# Data C100/200- Midterm

#### Spring 2025

| Name:                                     |              |               |
|---|--------------|---------------|
| Email:                                    |              | @berkeley.edu |
| Student ID:                               |              |               |
| Name and SID of the person on your left:  |              |               |
| Name and SID of the person on your right: |              |               |
| Exam Room:                                | Seat Number: |               |

#### **Instructions:**

This exam consists of **45 points** spread out over **4 questions** and the **Honor Code certification**. The exam must be completed in **110 minutes** unless you have accommodations supported by a DSP letter. Note that you should:

- Each true/false question and multiple choice question has **exactly one** correct answer. Please **fully** shade in the circle to mark your answer.
- Blank answers and incorrect answers are graded identically, so it's in your best interest to answer every question.
- For all math questions, **please simplify your answer**. Please also **show your work** if a large box is provided.
- For all coding questions, you may use commas and/or one or more function calls in each blank.
- You MUST write your Student ID number at the top of each page.
- You should not use a calculator, scratch paper, or notes you own other than the reference sheets distributed at the beginning of the exam.

For all Python questions, you may assume Pandas has been imported as pd, NumPy as np, the Python RegEx library as re, matplotlib.pyplot as plt, and seaborn as sns.

#### Honor Code [1 Pt]:

As a member of the UC Berkeley community, I act with honesty, integrity, and respect for others. I am the person whose name is on the exam, and I completed this exam in accordance with the Honor Code.

Signature: \_

This page has been intentionally left blank.

### 1 Awesome DATA 100 Staff [19 Pts]

At the end of the semester, Data 100 instructors analyze staff performance using a DataFrame called performance. The columns of performance are described below:

- name: Name of the staff member (type = str).
   Note: No two staff members have exactly the same name.
- role: One of three possible staff roles: "GSI", "TA", or "Tutor" (type = str).
- is\_graduating: True if the staff member is graduating this semester and False otherwise (type = bool).
- disc\_day: The staff member's discussion day. One of three values: "W" (Wednesday), "Th" (Thursday), or NaN if the staff member does not hold any discussions (type = str).
- oh\_tickets: Total office hours (OH) tickets the staff resolved during the semester (type = np.int64).
- ed\_hours: Total hours the staff spent resolving questions on EdStem during the semester (type = np.float64).
- grad\_day: The staff member's final day as an enrolled UC Berkeley student, in "yyyy-mm-dd" format (type = str).

The first five rows of performance are shown below:

|   | name      | role  | is_graduating | disc_day | oh_tickets | ed_hours | grad_day   |
|---|-----------|-------|---------------|----------|------------|----------|------------|
| 0 | Dan       | TA    | False         | W        | 111        | 11.9     | 2026-05-16 |
| 1 | Gisella   | Tutor | True          | NaN      | 117        | 13.1     | 2025-05-17 |
| 2 | Steven    | Tutor | False         | NaN      | 109        | 14.6     | 2026-08-15 |
| 3 | Malavikha | TA    | True          | NaN      | 122        | 19.0     | 2025-05-17 |
| 4 | Rose      | GSI   | True          | W        | 4          | 23.5     | 2025-05-17 |

- (a) For each performance column, choose the best variable type.
  - (i) [0.5 Pts] The best variable type for the column name is:
    - Qualitative ordinal
    - **Qualitative nominal**
    - Quantitative
  - (ii) [0.5 Pts] The best variable type for the column oh\_tickets is:
    - O Qualitative ordinal
    - O Qualitative nominal
    - **Quantitative**

(b) [0.5 Pts] What is the granularity of the performance DataFrame? Answer with a brief sentence or phrase.

