In []: # Initialize OK from client.api.notebook import Notebook ok = Notebook('final sp20.ok') In [1]: # Run this cell to set up the notebook, but please don't change it. import numpy as np import math from datascience import * # These lines set up the plotting functionality and formatting. import matplotlib matplotlib.use('Agg', warn=False) %matplotlib inline import matplotlib.pyplot as plots plots.style.use('fivethirtyeight') import warnings warnings.simplefilter(action="ignore", category=FutureWarning) # These lines load the tests. from client.api.notebook import Notebook ok = Notebook('final_sp20.ok') _ = ok.auth(inline=**True**) **Data 8 Final Exam Spring 2020** This exam is worth 150 points. You have until Friday, May 15 at 6PM to complete it. Late submissions will not be accepted. Do not use features of the Python language that have not been described in this course: we will not accept regrade requests for any issues caused by this. A few questions are marked Just for fun. These questions are optional and will not be graded. This exam is open notes: you may use any resources including the textbook, lecture, lab, homeworks, and projects, and old Piazza posts. Piazza will be closed to public posts. For clarifications, typos, and other similar issues with the exam only, you may make a private post. • The collaboration policy for this exam is similar to the homeworks: You may discuss the questions with other students, but any code, explanations, or text you write in this notebook must be your own. You must list the names of any students you discussed the questions with below in Question 0. Any clarifications to the exam will be made in this Google doc. In cases where we make a correction, we'll also make a Piazza announcement that you'll receive in your email. Please check the document and your email before you start or resume working on the exam. Unless otherwise stated, you may add cells or additional lines of code anywhere you want. Make sure you do not delete any existing cells. Also, do not reuse any variable names. Just like most assignments, unless otherwise stated, the public tests will not check for correctness. When budgeting your time, note that the later questions are longer and more difficult: make sure you leave enough time to attempt the entire exam. **Useful functions** These are useful functions from lecture and the textbook. They've all been defined correctly. You may use any of them for any question in the exam, Warning: Make sure you don't create any variables or functions that have the same names as any of these, or you may not receive credit on parts of the exam. Linear regression These functions are described in lectures 30-34 (Least Squares through Regression Wrapup), and also in Chapter 15 of the textbook. In [2]: def standard units(arr): return (arr - np.average(arr))/np.std(arr) def correlation(t, x, y): x_standard = standard_units(t.column(x)) y_standard = standard_units(t.column(y)) return np.average(x_standard * y_standard) def slope(t, x, y): r = correlation(t, x, y) $y_sd = np.std(t.column(y))$ $x_sd = np.std(t.column(x))$ return r * y_sd / x_sd def intercept(t, x, y): $x_{mean} = np.mean(t.column(x))$ $y_{mean} = np.mean(t.column(y))$ return y_mean - slope(t, x, y)*x_mean def fitted_values(t, x, y): a = slope(t, x, y)b = intercept(t, x, y)return a*t.column(x) + b def fitted_value(t, x, y, new_x): a = slope(t, x, y)b = intercept(t, x, y)return a*new_x + b def residuals(t, x, y): predictions = fitted_values(t, x, y) return t.column(y) - predictions Question 0 Please list the names and email addresses of everyone you collaborated with when taking this exam. As a reminder, you are allowed to discuss the questions with other students, but you must write all code and fill in all answers by yourself. Write your answer here, replacing this text. **Question 1** For each of the following cells, assign the specified variable so that the output of the entire cell is as specified. You may only change the first line of each cell: **don't change anything else**. For example, if the cell said to assign the variable x to produce the output 10, and it looked like this: $x = \dots$ x + 7then the correct answer would be to assign x = 3, so that the output of the last line would be 10. Question 1.1 (2 points) Assign the variable a so that the output of a + b is an array of two integers, whose first element is 2020 and whose second element is 8. In [3]: a = ...# Do not change anything below this line. b = make array(2017, 0)a + bIn []: | ok.grade("q1_1"); Question 1.2 (2 points) Assign the variable a2 so that the output of running the below cell is an array of integers whose first element is 2020 and whose second element is 8. In [6]: a2 = ...# Do not change anything below this line. t = Table().with columns('year', np.arange(2000, 2030, 10), 'course', make_array(8, 100, 102), first_col_index = a2.item(0) second col index = a2.item(1)year = t.column('year').item(first_col_index) course = t.column('course').item(second_col_index) make_array(year, course) In []: | ok.grade("q1_2"); **Question 2: Books and Graphs** Zeynep goes through her family's entire collection of 370 books and records data for each one. She records her data in a table called books . The lengths of the books (in pages) are in a column labeled Length . For this question, you won't be able to use the table books: you must answer based only on the information provided. She notices that all of the books are between 50 and 400 pages, and uses the following line of code to generate a histogram showing the lengths: books.hist('Length', unit='page', bins=np.arange(50, 400, 25)) 1 Percent ber bage
