DATA 6 Fall 2025

Lisa Yan Quiz 2

Solutions last updated: Monday, October 13, 2025							
Your name:							
Your student ID:							
Your Berkeley email:							
Your room location:							
Student ID of the person	to your left: ₋						
Student ID of the person	to your right	:					
You have 50 minutes. The	ere are 3 que	estior	ns of	varyi	ng cr	edit. (34	points total)
	Question:	НС	1	2	3	Total	
	Points:	1	12	21	0	34	
For questions with circul may select only one choice		you		-			square checkboxes , you nore choices.
O Unselected option (ounfilled)	Completely					an select	
On't do this (it will incorrect)	be graded as	8				ple squa t do this)	
Only one selected op filled)	otion (compl	etely					
multiple answers, your a	nswer is am st interpreta	bigud tion.	ous, o For c	or the	bub g que	ble/chec stions w	ot be graded. If you write kbox is not entirely filled ith blanks, you may write sthan provided.
As a member of the UC I others. I will follow the r	-			act v	ith h	onesty, i	integrity, and respect for
Honor Code (HC): I have 1	read and agr	ee to	the h	onor	code	above.	
(1 point) Sign your name	_						

Q1 Yelp Business Reviews

(12 points)

Yelp is a website where users can review businesses on a 5-star rating scale.

name	stars	review_count	category
Bellevue Nails	3	47	Beauty & Spas
TLC Dental Center	4	10	Health & Medical
Bowmans Tavern	4	278	Restaurants
Enzo Custom	4	54	Shopping
Winter The Dolphins Beach Club	2	22	Hotels & Travel

Table 1: The first 5 rows of the businesses table (110,678 rows total).

Each row of businesses represents a business. Variable descriptions:

- name: Business name. There are over 80,000 unique business names.
- stars: Rating as an *integer*, 1 (lowest) to 5 (highest).
- review_count: Number of reviews for this business.
- category: Category of business: Restaurants, Shopping, Health & Medical, Beauty & Spas, Automotive, Hotels & Travel. Each business is in exactly one category.

	,	v S v	
Q1.1	(1 point) Which of the below work category?	ıld be best for visualizing how m	nany businesses are in each
	Bar chart	O Scatter plot	O None of the above
	○ Histogram	O Line plot	
Q1.2	(1 point) Which of the below wo business?	ould be best for visualizing the	distribution of reviews per
	O Bar chart	O Scatter plot	O None of the above
	Histogram	O Line plot	
Q1.3	(1 point) Which of the below wor	ald be best for visualizing the na	ames of businesses?
	O Bar chart	O Scatter plot	None of the above
	O Histogram	O Line plot	
Q1.4	(1 point) The review counts automotive_counts = make_arra		sinesses are defined as
	The p-th percentile of an array values. What is the 65th percent		least as large as p% of the
	O 10	○ 30	○ 50
	○ 25	• 40	O None of the above

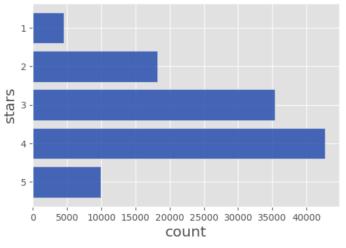
Consider the stars table, which counts how many businesses received each star rating:

stars	count
1	4506
2	18197
3	35332
4	42689
5	9954

Table 2: The stars table.

Q1.5 (2 points) Which line of code correctly uses businesses to create stars? Select all that apply.

Q1.6 (3 points) Assume **stars** is correctly assigned above. Which line of code produces the following visualization? Select all that apply.



stars.barh()	<pre>stars.barh("stars", "count")</pre>	businesses.hist("stars")
stars.barh("stars")	☐ businesses.barh("stars")	O None of the above

Q1.7 (3 points) Is the above visualization a histogram? Why or why not?

Solution: No; it is a bar chart, because stars are ordinal categorical variables. The spacing between stars on the y-axis doesn't mean anything.

SID: _____

Q2 The Iris Dataset (21 points)

The iris dataset measures flowers of three different species of the iris plant: *iris versicolor*, *iris setosa*, and *iris virginica*.

The petal and sepal are different parts of a flower. Each row in the iris table records one flower's species, its petal dimensions, and its sepal dimensions.