Solution: Staff member or name

(c) [1.5 Pts] Instructors want to visualize the distribution of ed\_hours. Mark True if the plot type is appropriate for this visualization task, and False otherwise.

| $\bigcirc$ | True | $\bigcirc$ | False | Contourplot |
|------------|------|------------|-------|-------------|
| $\bigcirc$ | True | $\bigcirc$ | False | Histogram   |
| $\bigcirc$ | True | $\bigcirc$ | False | KDE Plot    |

**Solution:** The first option is incorrect because contourplots plots visualize the joint distribution of two quantitative variables by representing density through contour lines. For single, continuous numeric variables like ed\_hours, a histogram or KDE plot is more suitable.

- (d) [1.5 Pts] Instructors want to visualize the relationship between oh\_tickets and ed\_hours. Mark True if the plot type is appropriate for this visualization task, and False otherwise.
  - True False Jointplot
     True False Overlaid histograms
     True False Hexplot

**Solution:** The first option is correct because jointplots plots effectively visualize the joint distribution of two quantitative variables. The second option is incorrect, as overlaid histograms display the distribution of individual variables but do not show the relationship between them. The third option is correct for the same reason as the first- hexplots visualize the joint distribution by grouping data points into hexagonal bins, with shading indicating density.

(e) [1 Pt] Instructors compare oh\_tickets for Tutors and TAs using a boxplot shown below:



Select True or False for the statements below.

○ True ○ False The median OH tickets resolved is higher for Tutors than for TAs.

○ True ○ False About 75% of TAs each resolved more OH tickets than about 75% of Tutors.

**Solution:** The first option is incorrect because the median number of OH tickets resolved by TAs, represented by the line inside the TAs' boxplot above, is higher than that of Tutors'. The second option is correct because the 75th percentile (Q3) of the Tutor distribution aligns roughly with the 25th percentile (Q1) of the TA distribution.

The first five rows of performance are shown again here for your convenience:

|   | name      | role  | is_graduating | disc_day | oh_tickets | ed_hours | grad_day   |
|---|-----------|-------|---------------|----------|------------|----------|------------|
| 0 | Dan       | ТА    | False         | W        | 111        | 11.9     | 2026-05-16 |
| 1 | Gisella   | Tutor | True          | NaN      | 117        | 13.1     | 2025-05-17 |
| 2 | Steven    | Tutor | False         | NaN      | 109        | 14.6     | 2026-08-15 |
| 3 | Malavikha | ТА    | True          | NaN      | 122        | 19.0     | 2025-05-17 |
| 4 | Rose      | GSI   | True          | W        | 4          | 23.5     | 2025-05-17 |

(f) Fill in the blanks to create a side-by-side violin plot to visualize the distribution of ed\_hours for each combination of role type and graduation status, as shown below:



Note: The order of the arguments does not matter as long as each is named.

```
sns.violinplot(
data=performance[performance["role"].isin(["TA","Tutor"])],
_____(i)____,
____(ii)____,
____(iii)____,
```

(i) [0.5 Pts] Fill in blank (i):

**Solution:** x = "is\_graduating"

(ii) [0.5 Pts] Fill in blank (ii):

Solution: y = "ed\_hours"

(iii) [0.5 Pts] Fill in blank (iii):

Solution: hue = "role"

0 1

The first five rows of performance are shown again here for your convenience:

| name    | role  | is_graduating | disc_day | oh_tickets | ed_hours | grad_day   |
|---------|-------|---------------|----------|------------|----------|------------|
| Dan     | ТА    | False         | W        | 111        | 11.9     | 2026-05-16 |
| Gisella | Tutor | True          | NaN      | 117        | 13.1     | 2025-05-17 |
|         |       |               |          |            |          |            |

| 2 | Steven    | Tutor | False | NaN | 109 | 14.6 | 2026-08-15 |
|---|-----------|-------|-------|-----|-----|------|------------|
| 3 | Malavikha | TA    | True  | NaN | 122 | 19.0 | 2025-05-17 |
| 4 | Rose      | GSI   | True  | W   | 4   | 23.5 | 2025-05-17 |

(g) [2 Pts] Write a single line of code to return a String with the name of the staff member who has resolved the highest number of oh\_tickets.