8.0
4.0
4.0
5.0 0 50 100 200 250 300 350 Book length (page) Question 2.1 (3 points) What proportion of books are between 200 and 224 pages long (inclusive)? Assign the variable q21_proportion_books_200_224 to your answer. Your answer should be a number between 0 and 1, and should be correct to within .01. In [9]: | q21_proportion_books_200_224 = ... In []: | ok.grade("q2 1"); **Question 2.2** (4 points) Based on the histogram and the information above, which must be true? Assign book_interpretation_choices to an array of your numbered answer(s). 1. There are more long books (200 pages and over) than short books (under 200 pages) in her family's collection. 2. The percentage of books between 125 and 149 pages (inclusive) is about 20% (that is, 19-21%). 3. The books in her family's collection are like a random sample from the population of all books. 4. In her family's collection, nonfiction books are longer on average than fiction books. 5. The mean length is less than the median length. In [13]: book_interpretation_choices = ... **Question 2.3** (2 points) Which of the following is closest to the median length? Assign the variable median_length_choice to either 1, 2, or 3 depending on your choice. 1. 200 2. 260 3. 315 In [14]: median length choice = ... In []: ok.grade("q2_3"); Zeynep finds review data for each book online, and adds it to her books table in a column called Rating. The values in this column are floating point numbers between 1 and 5, indicating the number of "stars" the book got, on average. She computes the following quantities, which are reproduced for you in the cell below. Remember, you can't use the books table, only the quantities defined for you. In [50]: correlation(books, 'Length', 'Rating') Out[50]: 0.5960786546537442 In [51]: lengths = books.column('Length') np.mean(lengths), np.std(lengths) Out[51]: (237.6054054054054, 61.688108039429245) In [52]: ratings = books.column('Rating') np.mean(ratings), np.std(ratings) Out[52]: (2.675278572606432, 0.3234389884574323) In [18]: r = 0.596length mean = 237.605length sd = 61.7 $rating_mean = 2.675$ $rating_sd = 0.323$ **Question 2.4** (2 points) Using linear regression, what would she predict for the average rating (in stars) of a book with 300 pages? Assign the variable book_with_300_pages_avg_rating to your answer. Hint: try to avoid names like correlation, slope, intercept for your variables, since those are the names of functions we've defined for you at the top of the notebook. In [19]: book_with_300_pages_avg_rating = ... book with 300 pages avg rating In []: | ok.grade("q2_4"); **Question 2.5** (2 points) Using linear regression, what would she predict for the length of a book (in pages) with an average rating of 4 stars? Assign the variable book_with_4_stars_length to your answer. Hint: try to avoid names like correlation, slope, intercept for your variables, since those are the names of functions we've defined for you at the top of the notebook. In [22]: book with 4 stars length = ... book_with_4_stars_length In []: | ok.grade("q2_5"); Question 2.6 (2 points) Fill in the blank with the smallest possible value that's guaranteed to be correct. Assign the variable book rating blank to your answer. Without knowing anything else about the distribution of ratings, we can guarantee that 75% of the ratings will be between 2.029 stars and ___ stars. In [25]: book_rating_blank = ... book_rating_blank In []: | ok.grade("q2_6"); **Question 2.7** (4 points) Fill in the blank with the smallest value that's guaranteed to be correct. Assign the variable range_of_accurate_book_predictions to your answer. When predicting average rating from book length, at least 93.75% of the predictions will be correct to within ___ stars of the true value. *Hint*: 93.75% is the same as $1-\frac{1}{16}$. In [28]: range_of_accurate_book_predictions = ... range_of_accurate_book_predictions In []: | ok.grade("q2_7"); **Question 2.8** (3 points) Zeynep uses her data to conduct a hypothesis test, but refuses to tell you any of the details. All she tells you is: Larger values of the test statistic favor the alternative hypothesis. • Her *p*-value cutoff is 0.05. • Her histogram of 5,000 simulated values of her test statistic under the null hypothesis looks like this (assume all values are shown in the histogram): 10 Percent per unit 8 6 2 0 285 290 295 300 305 310 315 Test statistic Based only on this information, which of the following must be true? Assign mystery_hypothesis_test_conclusions to an array with your numbered answer(s). 1. If the test statistic in her data is 310, then the data are more consistent with the alternative hypothesis than the null hypothesis. 2. If the test statistic in her data is 270, then we should conclude the data are more consistent with the null hypothesis than the alternative hypothesis. 