Iris Versicolor

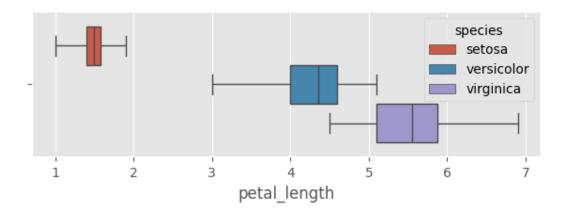
Iris Setosa

Iris Virginica

sepal_length	sepal_width	petal_length	petal_width	species
7.1	3	5.9	2.1	virginica
5.8	2.6	4	1.2	versicolor
5.8	4	1.2	0.2	setosa
4.9	2.5	4.5	1.7	virginica
6.1	2.8	4	1.3	versicolor

Table 3: The first 5 rows of the $\tt iris$ table (150 rows total). Measurements are in **centimeters (cm)**.

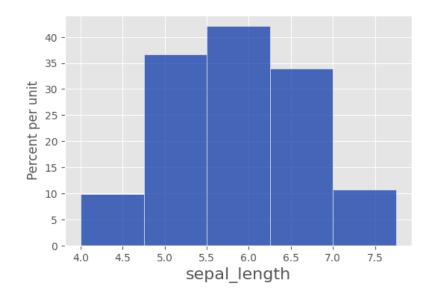
The iris table was used to generate the following boxplots of petal lengths (in cm) by species:



- Q2.1 (4 points) Using the boxplots above, which of the following are true statements about the flowers in the iris table?
 - The bottom 50% of *versicolor* flowers have shorter petals than all *virginica* flowers.
 - ☐ There are fewer *setosa* flowers than *virginica* flowers.
 - The distribution of *setosa* petal lengths is roughly symmetric.
 - ☐ For *versicolor* flowers, the mean petal length is greater than the median petal length.
 - O None of the above

SID:	

Q2.2 (5 points) Complete the code below so that supplying my_bins as the optional bins argument -Q2.6 to the provided hist method produces the histogram shown. The histogram plots sepal lengths in the iris table with bins: [4, 4.75), [4.75, 5.5), [5.5, 6.25), [6.25, 7), and [7, 7.75].



- Q2.7 (1 point) Consider the y-axis in the histogram above. What is the "unit" in "Percent per unit"?
 - O percent

O sepal width

O None of the above

O flower

- centimeter
- Q2.8 (1 point) Based on the histogram above, which bin contains the mean sepal length?
 - \bigcirc [4, 4.75)

(5.5, 6.25)

 \bigcirc [7, 7.75]

 \bigcirc [4.75, 5.5)

 \bigcirc [6.25, 7)

- O None of the above
- Q2.9 (1 point) Based on the histogram above, which bin contains the median sepal length?
 - \bigcirc [4, 4.75)

 \bullet [5.5, 6.25)

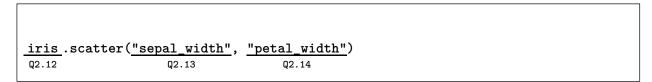
 \bigcirc [7, 7.75]

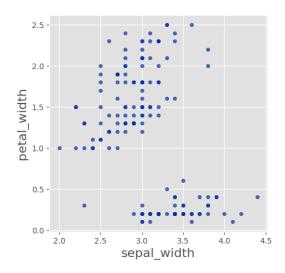
 \bigcirc [4.75, 5.5)

 \bigcirc [6.25, 7)

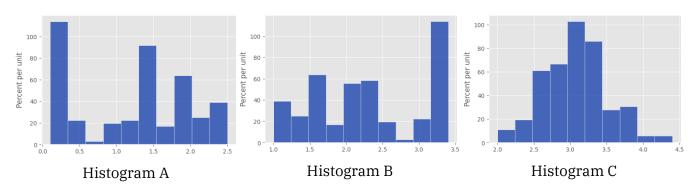
- O None of the above
- Q2.10 (2 points) Based on the histogram above, complete the code below to get all rows of the iris -Q2.11 table with sepal lengths in the bin with the most flowers.

Q2.12 (3 points) Complete the code below to produce the scatter plot shown, which plots petal -Q2.14 width by sepal width for all flowers in the iris table.





Still considering the scatter plot above, you are given three histograms. One corresponds to petal widths, one corresponds to sepal widths, and one corresponds to neither.



Q2.15 (2 points) If we run iris.hist('petal_width'), which histogram is produced?

Histogram A

O Histogram C

O Histogram B

O None of these

Q2.16 (2 points) If we run iris.hist('sepal_width'), which histogram is produced?

O Histogram A

Histogram C

O Histogram B

O None of these

ust for fun!	(0 point
raw something fun, or write a message for the staff! Or leave this bla	ank!
veryone had cool drawings/messages! :D	
everyone had coor drawings/messages: .D	

SID: _____