Note: There is exactly one staff member with the highest count.

Answer: (\_\_\_\_\_

```
Solution: There are many possible solutions!
Sample Answers:
  • performance.sort_values('oh_tickets', ascending=False)
    .iloc[0]['name']
  • performance.sort_values('oh_tickets', ascending=False)
    .iloc[0,0]
  • performance.sort_values('oh_tickets', ascending=False)
    .iloc[0].iloc[0]
  • performance.sort_values('oh_tickets', ascending=False)
    .['name'].values[0]
    Incorrect:
  • performance.sort_values('oh_tickets', ascending=False)
    ['name'][0] - this gets "label 0" not position 0.
Using functions not in the Reference Sheet:
  • performance.loc[performance['oh_tickets']
    .idxmax(), 'name']
  • performance.nlargest(1, 'oh_tickets')['name'].values[0]
```

```
performance[performance['oh_tickets'] ==
performance['oh_tickets'].max()]['name'].values[0]
performance.iloc[performance['oh_tickets']
```

```
.idxmax()]['name']
```

(h) Instructors want to assess office hours and EdStem performance by staff role.

Assign staff\_stats to a DataFrame where role is the index, and the values are the **minimum** number of oh\_tickets and **maximum** number of ed\_hours for staff members in each role category. Fill in the blanks to achieve this.

Note: The resulting DataFrame should contain only two columns. The order of the columns does not matter.

```
staff_stats = performance.groupby(_____(i)____)
.agg( {_____(ii)____, ____(iii)____})
```

(i) [0.5 Pts] Fill in blank (i):

Solution: "role"

(ii) [0.5 Pts] Fill in blank (ii):

```
Solution: "oh_tickets": "min" or np.min or min
```

(iii) [0.5 Pts] Fill in blank (iii):

Solution: "ed\_hours": "max" or np.max or max

The first five rows of performance are shown again here for your convenience:

|   | name      | role  | is_graduating | disc_day | oh_tickets | ed_hours | grad_day   |
|---|-----------|-------|---------------|----------|------------|----------|------------|
| 0 | Dan       | ТА    | False         | W        | 111        | 11.9     | 2026-05-16 |
| 1 | Gisella   | Tutor | True          | NaN      | 117        | 13.1     | 2025-05-17 |
| 2 | Steven    | Tutor | False         | NaN      | 109        | 14.6     | 2026-08-15 |
| 3 | Malavikha | ТА    | True          | NaN      | 122        | 19.0     | 2025-05-17 |
| 4 | Rose      | GSI   | True          | W        | 4          | 23.5     | 2025-05-17 |

(i) [1.5 Pts] Instructors want to recognize top performers on EdStem who **do not hold discussion sections**. A top-performer is any staff who spent **at least 15 hours** on EdStem.

Note: For staff without discussion sections, disc\_day is NaN.

Fill in the following blank to assign top\_ed to a modified version of the performance DataFrame that includes all rows corresponding to these staff members:

top\_ed = performance[\_\_\_\_\_(A)\_\_\_\_]

Fill in blank (A) below.

Answer: (\_\_\_\_\_

```
Solution: (performance["ed_hours"] >= 15) &
(performance["disc_day"].isna()) or
(pd.isnull(performance["disc_day"]))
(~pd.notna(performance["disc_day"]))
```

Note: Checking == np.nan is wrong because this doesn't compare equal to itself.

