

EECS 398 W25 Midterm Review

February 23, 2025

Announcements

- The Midterm Exam is on **Tuesday, February 25th from 7-9PM.**
 - It covers Lectures 1-12, Homeworks 1-6, and Discussions 1-7.
- Midterm Review in lecture tomorrow (going over F24 Final 1-8.2).
- Study Tips
 - Go through lecture notebooks & homeworks to help make cheat sheet (one page, double-sided, **handwritten**).
 - Do [discussion problems](#).
 - Take F24 Midterm and Problems 1-8.2 of F24 Final (besides SQL question).

Agenda

- We'll be working through <https://study.practicaldsc.org/mt-review-sunday/index.html>.
- We'll post these annotated slides and the recording after, along with enabling solutions on the study site for this worksheet.

Grouping, Querying, and Merging - Akanksha

The EECS 398 staff are looking into hotels — some in San Diego, for their family to stay at for graduation (and to eat Mexican food), and some elsewhere, for summer trips.

Each row of `hotels` contains information about a different hotel in San Diego. Specifically, for each hotel, we have:

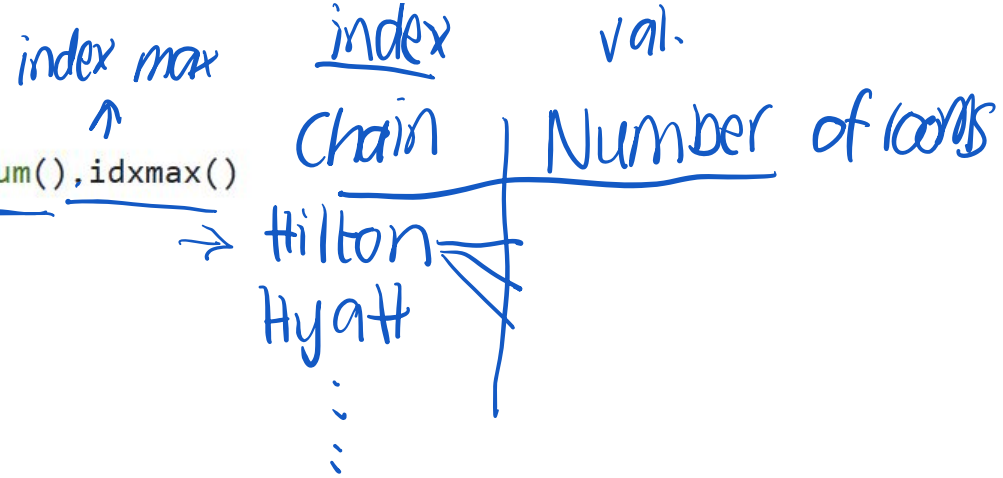
- `"Hotel Name" (str)`: The name of the hotel. **Assume hotel names are unique.**
- `"Location" (str)`: The hotel's neighborhood in San Diego.
- `"Chain" (str)`: The chain the hotel is a part of; either `"Hilton"`, `"Marriott"`, `"Hyatt"`, or `"Other"`. A hotel chain is a group of hotels owned or operated by a shared company.
- `"Number of Rooms" (int)`: The number of rooms the hotel has.

The first few rows of `hotels` are shown below, but `hotels` has many more rows than are shown.

	Hotel Name	Location	Chain	Number of Rooms
0	Hotel del Coronado	Coronado	Hilton	680
1	Manchester Grand Hyatt	Downtown	Hyatt	1628
2	Hilton San Diego Bayfront	Downtown	Hilton	1190
3	Pendry San Diego	Gaslamp Quarter	Other	317
4	The Westin San Diego Gaslamp Quarter	Gaslamp Quarter	Marriott	450
5	San Diego Marriott La Jolla	La Jolla	Marriott	462
6	La Valencia Hotel	La Jolla	Other	114
7	Coronado Island Marriott Resort & Spa	Coronado	Marriott	310

Now, consider the variable `summed`, defined below.

```
summed = hotels.groupby("Chain")["Number of Rooms"].sum(), idxmax()
```



Problem 1.1

What is `type(summed)`?

