limit does not exist



10. (2 points) For the same logistic regression model given in the previous problem, what is the decision boundary?

- A. $x = -3/2$
- B. $x = -2/3$
- C. $x = -1/2$
- D. $x = 0$
- E. $x = 1/2$
- F. $x = 2/3$
- G. $x = 3/2$



11. (2 points) Let $f(x)$ be the PDF of a normal distribution with mean 1 and variance 1. Compute

$$\int_{-0.4}^{2.2} f(x) dx$$

- | | |
|----------|----------|
| A. 0 | F. 0.642 |
| B. 0.081 | G. 0.804 |
| C. 0.115 | H. 0.885 |
| D. 0.196 | I. 0.919 |
| E. 0.331 | J. 1 |



12. (2 points) Suppose you generate 1,000 confidence intervals for the mean of a population, using fixed significance level α . You discover that 798 of them in fact do contain the true mean. Which of the following is the most appropriate estimate of the significance level α ?

A. 0.01
 B. 0.02
 C. 0.05
 D. 0.1
 E. 0.2
 F. 0.8



13. (2 points) Suppose you compute a sample mean, drawn from a population that is normally distributed with known variance σ^2 . Which combination of significance level α and sample size n produces the **narrowest** confidence interval for the mean?

A. $\alpha = 0.1$ and $n = 25$
 B. $\alpha = 0.1$ and $n = 36$
 C. $\alpha = 0.05$ and $n = 25$
 D. $\alpha = 0.05$ and $n = 36$
 E. $\alpha = 0.01$ and $n = 25$
 F. $\alpha = 0.01$ and $n = 36$



14. (2 points) The average weight of a Daurian hedgehog is about 590 g. The weight of these hedgehogs is known to be normally distributed, but an enterprising young data scientist/hedgehog enthusiast is concerned that the variation in the hedgehog weights is larger than acceptable. In an attempt to estimate the variance, she selects $n = 9$ hedgehogs at random and weighs them. The sample yields a sample variance of 30 g². Which of the following gives a 95% CI for the variance?

Recall that $\chi^2_{\alpha,\nu}$ and $t_{\alpha,\nu}$ are the χ^2 and t statistics (respectively) with α probability above it and ν degrees of freedom.

A. $\frac{8 \cdot 30}{\chi^2_{0.025,8}} \leq \sigma^2 \leq \frac{8 \cdot 30}{\chi^2_{0.975,8}}$

B. $\frac{8 \cdot 30}{\chi^2_{0.05,8}} \leq \sigma^2 \leq \frac{8 \cdot 30}{\chi^2_{0.95,8}}$

C. $\frac{8 \cdot 30}{-\chi^2_{0.025,8}} \leq \sigma^2 \leq \frac{8 \cdot 30}{\chi^2_{0.025,8}}$

D. $30 - \chi^2_{0.025,8} \cdot \sqrt{\frac{30}{9}} \leq \sigma^2 \leq 30 + \chi^2_{0.025,8} \cdot \sqrt{\frac{30}{9}}$

E. $590 - \chi^2_{0.025,8} \cdot \sqrt{\frac{30}{9}} \leq \sigma^2 \leq 590 + \chi^2_{0.025,8} \cdot \sqrt{\frac{30}{9}}$

F. $590 - t_{0.025,8} \cdot \sqrt{\frac{30}{9}} \leq \sigma^2 \leq 590 + t_{0.025,8} \cdot \sqrt{\frac{30}{9}}$

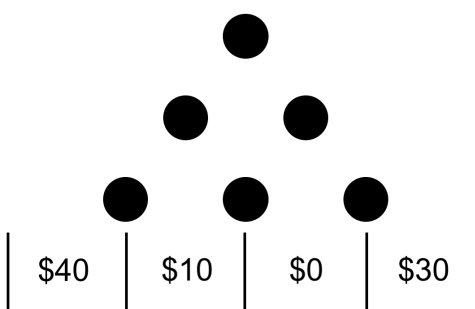


15. (2 points) Suppose that the random variable X has mean 5 and standard deviation 3. Let Y be the random variable given by $Y = X^2 - 2$. What is the expected value of Y ?

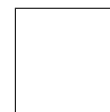
A. 5
 B. 17
 C. 23
 D. 25
 E. 32
 F. 36



16. (2 points) A game of **Plinko** is to be played on the board shown below. The pegs are unbiased, meaning that the disc has equal probability of moving left or right at each peg. Furthermore, the disc can only be dropped from directly above the top-most peg. What is the expected value of your winnings with a single disc?



A. \$8
 B. \$10
 C. \$12.5
 D. \$20
 E. \$26.67
 F. \$80



The rest of this page may be used for scratch work.

Short answer problems: If your answer does not fit in the box provided, make a note of where it is continued!

17. (4 points) What is a p -value? Describe, in words, what a p -value is meant to quantify. Then, provide a definition in terms of a probability (you may use words in your probability terms). What do we use it for?

18. (4 points) Amy, the famous hedgehog data scientist, likes to perform hypothesis tests to soothe her weary bones at the end of a long day of science. On this particular day, Amy draws a sample of size 100 from the set of all CU students, and computes the sample mean height. Amy decides that if her sample mean exceeds 70 inches (5 ft 10 in), then that is sufficient evidence to convince her that CU students are significantly taller than 67 inches (5 ft 7 in), which is about average for adults in the United States. From her previous work, Amy knows that the standard deviation of the distribution of all CU students' heights is 20 inches.