(j) Instructors change their mind about which rows of performance to keep. They assign top\_role to a modified performance DataFrame that contains rows where the corresponding role is classified as top-performing. A role is top-performing if the average ed\_hours for that role is at least 15. Fill in the blanks to achieve this.

```
top_role = performance.groupby(____(A)___).___(B)____
```

(i) [0.5 Pts] Fill in blank (A):

Solution: "role"

(ii) [1.5 Pts] Fill in blank (B):

Solution: filter(lambda sf: sf["ed\_hours"].mean() >= 15)
or filter(lambda sf: np.mean(sf["ed\_hours"]) >= 15)

The first five rows of performance are shown again here for your convenience:

|   | name      | role  | is_graduating | disc_day | oh_tickets | ed_hours | grad_day   |
|---|-----------|-------|---------------|----------|------------|----------|------------|
| 0 | Dan       | ТА    | False         | W        | 111        | 11.9     | 2026-05-16 |
| 1 | Gisella   | Tutor | True          | NaN      | 117        | 13.1     | 2025-05-17 |
| 2 | Steven    | Tutor | False         | NaN      | 109        | 14.6     | 2026-08-15 |
| 3 | Malavikha | ТА    | True          | NaN      | 122        | 19.0     | 2025-05-17 |
| 4 | Rose      | GSI   | True          | W        | 4          | 23.5     | 2025-05-17 |
|   |           |       |               |          |            |          |            |

(k) Instructors want to find out how much time graduating staff members spend on EdStem. To investigate, they created the DataFrame shown below, where they aggregated ed\_hours using the median and imputed Null values with 0.

| role          | GSI  | TA   | Tutor |  |  |  |
|---------------|------|------|-------|--|--|--|
| is_graduating |      |      |       |  |  |  |
| False         | 13.3 | 11.9 | 12.35 |  |  |  |
| True          | 23.5 | 17.8 | 12.30 |  |  |  |

Fill in the blanks to replicate the above DataFrame.

| <pre>performance.pivot_table(</pre> | (i),   |
|-------------------------------------|--------|
|                                     | (ii),  |
|                                     | (iii), |
|                                     | (iv),  |
|                                     | (V))   |

Note: The order of the arguments does not matter as long as each is named.

(i) [0.5 Pts] Fill in blank (i):

Solution: index = "is\_graduating"

(ii) [0.5 Pt] Fill in blank (ii):

```
Solution: columns = "role"
```

(iii) [0.5 Pts] Fill in blank (iii):

**Solution:** values = "ed\_hours"

(iv) [0.5 Pts] Fill in blank (iv):

**Solution:** aggfunc = "median" or np.median

(v) [0.5 Pts] Fill in blank (v):

**Solution:** fill\_value = 0

The first five rows of performance are shown again here for your convenience:

|   | name      | role  | is_graduating | disc_day | oh_tickets | ed_hours | grad_day   |
|---|-----------|-------|---------------|----------|------------|----------|------------|
| 0 | Dan       | ТА    | False         | W        | 111        | 11.9     | 2026-05-16 |
| 1 | Gisella   | Tutor | True          | NaN      | 117        | 13.1     | 2025-05-17 |
| 2 | Steven    | Tutor | False         | NaN      | 109        | 14.6     | 2026-08-15 |
| 3 | Malavikha | ТА    | True          | NaN      | 122        | 19.0     | 2025-05-17 |
| 4 | Rose      | GSI   | True          | W        | 4          | 23.5     | 2025-05-17 |

(1) Finally, the instructors want to remove all rows where any value is missing and retain only those corresponding to staff members graduating in 2026.

Follow the steps and fill in the blanks to achieve this:

# Step 1: Remove all rows with missing values (NaNs).

perf\_no\_missing = performance.\_\_\_\_(i)\_\_\_\_\_

# Step 2: Extract year from the "grad\_day" column and convert to type "int".

perf\_no\_missing["year"] = \_\_\_\_(ii)\_\_\_\_\_

# Step 3: Filter the performance DataFrame to keep staff members who are graduating in 2026.

filtered\_perf = perf\_no\_missing[\_\_\_\_(iii)\_\_\_]

(i) [0.5 Pts] Fill in blank (i):