3. If this was an A/B test where the test statistic was the difference in the means of two groups, and the test statistic in her data is 312, then she can conclude that there's a causal link between the groups and the value she's measuring for her test statistic. In [31]: | mystery_hypothesis_test_conclusions = ... **Question 3: Gardening** Silla wants to get better at keeping his plants alive. So, he decides to try using a classifier to predict whether the plants into his garden will survive. He spends 8 weeks, from March 17 until May 12, collecting data for each plant on: • Water: a float, the average amount of water he gave it per week (in ounces), • Pot Size: a float, the size of the pot the plant was in (in gallons) • Survived: a boolean, whether or not the plant was alive on May 12 He trains a k-nearest neighbor classifier on his data, and uses k=3. He also draws the following scatterplot, marking plants that survived with blue circles and plants that died with red squares. Assume all his plants are shown here: Survived=True Survived=False 5 Wai 3 2 1 0 3 Pot Size Silla wants to use the classifier to make a predictions. For each of the following questions, you should answer one of: True , if the prediction is that the plant would survive False, if the prediction is that the plant would die, • A string with an explanation if it isn't appropriate to use his classifier to answer the question. If you answer True or False, then you should *not* include any explanation: it will not be graded. For example, if a question asked "How tall will his tallest plant be?", then your answer should be "The classifier can't predict how tall plants are, only whether or not they survive". Question 3.1 (3 points) Suppose Silla had planted an extra plant in a 0.25-gallon pot that received 3.5 ounces of water per week. Would the classifier predict that this plant will survive? Write your answer here, replacing this text. **Question 3.2** (3 points) Silla was taking care of his roommate's plant (which isn't shown in the graph above). It was in a 2.5-gallon pot, and he gave it 5 ounces of water per week. Would the classifier predict that this plant will survive? Write your answer here, replacing this text. **Question 3.3** (3 points) Silla was taking care of his other roommate's plant (which also isn't shown in the graph above). It was in a 6gallon pot, and he gave it 6 ounces of water per week. Would the classifier predict that this plant will survive? Write your answer here, replacing this text. Question 3.4 (3 points) Silla graduates and moves from Berkeley to Miami, Florida. After moving in to his new apartment, he buys some new plants. He puts one in a 1-gallon pot, and gives it 2 ounces of water per week for 8 weeks (from June 1 to July 27). Would the classifier predict that this plant will survive? Write your answer here, replacing this text. **Question 4: Electricity** The table electricity contains data for 200 randomly sampled energy utilities from the US in 2017. In [32]: electricity = Table.read_table('electricity2017_sample.csv').drop('Power in') electricity.show(3) It has the following columns: Name: a string, the name of the utility • **State**: a string, the two-letter abbreviation for the state the utility operates in • **Type**: a string, the type of utility • Residential customers: an int, the number of residential customers the utility serves • **Revenue**: a float, the total revenue for the utility in 2017 measured in thousands of dollars Power generated: a float, the amount of power the utility generated itself (in megawatt-hours) • Power bought: a float, the amount of power the utility bought or exchanged from other utilities (in megawatt-• **Summer demand**: a float, representing peak demand in the summer (in megawatts) • **Winter demand**: a float, representing peak demand in the winter (in megawatts) Question 4.1 (3 points) Janea wants to use this random sample to understand energy utilities more broadly. In particular, she wants to estimate the following quantities from this particular sample: 1. The maximum number of residential customers served by any energy utility in the US 2. The average (mean) demand in the summer across all energy utilities in the US 3. The median revenue of all energy utilities in Hawaii (*Hint*: the two-letter abbreviation for Hawaii is HI). 4. The slope of a linear prediction of winter demand from summer demand, across all energy utilities in the US If she decides to use the bootstrap to construct a 95% confidence interval from this sample, which of the quantities above are good choices for this technique? Assign good_electricity_bootstrap_candidates to an array of your numbered answer(s). In [33]: good electricity bootstrap candidates = ... Question 4.2 (2 points) Javier is on the board of a US energy utility, and wants to predict what his revenue will be based on summer demand. He isn't sure whether linear regression is a good fit: create a plot to help convince him that it is. *Hint*: you shouldn't need more than 3-4 lines to solve this, and our solution took fewer than that. In [34]: # Code to create a plot goes here. **Question 4.3** (5 points) Javier's utility had 90 megawatts of demand this summer. He uses the bootstrap to construct a 95% confidence interval for the prediction of its revenue based on linear regression, and finds that the interval is \$40.6 to \\$45 million dollars (remember, the Revenue column is in thousands of dollars). Which of the following must be true? Assign revenue_demand_prediction_choices to an array of your numbered answer(s). 1. 95% of utilities in the sample had revenue between \$40.6 and \\$45.0 million dollars. 2. In 95% of Javier's bootstrap samples, the height of the regression line at x=90 was between 40600 and 45000 (that is, \$40.6 and \\$45 million dollars). 3. In the population of all US energy utilities, 95% of the utilities with 90 megawatts of summer demand have between \$40.6 and \\$45 million dollars of revenue. 4. There is a 95% chance that Javier's utility's revenue this year will be between \$40.6 and \\$45 million dollars. 