- `int`
- `str`
- `Series`
- `DataFrame`
- `DataFrameGroupBy`

	Hotel Name	Location	Chain	Number of Rooms
0	Hotel del Coronado	Coronado	Hilton	680
1	Manchester Grand Hyatt	Downtown	Hyatt	1628
2	Hilton San Diego Bayfront	Downtown	Hilton	1190
3	Pendry San Diego	Gaslamp Quarter	Other	317
4	The Westin San Diego Gaslamp Quarter	Gaslamp Quarter	Marriott	450
5	San Diego Marriott La Jolla	La Jolla	Marriott	462
6	La Valencia Hotel	La Jolla	Other	114
7	Coronado Island Marriott Resort & Spa	Coronado	Marriott	310

Problem 1.2

In one sentence, explain what the value of `summed` means. Phrase your explanation as if you had to give it to someone who is not a data science major; that is, don't say something like "it is the result of grouping `hotels` by `"Chain"`, selecting the `"Number of Rooms"` column, ...", but instead, give the value context.

```
summed = hotels.groupby("Chain")["Number of Rooms"].sum().idxmax()
```

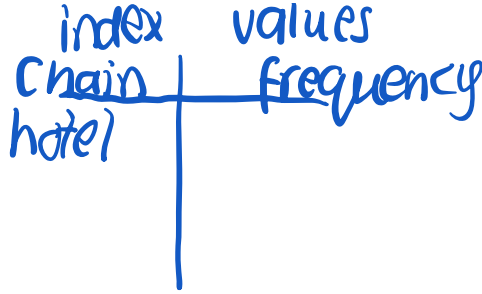
hotel chain with the highest total number of rooms

	Hotel Name	Location	Chain	Number of Rooms
0	Hotel del Coronado	Coronado	Hilton	680
1	Manchester Grand Hyatt	Downtown	Hyatt	1628
2	Hilton San Diego Bayfront	Downtown	Hilton	1190
3	Pendry San Diego	Gaslamp Quarter	Other	317
4	The Westin San Diego Gaslamp Quarter	Gaslamp Quarter	Marriott	450
5	San Diego Marriott La Jolla	La Jolla	Marriott	462
6	La Valencia Hotel	La Jolla	Other	114
7	Coronado Island Marriott Resort & Spa	Coronado	Marriott	310

Problem 1.3

Consider the variable `curious`, defined below.

```
curious = frame["Chain"].value_counts().idxmax()
```



Fill in the blank: `curious` is guaranteed to be equal to `summed` only if `frame` has one row for every ___ in San Diego.

- hotel
- hotel chain
- hotel room
- neighborhood

	Hotel Name	Location	Chain	Number of Rooms
0	Hotel del Coronado	Coronado	Hilton	680
1	Manchester Grand Hyatt	Downtown	Hyatt	1628
2	Hilton San Diego Bayfront	Downtown	Hilton	1190
3	Pendry San Diego	Gaslamp Quarter	Other	317
4	The Westin San Diego Gaslamp Quarter	Gaslamp Quarter	Marriott	450
5	San Diego Marriott La Jolla	La Jolla	Marriott	462
6	La Valencia Hotel	La Jolla	Other	114
7	Coronado Island Marriott Resort & Spa	Coronado	Marriott	310

Problem 1.4

Fill in the blanks so that `popular_areas` is an array of the names of the unique neighborhoods that have at least 5 hotels and at least 1000 hotel rooms.

```
f = lambda df: __ (i) __
```

```
popular_areas = (hotels
```

```
.groupby(__ (ii) __)
```

```
.__ (iii) __
```

```
.__ (iv) __)
```

location
filter

(rows, cols)

$df.shape[0] \geq 5$

and $df["Number of Rooms"].sum() \geq 1000$



1. What goes in blank (i)?

"Hotel Name"

"Location"

"Chain"

"Number of Rooms"

3. What goes in blank (iii)?

`agg(f)`

`filter(f)`

`transform(f)`

4. What goes in blank (iv)?

`["Location"].unique`

	Hotel Name	Location	Chain	Number of Rooms
0	Hotel del Coronado	Coronado	Hilton	680
1	Manchester Grand Hyatt	Downtown	Hyatt	1628
2	Hilton San Diego Bayfront	Downtown	Hilton	1190
3	Pendry San Diego	Gaslamp Quarter	Other	317
4	The Westin San Diego Gaslamp Quarter	Gaslamp Quarter	Marriott	450
5	San Diego Marriott La Jolla	La Jolla	Marriott	462
6	La Valencia Hotel	La Jolla	Other	114
7	Coronado Island Marriott Resort & Spa	Coronado	Marriott	310

Problem 1.5

Consider the code below.

```
cond1 = hotels["Chain"] == "Marriott"
cond2 = hotels["Location"] == "Coronado"
combined = hotels[cond1].merge(hotels[cond2], on="Hotel Name", how="???")
```

hotels

Chain	Location	cond1	cond2
Marriott	Coronado	T	T
Marriott	Coronado	F	
Marriott	AnnArbor	T	F

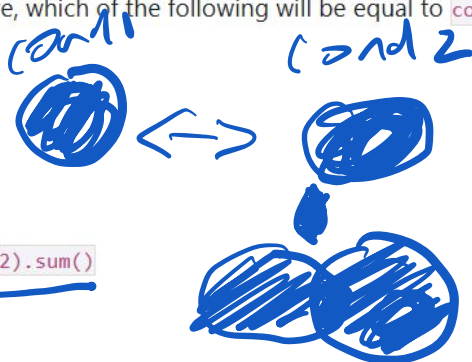
1. If we replace `???` with `"inner"` in the code above, which of the following will be equal to `combined.shape[0]`?