Solution: dropna(), optimally can specify axis=0 as an argument

(ii) [1 Pt] Fill in blank (ii):
Note: For Step 2, you <u>must</u> use string slicing to complete the task.
You may <u>not</u> use the . dt accessor.

```
Solution: perf_no_missing["grad_day"].str[:4] or
.str.split("-").str[0] with .astype(int) or astype("int").
.int() is invalid.
```

(iii) [1 Pt] Fill in blank (iii):

**Solution:** perf\_no\_missing["year"]==2026

## 2 CHARprinter's Intro-Spection [6 Pts]

Sabrina Charprinter obtained text from the "Staff" page of the Data 100 website. She wants to use RegEx to extract certain pieces of information (for her upcoming Data 100 parody song).

**Note:** For all parts, you will only need to consider the example strings given to you. You may assume that these examples cover all edge cases.

(a) [3 Pts] Sabrina found metadata in each staff introduction to process as strings.

Suppose Sabrina has already created a pattern and runs the code below. The character "\_" represents a single space, and you may use it in your response.

```
metadata_=_"Name:_Oski;_Age:_159;_Courses_taken:_Data8,
_Data100,_CS61A,_Stat134,_Data88s;_Phone Number:
_555-100-5555;_Likes:_[Bears,_Strawberries,_Data Science]"
```

```
pattern = r'' + : ([^A-Za-z ]+)''
```

```
re.findall(pattern, metadata)
```

List matches in the order returned by re.findall(pattern, metadata), with the first match next to Match 1.

Note: You may have less than 6 total matches. For any unmatched slots, write No Match instead.

| Match | 1: | ۲          |   |
|-------|----|------------|---|
| Match | 2: | ۲          |   |
| Match | 3: | ۱ <u> </u> | • |
| Match | 4: | 1          | • |
| Match | 5: | ۲          | • |
| Match | 6: | ۱ <u> </u> | ' |

#### Solution:

```
Match 1: '159;'
Match 2: '555-100-5555;'
Match 3: '['
Match 4: 'No Match'
Match 5: 'No Match'
Match 6: 'No Match'
```

Lenient grading: If students forget to specify No Match (and leave the last 3 Matches blank), it is a valid solution.

- (b) Help Sabrina create a RegEx pattern to extract the course IDs for all Data Science courses when running the code below. The IDs appear immediately after the subject name "Data" (with no spaces), and they follow these rules:
  - They consist of exactly three digits, or
  - They consist of one or two digits, optionally followed by a letter from the set ["c", "s", "x"].

For example, the output of running the following code block should be ['8', '100', '88s'].

```
metadata_=_"Name:_Oski;_Age:_159;_Courses_taken:_Data8,
_Data100,_CS61A,_Stat134,_Data88s;_Phone Number:
_555-100-5555;_Likes:_[Bears,_Strawberries,_Data Science]"
```

pattern = r"Data(\_\_\_(i)\_\_\_|\_\_(ii)\_\_\_)"

```
re.findall(pattern, metadata)
```

(i) [1 Pt] Fill in blank (i):

**Solution:**  $d{3} \text{ or } d{3} \text{ b or } d/d/d$ 

(ii) [2 Pts] Fill in blank (ii):

Solution: \d{1,2}[csx]? or \d\d?[csx]? or \d\d?[csx]?\b. \* may be used in place of ? in all solutions listed above.

Full points if conditions are correct but the order is switched.

# 3 Bay Area Rapid Studies (BARS) [12.5 Pts]

Rachel, a data scientist at the UC Berkeley Transportation Department, oversees the **BayPass Program**. The BayPass Program is a study that analyzes the effects of giving UC Berkeley students free access to all Bay Area transportation services through a BayPass transit card.