5. When drawing a random sample of 200 US energy utilities from the population of all US energy utilities, there is a 95% chance that in the sample, the linear regression prediction for utilities with 90 megawatts of summer demand will be between \$40.6 and \\$45 million dollars. In [35]: revenue_demand_prediction_choices = ... **Question 4.4** (4 points) Javier wants to estimate the difference between the average size (number of customers) of municipal utilities (owned by local governments) and cooperative utilities (owned by customers). In other words, he wants to compute the municipal average number of customers minus the cooperative average number of customers. Help him by completing the function below. For this question, the public tests check for correctness. In [36]: def compute_municipal_cooperative_difference(tbl): Given a table with columns 'Type' and 'Residential customers', returns the difference between the average size of municipal utilities and cooperative utilities. . . . In []: ok.grade("q4 4"); **Question 4.5** (4 points) Javier wonders: if his sample had been different, could he have gotten a different answer? Complete the function below, which computes 1,000 bootstrap samples of the statistic from the previous part and saves them in the array bootstrap utility diff which is returned. In [39]: def bootstrap_u_d(): bootstrap_utility_diff = ... resample = ...bootstrap utility diff = ... return bootstrap utility diff # Don't edit this code below bootstrap utility differences = bootstrap u d() In []: ok.grade("q4 5"); **Question 4.6** (3 points) Javier's favorite number is 92, so he wants to make a 92% confidence interval to estimate the value of the statistic above in the population. Use the bootstrap utility differences to assign bootstrap utility left and bootstrap utility right: In [42]: bootstrap_utility_left = ... bootstrap_utility_right = ... (bootstrap_utility_left, bootstrap_utility_right) In []: ok.grade("q4 6"); **Question 5: The Office** The Data 8 instructors start arguing over one of their favorite TV shows, and decide to resolve their debates using data. They find data on how many words each character speaks in each episode of the show The Office in the table In [46]: office = Table.read table('the office.csv') office.show(3)The table contains the following columns: season: an int, indicating which season of the show episode: an int, indicating which episode within that season the data is for overall_episode: an int, indicating the episode number within the entire show (for example, since the first season has 6 episodes, the first episode of the second season is the seventh episode overall) • speaker: a string, the name of the character (or the string TOTAL, indicating the row contains the total number of words for the episode) num_words: an int, the number of words that character spoke in that episode For example, the second row indicates that the character Michael said 1598 words in the first episode of the first season. (4 points) In order to make line and scatter graphs of how much the characters speak over the course of the show, we'll need to write a function that takes an array of character names and returns a table with one row per overall episode number. It should have one column for the episode numbers, and a column for each character in the array, where the values are the total number of words spoken by that character in that episode. For example, calling get character episode table(make array('Jim', 'Dwight', 'Pam')) should give you a table whose first few rows look like: overall_episode Dwight Jim Pam 295 296 335 2 185 366 113 782 342 218 609 588 236 Question 5.1 This means that Dwight spoke 296 words in the first episode, and Jim and Pam spoke 335 and 295 respectively. (*Hint*: the predicate are.contained_in might be helpful.) In [47]: | def get_character_episode_table(character_names): Takes an array of character names, and returns a table like the one shown above. . . . In []: | ok.grade("q5_1"); Just for fun: Use the function you just wrote to visualize the relationships between how much characters speak, and how that changes over the course of the show. # This cell is for answering the "Just for fun" question above. Nothing in it will be graded. In [50]: Ramesh and Swupnil start arguing over Season 8's popularity, and find the following data to help support their claims: In [51]: office_viewers = Table.read_table('office_wikipedia.csv') office viewers.show(3) The table has the following columns: • **Overall number**: an int, indicating the episode number within the entire show (for example, since the first season has 6 episodes, the first episode of the second season is the seventh episode overall) • Name: a string, the episode name • Writer: a string, the person or people who wrote the script for that episode Millions of viewers: a float, the number of viewers, in millions, who watched that episode when it aired on live TV (on NBC). Just for fun: Which episode had the most viewers and why? # This cell is for answering the "Just for fun" question above. Nothing in it will be graded. In [52]: Question 5.2 (3 points) Using the two tables above (office viewers and office), make a new table that has one row for each episode of the show, and three columns: season, overall episode, and Millions of viewers. The columns may be in any order, but they must contain the correct values as described by their names. In [53]: # The overall and season table was helpful in our solution, # but you don't have to use it. overall_and_season = office.group(['overall_episode', 'season']).select(0, 1) office_season_viewers = ... office_season_viewers.show(3) In []: | ok.grade("q5 2"); Question 5.3 (2 points) Ramesh thinks that Season 8 is significantly less popular (measured by number of viewers) than any other season of the show. Swupnil disagrees. They