- `min(cond1.sum(), cond2.sum())`
- `(cond1 & cond2).sum()`
- `cond1.sum() + cond2.sum()`
- `cond1.sum() + cond2.sum() - (cond1 & cond2).sum()`
- `cond1.sum() + (cond1 & cond2).sum()`

outer # rows:
 num rows in inner merge
 + non-matches in left df
 + non-matches in right df

2. If we replace `???` with `"outer"` in the code above, which of the following will be equal to `combined.shape[0]`?

- `min(cond1.sum(), cond2.sum())`
- `(cond1 & cond2).sum()`
- `cond1.sum() + cond2.sum()`
- `cond1.sum() + cond2.sum() - (cond1 & cond2).sum()`
- `cond1.sum() + (cond1 & cond2).sum()`



	Hotel Name	Location	Chain	Number of Rooms
0	Hotel del Coronado	Coronado	Hilton	680
1	Manchester Grand Hyatt	Downtown	Hyatt	1628
2	Hilton San Diego Bayfront	Downtown	Hilton	1190
3	Pendry San Diego	Gaslamp Quarter	Other	317
4	The Westin San Diego Gaslamp Quarter	Gaslamp Quarter	Marriott	450
5	San Diego Marriott La Jolla	La Jolla	Marriott	462
6	La Valencia Hotel	La Jolla	Other	114
7	Coronado Island Marriott Resort & Spa	Coronado	Marriott	310

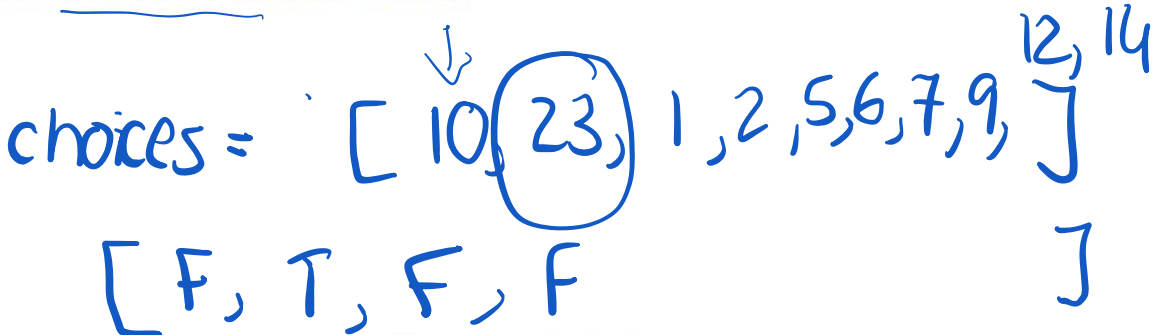
Random Simulations - Akanksha

Problem 2

Billina Records, a new record company focused on creating new TikTok audios, has its offices on the 23rd floor of a skyscraper with 75 floors (numbered 1 through 75). The owners of the building promised that 10 different random floors will be selected to be renovated.

Below, fill in the blanks to complete a simulation that will estimate the probability that Billina Records' floor will be renovated.

```
total = 0
repetitions = 10000
for i in np.arange(repetitions):
    choices = np.random.choice(__(a)__, 10, ____(b)__)
    if ____(c)__:
        total = total + 1
prob_renovate = total / repetitions
```



What goes in blank (a)?

- np.arange(1, 75)
- np.arange(10, 75)
- np.arange(0, 76)
- np.arange(1, 76)

↑
range is exclusive

What goes in blank (b)?

- replace=True
- replace=False

What goes in blank (c)?