- (a) Rachel selects 12,000 distinct UC Berkeley students uniformly at random from the UC Berkeley enrollment database and gives them a BayPass transit card. She collects usage data from these cards for analysis.
  - (i) [0.5 Pts] What is the population of interest?
    - **UC Berkeley students**
    - O BayPass card holders
    - UC Berkeley students who use Bay Area transit services
    - People who use Bay Area transit services

**Solution:** Rachel wants to study "effects of providing *UC Berkeley students*" with...", so the population of interest is all UC Berkeley students.

(ii) [0.5 Pts] What is the sampling frame in Rachel's study?

**Solution:** All UC Berkeley students (enrolled currently)

- (b) [2 Pts] Rachel invites UC Berkeley students through Data 100 EdStem to take a survey about their experiences with Bay Area public transportation for 5% midterm extra credit. Which statements about her sampling process are true?
  - True
     False
     Since Rachel provided a generous incentive, her survey results will no longer be affected by non-response bias.
     True
     False
     Rachel's sampling method may suffer from selection bias.
     True
     False
     Rachel's sampling method may suffer from response bias.
     True
     False
     The respondents are guaranteed to be representative of the target population.

**Solution:** A is incorrect because it is possible to incur non-response bias for reasons other than a lack of incentive. For example, if students have to use the annoying Duo Mobile Two-Factor authentication to access the survey, busy students are likely to just give up and not respond.

B is correct because only Data 100 students who check EdStem have the chance to be included in the sample. Others are excluded.

C is correct because no sampling method prevents response bias, as participants are always subject to some degree of interference from answering with their true opinions. This can be due to multiple reasons, including but not limited to an unwillingness to answer the survey and non-anonymity.

D is incorrect due to the bias in the sample - the most glaring bias being the selection bias detailed in answer choice B.

- (c) [0.5 Pts] As part of a separate transportation analysis, Rachel wants to survey a sample of UC Berkeley students on how many times they traveled out of Berkeley in the past month. She decides to divide the population into groups based on their year (e.g., freshman, sophomore, junior, senior, and graduate students) and then conduct a simple random sample of size  $n_g$  among each group g, where  $n_g$  is proportional to the number of enrolled students at UC Berkeley in group g. What type of sample is this?
  - Convenience sample
  - Stratified random sample
  - Uniform random sample with replacement
  - Post stratification

**Solution:** Rachel is dividing the population into subgroups (strata) based on their year and then randomly sampling from each group. This is the definition of a stratified random sample.

(d) [1 Pt] Rachel is trying to understand how long BayPass users stay on campus. She collects data on hours  $(x_i)$  for three students:  $\{2, 4, 6\}$ . She is using a boxcar kernel with bandwidth  $\alpha = 0.5$  to estimate the density distribution of the data.

A boxcar kernel is defined as follows:

$$K_{\alpha}(x, x_i) = \begin{cases} \frac{1}{\alpha}, & \text{if } |x - x_i| \le \frac{\alpha}{2} \\ 0, & \text{otherwise.} \end{cases}$$

Which of the following KDE plots correctly represents the estimated density?



**Solution:** Option A! The width of each kernel is 0.5 since  $\alpha$  represents the total width. Given  $\alpha = 0.5$ , the correct choice is Option A. In comparison, Option B has  $\alpha = 1$ , making the kernels too wide, and Option C has  $\alpha = 2$ , making them even wider!

(e) Rachel is now examining a small sample dataset representing the number of rides taken in a month by three BayPass users:  $\{1, 6, 17\}$ . She wants to pick a single summary statistic  $\theta$  to describe the data. Determine the value  $\hat{\theta}$  that minimizes each of the following objective functions.