decide on the following null and alternative hypotheses for their test: **Null hypothesis:** The average viewer count for Season 8 is like the average viewer count for the same number of episodes picked at random from the entire show. Alternative hypothesis: No, the average viewer count for Season 8 is lower. Which hypothesis supports Swupnil's argument? Assign the variable season_8_swupnil to either 'null' or 'alternative': In [57]: season_8_swupnil = ... In []: | ok.grade("q5_3"); Question 5.4 (3 points) Based on the null and alternative hypothesis above, describe a test statistic they could use to conduct a hypothesis test: Write your answer here, replacing this text. **Question 5.5** (3 points) They decide on a p-value cutoff of 0.01, and carry out the test using a correct test statistic. They obtain a p-value of 0. Based only on the information provided to you, which must be true? Assign season_8_test_conclusions to an array with your numbered answer(s). 1. If the null hypothesis were true, the probability of observing a result that supports the alternative hypothesis is 0.01.2. If the alternative hypothesis were true, the probability of observing their test statistic is 1. 3. The data support the alternative hypothesis. 4. The data support the null hypothesis. In [61]: | season_8_test_conclusions = ... **Question 6: Actors** Recall the actors table from lecture: In [62]: actors = Table.read_table('actors.csv').where('Number of Movies', are.above(10)) actors.show(3)You can find a description of the columns in Chapter 7 of the textbook. Note that just like in lecture and the textbook, we've removed Anthony Daniels since he's an outlier. For this question, we'll focus on using the box office gross from the actor's top movie (that is, the Gross column) to predict the total box office receipt from all the actor's movies (that is, the Total Gross column). Recall that when learning about linear regression, we used root mean square error (RMSE). Here's a function that computes it: In [63]: # You don't have to do anything in this cell: this is the same function that we # defined in lecture, reproduced just so you can see it. def rmse(predictions, actual_values): Takes an array of predictions and an array of actual observed values, and returns the root mean squared error. error = actual_values - predictions squared_error = error ** 2 mean_squared_error = np.mean(squared_error) return np.sqrt(mean_squared_error) Question 6.1 (3 points) Suppose instead of using RMSE, we decide to use a weighted error. We want our prediction to be more accurate for actors who've been in fewer movies, and we don't care as much about actors who've been in lots of movies. So, we're going to compute the *weighted RMSE* for a prediction line, using this procedure: 1. Compute the squared error, just like before. 2. For each actor, divide the squared error by the weights (i.e., the number of movies the actor has been in). The result of this division will be called the weighted squared error for each actor. 3. Compute the average weighted squared error. 4. Take the square root. Complete the function below to implement the procedure described here. In [64]: **def** weighted rmse(predictions, actual values, weights): error = actual_values - predictions In []: | ok.grade("q6_1");

Question 6.2 (4 points) Complete the weighted_prediction_error function below, which takes in the slope and intercept for any line, and computes the weighted RMSE for predicting these actors' total box office gross using that line. As a reminder, the line predicts total box office gross (the Total Gross column) from the box office gross from the actor's top movie (that is, the Gross column), and uses the number of movies the actor has been in (the Number of Movies column) for the weights in computing the weighted RMSE. *Hint*: your solution should use the weighted_rmse function you defined above. In [67]: def weighted_prediction_error(any_slope, any_intercept): In []: | ok.grade("q6 2"); Question 6.3 (4 points) Find the slope and intercept of the line whose predictions have the smallest weighted RMSE. In [70]: best weighted slope = ... best_weighted_intercept = ... In []: | ok.grade("q6_3"); The code below plots the original linear regression line (blue) and the line with the smallest weighted RMSE (orange) using your answer from above. The size of each dot represents the number of movies that actor has been in. If you see errors trying to run this cell, it might be because you redefined the slope and intercept functions from the beginning of this notebook. Try changing the variable names you use, re-running the cell at the top, and running this cell again. In [74]: # Do not change any of the code in this cell. # You don't have to answer any questions here or understand how this cell works: # it just shows you the result of your work. x = actors.column('Gross') y = actors.column('Total Gross') weights = actors.column('Number of Movies') $x_plot = make_array(169, 937)$ regression_line = x_plot * slope(actors, 'Gross', 'Total Gross') + intercept(actors, 'Gross', 'Total l Gross') weighted_line = x_plot * best_weighted_slope + best_weighted_intercept plots.figure() plots.scatter(x, y, s=(weights ** 2) / 30, c='black') plots.plot(x_plot, regression_line, color='tab:blue') plots.plot(x_plot, weighted_line, color='tab:orange') plots.xlabel('Gross') plots.ylabel('Total Gross'); **Question 7: US Counties and Food Deserts** A food desert is an area that has limited access to nutritious food. In this question, we'll look at food access data for US counties from 2010 to 2015. In [75]: food_access = Table.read_table('food_access_2010_2015.csv') food_access.show(3) The food_access table contains one row for each county in the 50 states of the US. It has the following columns: • county: a string, the name of the county **population**: an int, the number of people living in that county • **state**: a string, the two-letter abbreviation for the state the county is in housing_units: an int, the number of housing units available in the county • urban_pct : a float between 0 and 100, the percentage of housing units that are urban in the county • low access 20: an int, the number of people in that county living more than twenty miles away from their nearest grocery store • low access 10: an int, the number of people in that county living more than ten miles away from their nearest grocery store (includes people living more than twenty miles away) • low_access_1: an int, the number of people in that county living more than one mile away from their nearest grocery store (includes people living more than 10 and 20 miles away) • carless_pct: a float between 0 and 100, the percentage of people in that county who don't have a car in their household In this question, you'll also work with the states table: In [76]: states = Table.read_table('states.csv') states.show(3)For each state, the table has: State: a string, the name of the state • State Code: a string, the two-letter abbreviation for the state Region: a string, which region the state belongs to • **Division**: a string, which division the state belongs to If you're curious about what the regions and divisions are (and what the difference between them is), see this map from the US Census Bureau. Part 1: Table manipulation For each of the following questions, write Python code that computes (or draws) the specified quantity, table, or graph. Question 7.1.1 (2 points) The number of people living more than 10 miles away from a grocery store in Alameda County, where UC Berkeley is. In [77]: | num_people_more_than_10_miles_away_alameda_county = ... num_people_more_than_10_miles_away_alameda_county In []: | ok.grade("q7_1_1"); Just for fun: Same as the above question, but for any other county in California that you've lived in or spent time in. Do the results surprise you? In [80]: # This cell is for answering the "Just for fun" question above. Nothing in it will be graded. Question 7.1.2 (3 points) The county in California that has the most people living far away (>10 miles) from a grocery store. The abbreviation for California is CA. In [81]: california_county_most_food_deserted = ... california county most food deserted In []: ok.grade("q7_1_2"); Question 7.1.3 (6 points) A table biggest_state_for_each_division with the largest state (by population) in each **division** (not region). It should have **one row for each division**, and two columns: one with the name of the division, and one with the name (not the two letter abbreviation) of the most populated state in that division. The names of the columns don't matter, but they must be in that order (division column first, state column second). For example, the most populated state in the Pacific division is California, so one of the rows in your table should have as its first item the string Pacific and as its second item the string California. (*Hint*: you may find the first function helpful: we've defined it for you here.) In [84]: def first(arr): return arr.item(0) state populations = ... biggest_state_for_each_division = ... biggest_state_for_each_division In []: ok.grade("q7 1 3"); Part 2: relationships between quantities For the rest of this question, instead of looking at the number of people without access to a grocery store, we'll look at the percentage. We'll use the following terms: • low access at 1 mile: the percentage of people in a county whose nearest grocery store is at least 1 mile away • low access at 10 miles: the percentage of people in a county whose nearest grocery store is at least 10 miles away In [88]: food access with pcts = food access.with columns('low_access_1_pct', food_access.column('low_access_1') / food_access.column('population') * 100 'low_access_10_pct', food_access.column('low_access_10') / food_access.column('population') * 1 00, 'low_access_20_pct', food_access.column('low_access_20') / food_access.column('population') * 1 00,).drop('low_access_1', 'low_access_10', 'low_access_20') food_access_with_pcts.show(3) Question 7.2.1 (2 points) Create one plot that you'd use to answer the following question. You'll use your plot to answer the related multiple choice question below, but you do **not** have to answer the question directly, only create a plot. Is there any association between the percentage of urban housing and low access at **10 miles** in US counties? In [89]: # Code to generate your plot goes here Question 7.2.2 (3 points) Based only on the data in the table (and the graph you created using that data), which must be true? Assign low_access_10_urban_choices to an array with your numbered answer(s). 1. Knowing the percentage of urban housing in a county cannot help us predict low access at 10 miles. 2. There is a strong linear association between percentage of urban housing and low access at 10 miles. 3. The distribution of low access at 10 miles is different between non-urban counties (<5% urban housing) and other counties. 