What significance level is Amy using in this hypothesis testing experiment? Show all work.

Free response problems: If your answers do not fit on the page and the one right after it, make a note of where the work is continued!

19. (20 points) Suppose you have collected a sample from a population whose probability density function is governed by some unknown parameter θ . You compute a test statistic, X , which has the following probability density function, where C is a constant.

$$f(x) = \begin{cases} \frac{1}{2} & \theta - C \leq x \leq \theta + C \\ 0 & \text{otherwise} \end{cases}$$

Please answer the following questions, and be sure to show your work—for sufficient space, a blank page follows this one.

- (a) (4 points) For what value of the constant C is $f(x)$ a valid probability density function? Remember to show all work!
- (b) (4 points) You want to see if there is sufficient evidence to conclude that the true parameter value is *greater than* $\theta = 1$. State the relevant null and alternative hypotheses.
- (c) (6 points) Sketch the probability density function $f(x)$, assuming that the null hypothesis from part (b) is true. Label your axes, the density function f , a few important x - and y -tick marks along your axes.
Suppose the test statistic that you compute is $X = 3/2$. Depict this in your sketch as a vertical line and label with “TS” or “Test Statistic”. Clearly mark and label in your sketch what the area/value/point that gives the p -value associated with the hypothesis test from part (b).
- (d) (4 points) Compute the p -value associated with the hypothesis test from part (b), assuming you have the test statistic $X = 3/2$. You must set this up **as an integral** first, then simplify and perform the calculation by hand to receive credit.
- (e) (2 points) What can we conclude, at the 10% significance level?

Additional Workspace

20. (20 points) Suppose you use statsmodels OLS to perform a simple linear regression of the form $y = \beta_0 + \beta_1 x$ on data consisting of $n = 23$ observations and obtain the following results:

OLS Regression Results						
Dep. Variable:	y	R-squared:	0.331			
Model:	OLS	Adj. R-squared:	0.299			
Method:	Least Squares	F-statistic:	10.39			
Date:	Sun, 39 Dec 2025	Prob (F-statistic):	0.00407			
Time:	3:45:67	Log-Likelihood:	-89.784			
No. Observations:	23					
	coef	std err	t	P> t	[0.025	0.975]
const	-1.6176	5.562	-0.291	0.774	-13.184	9.948
x	2.9037	0.901	?????	?????	?????	?????

Please answer the following questions, and be sure to show your work—for sufficient space, a blank page follows this one.

- (6 points) Compute the missing 95% confidence interval for the slope parameter.
- (2 points) Based on your confidence interval, do we have reason to believe that β_1 is different from zero? Fully justify your response.
- (6 points) Suppose we want to test the hypothesis that the slope coefficient β_1 is *greater than* 1. State the null and alternate hypotheses and perform a relevant hypothesis test at the 5% significance level. Clearly state your conclusions.
- (2 points) What fraction of the total variation in the response is explained by the simple linear regression model? Fully justify your response.
- (4 points) Suppose this linear regression model relates the feature X = “amount of time spent handling hedgehogs” (in minutes) and response Y = “exam scores”. Use your previous answers and the output above to evaluate the **strength** and **significance** of the relationship between hedgehog handling and exam scores. Fully justify your response.

Additional Workspace

21. (20 points) You are in Tony's CSCI 3022 Introduction to Data Science course, and you are talking to your favorite CA about your grade. She tells you about her hypothesis that students who laugh at Tony's bad jokes in class get better grades than those who do not. Being an aspiring data scientist yourself, you decide to test her hypothesis. You collect a data set consisting of 12 students' grades on a particular quiz. You separate the students into three groups: (1) students who have *No Reaction* to Tony's bad jokes, (2) students who give a *Pity Laugh* at the bad jokes, and (3) students who *Laugh Out Loud* at the bad jokes. You obtain the following data for the students' grades, and decide to do a **one-way ANOVA** test to compare the three groups of students.

No Reaction (N)	Pity Laugh (P)	Laugh Out Loud (L)
4	7	8
5	8	9
7	8	11
8	9	12

- (a) (4 points) Clearly state the null and alternative hypotheses for the one-way ANOVA test to compare the three groups of students and determine whether or not there is evidence that there is some difference in performance on this quiz.
- (b) (8 points) Compute the relevant **test statistic** to test the hypotheses from part (a). Put a box around your answer for the test statistic. Show all work!
- (c) (4 points) For a test at the $\alpha = 0.1$ significance level, perform a **rejection region** test for your test statistic from part (b). Be sure to clearly state (i) the distribution you are referencing (including any degrees of freedom), (ii) the critical value to which you are comparing your test statistic, and (iii) the conclusion of your test.
- (d) (4 points) For this part of the problem only, suppose you collected more data from the No Reaction and Pity Laugh groups. You collected 32 samples from the No Reaction group, with a sample mean of 6 and sample variance of 20. From the Pity Laugh group, you also collected 32 samples, with a sample mean of 8 and sample variance of 30. Compute a p -value associated with a hypothesis test to determine if there is evidence that the Pity Laugh and No Reaction population means are different. **Do not pick a significance level**, but comment on the strength of this evidence for a difference in group means.

Additional Workspace