(i) [1 Pt] 
$$R(\theta) = \frac{1}{3} \sum_{i=1}^{3} (x_i - \theta)^2$$
  
 $\bigcirc \frac{8}{3}$   
 $\bigcirc \frac{6}{3}$   
 $\bigcirc 6$   
 $\bigcirc 8$   
(ii) [1 Pt]  $R(\theta) = \frac{1}{1000} \sum_{i=1}^{3} |x_i - \theta|$   
 $\bigcirc \frac{8}{1000}$   
 $\bigcirc 6$   
 $\bigcirc 8$   
 $\bigcirc \frac{6}{1000}$ 

- (f) [1.5 Pts] Rachel is building a model to predict the number of rides taken by UC Berkeley students using their BayPass. Which of the following scenarios would Mean Squared Error (MSE) be preferred over Mean Absolute Error (MAE) for a linear regression task?
  - $\bigcirc$  **True**  $\bigcirc$  False When the model needs to be sensitive to outliers.
  - $\bigcirc$  **True**  $\bigcirc$  False When large errors should be penalized more heavily.

○ **True** ○ False When a smooth, differentiable loss function is required for finding the minimum average loss.

**Solution:** MSE is sensitive to outliers- they are penalized more. MSE will penalize large errors heavily because it squares the errors. MSE is a smooth differentiable function. MSE will always have a unique solution in this case, but MAE is not guaranteed to have a single unique solution.

(g) [0.5 Pts] Rachel fits a Simple Linear Regression (SLR) model to predict the number of rides taken by UC Berkeley students using their BayPass. She uses the hours spent on campus (hours) as a predictor and the number of rides taken (rides) as a response. The equation for the SLR model is given by:

$$\widehat{\text{rides}} = \hat{\theta}_0 + \hat{\theta}_1 \text{ hours}$$

What is the interpretation of the parameter  $\hat{\theta}_1$  in this model?

- Estimated number of rides taken when no hours are spent on campus.
- Estimated average ride increase per additional hour on campus.
- Estimated average rides taken by students.
- $\bigcirc$  Estimated total number of rides by students in the dataset.

**Solution:**  $\hat{\theta}_1$  represents the slope of the regression line, which is the average increase in the predicted number of rides for every additional hour spent on campus.

(h) Rachel continues analyzing the number of rides  $(y_i)$  based on the hours  $(x_i)$ . She models the relationship using SLR without an intercept:

$$\widehat{\mathsf{rides}} = \theta_1 \,\mathsf{hours}$$

Instead of using MSE, she decides to minimize the following custom objective function:

$$R(\theta_1) = \frac{1}{n} \sum_{i=1}^{n} x_i (y_i - \theta_1 x_i)^2$$

(i) [2 Pts] Find the derivative of  $R(\theta_1)$  with respect to  $\theta_1$ . Your answer should be in terms of  $x_i, y_i, \theta_1, n$ . To be eligible for partial credit, show all your work in the box below.

Solution: Using the chain rule:

$$\frac{dR}{d\theta_1} = \frac{1}{n} \sum_{i=1}^n (-2)(-x_i) x_i (y_i - \theta_1 x_i) = \boxed{-\frac{2}{n} \sum_{i=1}^n x_i^2 (y_i - \theta_1 x_i)}$$

(ii) [2 Pts] Find  $\hat{\theta}_1$  that minimizes the objective function. Your answer should be in terms of  $x_i, y_i, n$ . To be eligible for partial credit, show all your work in the box below. **Note:** Assume that the provided objective function is convex.

Solution: Setting the derivative to zero:

$$-\frac{2}{n}\sum_{i=1}^{n}x_{i}^{2}(y_{i}-\hat{\theta}_{1}x_{i})=0$$
$$\sum_{i=1}^{n}x_{i}^{2}(y_{i}-\hat{\theta}_{1}x_{i})=0$$
$$\sum_{i=1}^{n}x_{i}^{2}y_{i}-\hat{\theta}_{1}\sum_{i=1}^{n}x_{i}^{3}=0$$
$$\hat{\theta}_{1}=\frac{\sum_{i=1}^{n}x_{i}^{2}y_{i}}{\sum_{i=1}^{n}x_{i}^{3}}$$

### 4 We Miss Moffitt </3 [6.5 Pts]

Since the closure of Moffitt Library, course staff have struggled to find a study spot. One popular alternative is MLK Student Union. Sarah decides to use Ordinary Least Squares (OLS) to predict the number of students ( $\mathbb{Y}$ ) in MLK Student Union for a given humidity level and time of day.