4. For counties that are very urban (>80% urban housing), more than half of their population is within 10 miles of a grocery store. In [90]: low_access_10_urban_choices = ... Question 7.2.3 (2 points) Create one plot that you'd use to answer the following question. You do not have to answer the question directly, only create a plot. Is there any association between the percentage of urban housing and low access at 1 mile in US counties? In [91]: # Code to generate your plot goes here Question 7.2.4 (2 points) What is the correlation between percentage of urban housing and low access at 1 mile? (Hint: You may find it helpful to use some of the functions defined in the "Useful functions" section near the top of this notebook). In [92]: correlation_urban_low1 = ... correlation_urban_low1 In []: | ok.grade("q7_2_4"); **Question 8: Likes on Instagram** Natalia, an engineer at Instagram, conducts a randomized controlled experiment to evaluate whether social media anxiety is **reduced** when users can't see the number of likes shown on a post. When the next Instagram software update ships, she randomly assigns users in the city of Berkeley to two groups: Group A: 10,000 users who can no longer see the number of likes, and • **Group B**: 8,000 users who can still see the number of likes. One month after shipping the new update, she measures each user's happiness and general sentiment towards Instagram with a survey. She wants to test the claim that the treatment, hiding the number of likes, increases users' happiness scores. She prepares a table instagram_users, containing 18,000 rows, one for each user in her experiment. She can't share the full dataset with you since it's proprietary, but she chooses 3 random rows to show you, just so you can see what it looks like: In [95]: instagram_users = Table.read_table('instagram_sample.csv') instagram_users • **sex**: a string, the user's sex (Male or Female) • age : an int, the user's age • group: a string, the user's group for the experiment (A or B) **happiness**: an int, the user's happiness score from the survey (between 0-100) sentiment: a string, the user's sentiment from the survey (Positive, Neutral, or Negative) Natalia's null hypothesis is that there is no difference in average happiness between people who see the number of likes compared to those who don't, and that any difference observed in the sample is due to chance. For her test statistic, Natalia decides to use the difference between the average happiness scores of Group A and Group B (that is, the average of Group A minus the average of Group B). Help her come up with an alternative hypothesis: **Question 8.1** (3 points) State an alternative hypothesis that she should use for her test. Write your answer here, replacing this text. **Question 8.2** (6 points) Natalia asks her coworker for help simulating one value of the test statistic under the null hypothesis, but the code he gives her has some mistakes. Fix the code below so that it correctly computes the test statistic under the null hypothesis. In [96]: # This code contains mistakes that you need to fix. def compute instagram test statistic(): # Shuffle the data shuffled groups = instagram users.take('group').sample(with replacement = False) users = instagram_users.append('shuffled_groups', shuffled_groups.column(1)) # Two averages mean_A = np.average(users.where('shuffled_groups', 'A').select('group')) mean_B = np.average(users.where('shuffled_groups', 'B').select('group')) # Test statistic return mean_A - mean B **Question 8.3** (4 points) In the question above, why do we shuffle the data? Choose all that apply. Assign instagram shuffling reasons to an array with your numbered answer(s). 1. Under the null hypothesis, the label of being in group A or group B doesn't matter. 2. We want to randomize treatment and control to establish causation. 3. We want to simulate two groups of people whose expected happiness is identical under the null hypothesis. 4. We want to ensure that the users in the experiment are selected randomly. In [97]: instagram shuffling reasons = ... instagram_shuffling_reasons Natalia fixes the code, and simulates 10,000 values of the test statistic under the null hypothesis. She shows them in the histogram below. You should assume that the histogram shows all of the simulated values. 3.5 Percent per unit 3 2.5 2 1.5 1 0.5 Difference between avg happiness **Question 8.4** (4 points) Based on this histogram, which of the following must be true? Assign instagram_null_simulation_choices to an array with your numbered answer(s). 1. Seeing the number of likes on a post has a positive effect on user happiness. 2. Seeing the number of likes on a post has no effect on user happiness. 3. If the p-value cutoff for the test is 0.01 and the observed test statistic is 35, then we should conclude the data are more consistent with the alternative hypothesis. 4. If the p-value cutoff for the test is 0.05 and the observed test statistic is 10, then we should conclude the data are more consistent with the null hypothesis. 5. Seeing the number of likes on a post has a negative effect on user happiness. In [98]: instagram_null_simulation_choices = ... **Question 8.5** (5 points) Suppose Natalia tells you the observed test statistic was 42. Which conclusion(s) are the data consistent with? Assign instagram_test_conclusion_choices to an array with your numbered answer(s). 1. The difference in happiness scores between group A and group B is due to chance alone 2. The treatment has a negative association with happiness scores 3. The treatment has a positive effect on happiness scores 4. The treatment increases users' happiness scores by 42 points. 5. The data support the null hypothesis. 6. The data support the alternative hypothesis. 