- choices == 23
- choices is 23
- np.count_nonzero(choices == 23) > 0
- np.count_nonzero(choices) == 23
- choices.str.contains(23)

Joining 2 dataframes



Merging - Caleb

rows from a merge

Inner Merge: $\sum \# \text{match in df1} \times \# \text{match in df2}$



Left Merge: # rows from inner merge + # non-matches in left df



Outer Merge: # rows from inner merge + non-matches in

df1

df2

8	A
8	B
4	C

8	A
5	
3	



inner df1 df2

8	A	A
8	B	A

$2 \times 1 = 2$

left df1 df2 right and left df

8	A	A
8	B	A
4	C	Na

$2 \times 1 + 1 = 3$

outer

8
8
5
3

$2 \times 1 + 1 + 2$

5

Suppose the DataFrame `today` consists of 15 rows — 3 rows for each of 5 different `"artist_names"`. For each artist, it contains the `"track_name"` for their three most-streamed songs today. For instance, there may be one row for `"olivia rodrigo"` and `"favorite crime"`, one row for `"olivia rodrigo"` and `"drivers license"`, and one row for `"olivia rodrigo"` and `"deja vu"`.

Another DataFrame, `genres`, is shown below in its entirety.

	<u>artist_names</u>	genre
0	harry styles	Pop
1	olivia rodrigo	Pop
2	glass animals	Alternative
3	drake	Hip-Hop/Rap
4	doja cat	Hip-Hop/Rap

today

artist name

track

olivia

olivia

olivia

ed

ed

ed

ed

fav crime

~~~~~

~~~~~

~~~~~

~~~~~

Problem 3.1

Suppose we perform an **inner** merge between `today` and `genres` on `"artist_names"`. If the five `"artist_names"` in `today` are the same as the five `"artist_names"` in `genres`, what fraction of the rows in the merged DataFrame will contain `"Pop"` in the `"genre"` column? Give your answer as a simplified fraction.

6 rows with pop in genre
15 rows total
3 + 5
2/5

Handwritten notes showing a list of artist names and their corresponding genres:

Artist	Genre
harry	Pop
harry	Pop
harry	Pop
olivia	Pop
olivia	Pop
olivia	Pop
glass	Alt

	artist_names	genre
0	harry styles	Pop
1	olivia rodrigo	Pop
2	glass animals	Alternative
3	drake	Hip-Hop/Rap
4	doja cat	Hip-Hop/Rap

Suppose the DataFrame `today` consists of 15 rows — 3 rows for each of 5 different `"artist_names"`. For each artist, it contains the `"track_name"` for their three most-streamed songs today. For instance, there may be one row for `"olivia rodrigo"` and `"favorite crime"`, one row for `"olivia rodrigo"` and `"drivers license"`, and one row for `"olivia rodrigo"` and `"deja vu"`.

genres

Problem 3.2

Suppose we perform an **inner** merge between `today` and `genres` on `"artist_names"`. Furthermore suppose that the only overlapping `"artist_names"` between `today` and `genres` are `"drake"` and `"olivia rodrigo"`. What fraction of the rows in the merged DataFrame will contain `"Pop"` in the `"genre"` column? Give your answer as a simplified fraction.

3 rows with Pop in genre

6 rows total

↑
3x2

$$\frac{1}{2}$$

olivia	~	Pop
olivia	~	Pop
olivia	~	Pop
drake	~	Hip
drake	~	Hip
drake	~	Hip

	artist_names	genre
0	harry styles	Pop
1	olivia rodrigo	Pop
2	glass animals	Alternative
3	drake	Hip-Hop/Rap
4	doja cat	Hip-Hop/Rap

genres

Suppose the DataFrame `today` consists of 15 rows — 3 rows for each of 5 different `"artist_names"`. For each artist, it contains the `"track_name"` for their three most-streamed songs today. For instance, there may be one row for `"olivia rodrigo"` and `"favorite crime"`, one row for `"olivia rodrigo"` and `"drivers license"`, and one row for `"olivia rodrigo"` and `"deja vu"`.

Problem 3.3

Suppose we perform an **outer** merge between `today` and `genres` on `"artist_names"`. Furthermore, suppose that the only overlapping `"artist_names"` between `today` and `genres` are `"drake"` and `"olivia rodrigo"`. What fraction of the rows in the merged DataFrame will contain `"Pop"` in the `"genre"` column? Give your answer as a simplified fraction.

3 rows with 'Pop' from inner merge
1 row from Harry styles non-match
in genres, with Pop

18 total rows

6
Inner Merge
3 non-matches from genres

1 non-matched from today

$\frac{4}{18} \rightarrow \frac{2}{9}$

	artist_names	genre
0	harry styles	Pop
1	olivia rodrigo	Pop
2	glass animals	Alternative
3	drake	Hip-Hop/Rap
4	doja cat	Hip-Hop/Rap

genres

Suppose the DataFrame `today` consists of 15 rows — 3 rows for each of 5 different `"artist_names"`. For each artist, it contains the `"track_name"` for their three most-streamed songs today. For instance, there may be one row for `"olivia rodrigo"` and `"favorite crime"`, one row for `"olivia rodrigo"` and `"drivers license"`, and one row for `"olivia rodrigo"` and `"deja vu"`.