(a) [2 Pts] Sarah fits an OLS model to predict the number of students (num\_students) at MLK Student Union using the humidity level (humidity) and time of day (time) as predictors. Her fitted model is:

num\_students = 
$$\hat{\theta}_0 + \hat{\theta}_1 \times \text{humidity} + \hat{\theta}_2 \times \text{time}$$

After fitting the optimal model, Sarah examines the residual vector  $\vec{e} = \mathbb{Y} - \hat{\mathbb{Y}}$ . For each statement, indicate whether it is True or False.

True False *e* is orthogonal to all columns of the design matrix X.
 True False If θ̂<sub>0</sub> is positive, the average of the residuals is also positive.
 True False Assuming X is full column rank, the sum of squared residuals is minimized when θ̂ is calculated using the normal equation.

 $\bigcirc$  True  $\bigcirc$  False Residual vector  $\vec{e}$  has the same dimension as the parameter vector  $\hat{\theta}$ .

(b) [0.5 Pts] Sarah observes a non-linear relationship between humidity and num\_students, as shown below:



Let  $x_i$  represent humidity for the *i*-th data point. Let  $h_i$  represent the transformed value of the humidity level used in the model.

Which of the following transformations should Sarah apply to linearize the data?

- $\bigcirc h_i = \log(x_i)$
- $\bigcirc \log(h_i) = x_i$  $\bigcirc h_i = (x_i)^2$  $\bigcirc h_i = (x_i)^3$

- (c) Answer the following questions about OLS models:
  - (i) [2 Pts] We want to fit an OLS model with n observations and p + 1 features (including the intercept). Select **True** if the dimensions of the following matrices are correct in this task, and **False** otherwise.
    - $\begin{array}{c|c|c|c|c|c|c|c|c|} \hline \mathbf{True} & \bigcirc & \mathbf{False} & & \mathbb{X} : n \times (p+1) \\ \hline & \mathbf{True} & \bigcirc & \mathbf{False} & & \mathbb{X}^T \mathbb{X} : (p+1) \times (p+1) \\ \hline & & \mathbf{True} & \bigcirc & \mathbf{False} & & \hat{\theta} : p \times 1 \\ \hline & & & \mathbf{True} & \bigcirc & \mathbf{False} & & \mathbb{Y} : n \times 1 \end{array}$
  - (ii) [1 Pt] Which of the following best explains the condition required for the OLS model to produce a unique solution?
    - $\bigcirc$  X must not be full column rank.
    - $\bigcirc$  X must be a square matrix.
    - $\bigcirc \mathbb{X}^T \mathbb{X}$  must be invertible.
    - $\bigcirc$  X must have fewer rows than columns.
  - (iii) [1 Pt] Which of the following best describes the geometric interpretation of the OLS prediction vector  $\hat{\mathbb{Y}}$ ?
    - $\bigcirc \ \hat{\mathbb{Y}}$  is the difference between  $\mathbb{Y}$  and the orthogonal projection of  $\mathbb{Y}$  onto the span of  $\mathbb{X}.$
    - $\bigcirc$   $\hat{\mathbb{Y}}$  is the orthogonal projection of  $\mathbb{Y}$  onto the span of  $\mathbb{X}$ .
    - $\bigcirc$  When  $\mathbb{X}$  is not full column rank,  $\hat{\mathbb{Y}}$  is still uniquely determined because the normal equation always has a single solution.
    - $\bigcirc$   $\hat{\mathbb{Y}}$  is the orthogonal projection of  $\mathbb{X}$  onto the span of  $\mathbb{Y}$ .

#### You are done with the midterm- Congratulations!

Draw your favorite DATA 100/200 memory so far!