7. There isn't enough information to make a conclusion of any kind. In [99]: instagram_test_conclusion_choices = ... For the remainder of this question, instead of testing for an increase in happiness scores, Natalia would like to test whether the treatment (i.e. hiding the number of likes) changes users' sentiment: that is, she's only interested in seeing whether the sentiment is significantly different. Recall from the table above that sentiment is measured as either Positive, Neutral, or Negative. She comes up with the following alternative hypothesis: The distribution of user sentiment between treatment (hiding the number of likes) and control (showing the number of likes) is different. Provide a null hypothesis and a test statistic she could use for her test. **Question 8.6** (3 points) Null hypothesis: Write your answer here, replacing this text. **Question 8.7** (3 points) What would be a valid test statistic to tell the two hypotheses above apart? Write your answer here, replacing this text. **Question 8.8** (3 points) Natalia and her coworker are deciding whether to roll this feature out worldwide. Her coworker wants to collect more data, and develops a sentiment analysis algorithm that analyzes the content of users' posted images and captions to automatically determine their sentiment (positive, negative, or neutral). Describe, in two sentences or less, any privacy concerns that you would have around the collection of this data. *Hint*: this question is more open-ended than most of the rest of this exam: there isn't only one correct answer. Write your answer here, replacing this text. **Question 9: Bread** Part 1: Frozen slices Ilin bakes too much bread, so she decides to slice it and put it in her freezer. She has three loaves of bread: one rye, one sourdough, and one multigrain. She slices each loaf: the rye bread has 8 slices, the sourdough bread has 5 slices, and the multigrain bread has 12 slices. In each loaf, two of the slices are "heels" (the slice at each end). She puts all the slices into a giant bag and mixes them up. Each day, she pulls out one slice at random, toasts it, and eats it. Question 9.1.1 (2 points) What is the probability that her first slice is sourdough? In [100]: | slice_sourdough = ... slice_sourdough In []: | ok.grade("q9_1_1"); Question 9.1.2 (3 points) What is the probability that in her first two slices, she doesn't get any sourdough? In [103]: no_sourdough = ... no_sourdough In []: | ok.grade("q9_1_2"); Question 9.1.3 (3 points) What is the probability that the first three slices are all heels? In [106]: three heels = ... three heels In []: ok.grade("q9 1 3"); Question 9.1.4 (3 points) The first slice she pulls out is a heel. What is the probability that it is sourdough? In [109]: sourdough_heel = ... sourdough_heel In []: | ok.grade("q9_1_4"); **Part 2: Quality control** Jin is in charge of quality control at a very large bakery that makes tens of thousands of loaves of bread each day. Each loaf is advertised as weighing 500 grams, but the exact weight varies a little bit from loaf to loaf. Jin has had complaints from customers that their bread weighed less than 500 grams. He speaks with the bakers, who assure him that the average loaf weighs 500 grams, and that the standard deviation of the loaf weights is 5 grams. They show him the following histogram of distribution of weights that the loaves are supposed to follow. According to the bakers, the distribution of weights in the population of loaves is as follows: 0.12 0.10 0.08 0.06 0.04 0.02 0.00 480 490 510 500 520 Question 9.2.1 (4 points) To verify this, he plans to hire 10,000 auditors to come into the bakery, and each one will randomly sample 100 loaves and weigh them. Each auditor will then tell Jin the average weight of the 100 loaves that they weighed. Assuming the population distribution information above is correct, about how many of the 10,000 auditors will find an average weight below 499 grams? Your response should be a number between 0 and 10,000. *Hint*: you should compute your answer using arithmetic from the information given, not using a simulation. In [112]: auditors avg below 499 = ... auditors_avg_below_499 In []: ok.grade("q9 2 1"); Question 9.2.2 (4 points) When Jin discusses this idea with his peers, they tell him it's much too expensive. Instead, they suggest hiring one auditor, and asking the auditor to collect a slightly larger random sample of 225 loaves. Assuming the population distribution information above is correct, approximately what is this probability that the auditor will find an average weight below 499 grams? Your response should be a number between 0 and 1. You should compute your answer using arithmetic from the information given, not using a simulation. In [115]: prob auditors avg below 499 = ... prob_auditors_avg_below_499 In []: ok.grade("q9_2_2"); 2. Submission Once you're finished, select "Save and Checkpoint" in the File menu and then execute the submit cell below. The result will contain a link that you can use to check that your assignment has been submitted successfully. If you submit more than once before the deadline, we will only grade your final submission. If you mistakenly submit the wrong one, you can head to okpy.org and flag the correct version. To do so, go to the website, click on this assignment, and find the version you would like to have graded. There should be an option to flag that submission for grading! In [119]: = ok.submit() In [120]: # For your convenience, you can run this cell to run all the tests at once! print("Running all tests...") _ = [ok.grade(q[:-3]) for q in os.listdir("tests") if q.startswith('q') and len(q) <= 10]</pre> print("Finished running all tests.")