Missing Value Imputation - Caleb

Listwise deletion: drop missing values from calculation

Mean imputation \rightarrow fill missing values with mean

Conditional mean imputation: fill missing values with mean of some group

Probabilistic Imputation: Replace missing values with a random sample of data.

The DataFrame `random_10` contains the `"track_name"` and `"genre"` of 10 randomly-chosen songs in Spotify's Top 200 today, along with their `"genre_rank"`, which is their rank in the Top 200 **among songs in their "genre"**. For instance, "the real slim shady" is the 20th-ranked Hip-Hop/Rap song in the Top 200 today. `random_10` is shown below in its entirety.

	track_name	genre rank	genre
0	good looking	7.0	Alternative
1	drowning (feat. kodak black)	NaN	Hip-Hop/Rap
2	the real slim shady	20.0	Hip-Hop/Rap
3	worldwide steppers	2.0	Hip-Hop/Rap
4	2055	5.0	Hip-Hop/Rap
5	drivers license	9.0	Pop
6	cinema	2.0	Pop
7	dos mil 16	4.0	Pop
8	happier than ever	NaN	Pop
9	bam bam (feat. ed sheeran)	NaN	Pop

The "genre_rank" column of random_10 contains missing values. Below, we provide four different imputed "genre_rank" columns, each of which was created using a different imputation technique. On the next page, match each of the four options to the imputation technique that was used in the option.

Option A

genre rank	genre
7.0	Alternative
5.0	Hip-Hop/Rap
20.0	Hip-Hop/Rap
2.0	Hip-Hop/Rap
5.0	Hip-Hop/Rap
9.0	Pop
2.0	Pop
4.0	Pop
2.0	Pop
2.0	Pop

Option B

genre rank	genre
7.0	Alternative
7.0	Hip-Hop/Rap
20.0	Hip-Hop/Rap
2.0	Hip-Hop/Rap
5.0	Hip-Hop/Rap
9.0	Pop
2.0	Pop
4.0	Pop
7.0	Pop
7.0	Pop

Option C

genre rank	genre
7.0	Alternative
2.0	Hip-Hop/Rap
20.0	Hip-Hop/Rap
2.0	Hip-Hop/Rap
5.0	Hip-Hop/Rap
9.0	Pop
2.0	Pop
4.0	Pop
2.0	Pop
7.0	Pop

Option D

genre rank	genre
7.0	Alternative
9.0	Hip-Hop/Rap
20.0	Hip-Hop/Rap
2.0	Hip-Hop/Rap
5.0	Hip-Hop/Rap
9.0	Pop
2.0	Pop
4.0	Pop
5.0	Pop
5.0	Pop

Problem 4.1

In which option was unconditional mean imputation used?



Option A		Option B		Option C		Option D	
genre rank	genre	genre rank	genre	genre rank	genre	genre rank	genre
7.0	Alternative	7.0	Alternative	7.0	Alternative	7.0	Alternative
5.0	Hip-Hop/Rap	<u>7.0</u>	Hip-Hop/Rap	2.0	Hip-Hop/Rap	9.0	Hip-Hop/Rap
20.0	Hip-Hop/Rap	20.0	Hip-Hop/Rap	20.0	Hip-Hop/Rap	20.0	Hip-Hop/Rap
2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap
5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap
9.0	Pop	9.0	Pop	9.0	Pop	9.0	Pop
2.0	Pop	2.0	Pop	2.0	Pop	2.0	Pop
4.0	Pop	4.0	Pop	4.0	Pop	4.0	Pop
2.0	Pop	<u>7.0</u>	Pop	2.0	Pop	5.0	Pop
2.0	Pop	<u>7.0</u>	Pop	7.0	Pop	5.0	Pop

mean of filled
vals: >

	track_name	genre rank	genre
0	good looking	<u>7.0</u>	Alternative
1	drowning (feat. kodak black)	NaN	Hip-Hop/Rap
2	the real slim shady	<u>20.0</u>	Hip-Hop/Rap
3	worldwide steppers	<u>2.0</u>	Hip-Hop/Rap
4	2055	<u>5.0</u>	Hip-Hop/Rap
5	drivers license	<u>9.0</u>	Pop
6	cinema	<u>2.0</u>	Pop
7	dos mil 16	<u>4.0</u>	Pop
8	happier than ever	NaN	Pop
9	bam bam (feat. ed sheeran)	NaN	Pop

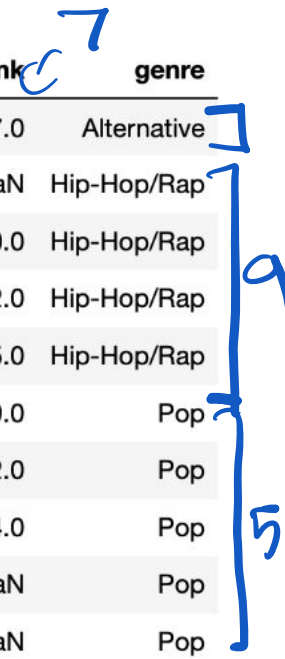
Problem 4.2

In which option was mean imputation conditional on "genre" used?



Option A		Option B		Option C		Option D	
genre rank	genre	genre rank	genre	genre rank	genre	genre rank	genre
7.0	Alternative	7.0	Alternative	7.0	Alternative	7.0	Alternative
5.0	Hip-Hop/Rap	7.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	9.0	Hip-Hop/Rap
20.0	Hip-Hop/Rap	20.0	Hip-Hop/Rap	20.0	Hip-Hop/Rap	20.0	Hip-Hop/Rap
2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap
5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap
9.0	Pop	9.0	Pop	9.0	Pop	9.0	Pop
2.0	Pop	2.0	Pop	2.0	Pop	2.0	Pop
4.0	Pop	4.0	Pop	4.0	Pop	4.0	Pop
2.0	Pop	7.0	Pop	2.0	Pop	5.0	Pop
2.0	Pop	7.0	Pop	7.0	Pop	5.0	Pop

	track_name	genre rank	genre
0	good looking	7.0	Alternative
1	drowning (feat. kodak black)	NaN	Hip-Hop/Rap
2	the real slim shady	20.0	Hip-Hop/Rap
3	worldwide steppers	2.0	Hip-Hop/Rap
4	2055	5.0	Hip-Hop/Rap
5	drivers license	9.0	Pop
6	cinema	2.0	Pop
7	dos mil 16	4.0	Pop
8	happier than ever	NaN	Pop
9	bam bam (feat. ed sheeran)	NaN	Pop



Problem 4.3

In which option was unconditional probabilistic imputation used?

unconditional

[7, 9, 2]
[20, 5, 4]

Option A		Option B		Option C		Option D	
genre rank	genre	genre rank	genre	genre rank	genre	genre rank	genre
7.0	Alternative	7.0	Alternative	7.0	Alternative	7.0	Alternative
5.0	Hip-Hop/Rap	7.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	9.0	Hip-Hop/Rap
<u>20.0</u>	Hip-Hop/Rap	20.0	Hip-Hop/Rap	<u>20.0</u>	Hip-Hop/Rap	20.0	Hip-Hop/Rap
2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap
5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap
9.0	Pop	9.0	Pop	9.0	Pop	9.0	Pop
2.0	Pop	2.0	Pop	2.0	Pop	2.0	Pop
4.0	Pop	4.0	Pop	4.0	Pop	4.0	Pop
2.0	Pop	7.0	Pop	2.0	Pop	5.0	Pop
<u>2.0</u>	Pop	7.0	Pop	<u>7.0</u>	Pop	5.0	Pop

	track_name	genre rank	genre
0	good looking	7.0	Alternative
1	drowning (feat. kodak black)	NaN	Hip-Hop/Rap
2	the real slim shady	20.0	Hip-Hop/Rap
3	worldwide steppers	2.0	Hip-Hop/Rap
4	2055	5.0	Hip-Hop/Rap
5	drivers license	9.0	Pop
6	cinema	2.0	Pop
7	dos mil 16	4.0	Pop
8	happier than ever	NaN	Pop
9	bam bam (feat. ed sheeran)	NaN	Pop

Problem 4.4

In which option was probabilistic imputation conditional on "genre" used?



Option A		Option B		Option C		Option D	
genre rank	genre	genre rank	genre	genre rank	genre	genre rank	genre
7.0	Alternative	7.0	Alternative	7.0	Alternative	7.0	Alternative
5.0	Hip-Hop/Rap	7.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	9.0	Hip-Hop/Rap
20.0	Hip-Hop/Rap	20.0	Hip-Hop/Rap	20.0	Hip-Hop/Rap	20.0	Hip-Hop/Rap
2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap	2.0	Hip-Hop/Rap
5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap	5.0	Hip-Hop/Rap
9.0	Pop	9.0	Pop	9.0	Pop	9.0	Pop
2.0	Pop	2.0	Pop	2.0	Pop	2.0	Pop
4.0	Pop	4.0	Pop	4.0	Pop	4.0	Pop
2.0	Pop	7.0	Pop	2.0	Pop	5.0	Pop
2.0	Pop	7.0	Pop	7.0	Pop	5.0	Pop

	track_name	genre rank	genre
0	good looking	7.0	Alternative
1	drowning (feat. kodak black)	NaN	Hip-Hop/Rap
2	the real slim shady	20.0	Hip-Hop/Rap
3	worldwide steppers	2.0	Hip-Hop/Rap
4	2055	5.0	Hip-Hop/Rap
5	drivers license	9.0	Pop
6	cinema	2.0	Pop
7	dos mil 16	4.0	Pop
8	happier than ever	NaN	Pop
9	bam bam (feat. ed sheeran)	NaN	Pop

Problem 6.2

Rahul wants to extract the 'instock availability' status of the book titled 'A Light in the Attic'. Which of the following expressions will evaluate to "In Stock"? Assume that Rahul has already parsed the HTML into a BeautifulSoup object stored in the variable named `soup`.

Code Snippet A

```
soup.find('p', attrs = {'class': 'instock availability'})\
.get('icon-ok').strip()
```

Code Snippet B

```
soup.find('p', attrs = {'class': 'instock availability'}).text.strip()
```

Code Snippet C

```
soup.find('p', attrs = {'class': 'instock availability'}).find('i')\
.text.strip()
```

Code Snippet D

```
soup.find('div', attrs = {'class': 'product_price'})\
.find('p', attrs = {'class': 'instock availability'})\
.find('i').text.strip()
```

```
<HTML>
<H1>The Book Club</H1>
<BODY BGCOLOR="FFFFFF">
Email us at <a href="mailto:support@thebookclub.com">
support@thebookclub.com</a>.
```

```
<div>
<ol class="row">
<li class="book_list">

<article class="product_pod">
<div class="image_container">

</div>

<p class="star-rating Three"></p>

<h3>
<a href="cat/index.html" title="A Light in the Attic">
A Light in the Attic
</a>
</h3>

<div class="product_price">
<p class="price_color">£51.77</p>

<p class="instock availability">
<i class="icon-ok"></i>
In stock
</p>

</div>
</article>
</li>
</ol>

</div>
</BODY>
</HTML>
```

'In Stock'

X

✓

X

||

||

<p class="instock availability">

In stock

↑

Regular Expressions - Angela

Problem 5

regex101.com

RegexEgg.com → [regex cheat sheet](#) [tables](#)

You want to use regular expressions to extract out the number of ounces from the 5 product names below.

Index	Product Name	Expected Output
0	Adult Dog Food 18-Count, 3.5 oz Pouches	3.5
1	Gardetto's Snack Mix, 1.75 Ounce	1.75
2	Colgate Whitening Toothpaste, 3 oz Tube	3
3	Adult Dog Food, 13.2 oz Cans 24 Pack	13.2
4	Keratin Hair Spray 2.6 oz	6

Regex finds a pattern, then captures it

The names are stored in a pandas Series called `names`. For each snippet below, select the indexes for all the product names that **will not** be matched correctly.

Pattern looks for: `--- oz`

For the snippet below, which indexes correspond to products that will **not** be matched correctly?

```
regex = r'(\d+(\.|\d)?) oz'
names.str.findall(regex)
```

extract the number of oz from each product name

- 0
- 1
- 2
- 3
- 4
- All names will be matched correctly.

`()`: groups sequence of digits and/or periods
capture group

`\d`: one digit from `[0-9]`
`[]`: character class, matches any one character inside `[]`
 → `[abc]`: matches either a, b, or c
 → `gr[ae]y`: matches "gray" or "grey"
 → NOT literal `"[]"`
 . : inside `[]`, so literal period "."

`+`: one or more of the preceding element

You want to use regular expressions to extract out the number of ounces from the 5 product names below.

Index	Product Name	Expected Output
0	Adult Dog Food 18-Count, 3.5 oz Pouches ✓	3.5
1	Gardetto's Snack Mix, 1.75 Ounce	1.75
2	Colgate Whitening Toothpaste, 3 oz Tube	3
3	Adult Dog Food, 13.2 oz Cans 24 Pack ✓	13.2
4	Keratin Hair Spray 2!6 oz captures "2!6" instead of 6	2!6

The names are stored in a pandas Series called `names`. For each snippet below, select the indexes for all the product names that **will not** be matched correctly. requires at least 1 digit, then some character, then one or more digits or "ounce"

For the snippet below, which indexes correspond to products that will **not** be matched correctly?

```
regex = r'(\d+(\.\d+)?)\d+\s(oz|ounce)'
```

0

1

2

3

4

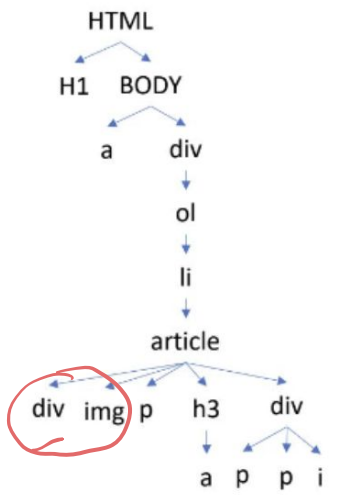
All names will be matched correctly.

- \dagger : matches one or more digits
- $?$: after the $+$ makes this quantifier lazy (non-greedy) → matches as few digits as possible "if it exists or not" OK with no extra digits.
- \cdot : not escaped "\." or in $[\]$ → matches any single character except a newline → probably intended to match a decimal point
- $()$: capture group extracts matched number separately
- $|$: "or" operator splits regex into two alternatives
 - $(\d+(\.\d+)?)\d+\s\text{oz}$
 - $\d+\s\text{ounce}$

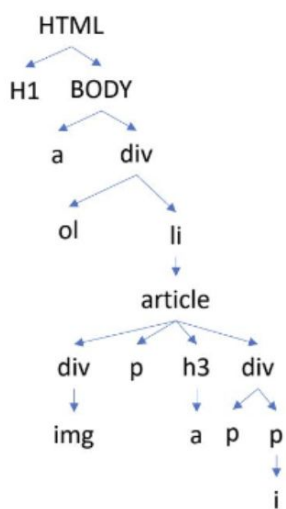
Web Scraping - Abhi

Which is the equivalent Document Object Model (DOM) tree of this HTML file?

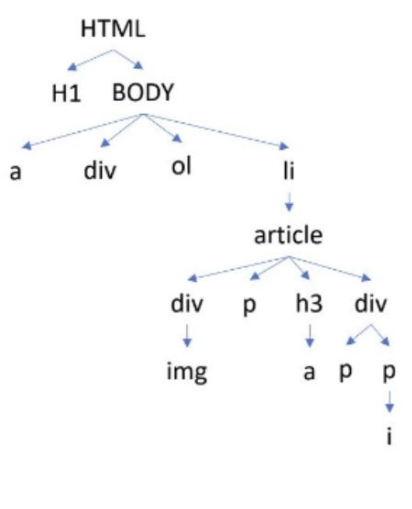
Tree A



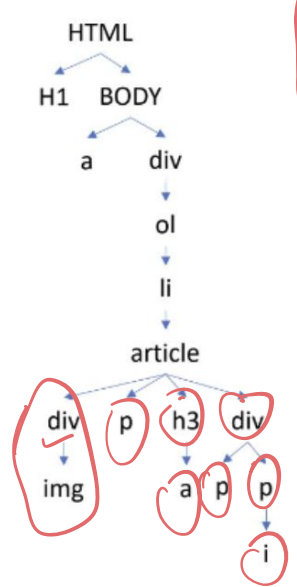
~~Tree B~~



~~Tree C~~



Tree D



div
↓
ol
↓
li

```

<HTML>
<H1>The Book Club</H1>
<BODY BGCOLOR="FFFFFF">
  Email us at <a href="mailto:support@thebookclub.com">
  support@thebookclub.com</a>.
  
```

```

<div>
  <ol class="row">
    <li class="book_list">
      <article class="product_pod">
        <div class="image_container">
          
        </div>
        <p class="star-rating Three"></p>
        <h3>
          <a href="cat/index.html" title="A Light in the Attic">
            A Light in the Attic
          </a>
        </h3>
        <div class="product_price">
          <p class="price_color">£51.77</p>
          <p class="instock availability">
            <i class="icon-ok"></i>
            In stock
          </p>
        </div>
      </article>
    </li>
  </ol>
</div>
</BODY>
</HTML>
  
```

Problem 6.2

Rahul wants to extract the 'instock availability' status of the book titled 'A Light in the Attic'. Which of the following expressions will evaluate to "In Stock"? Assume that Rahul has already parsed the HTML into a BeautifulSoup object stored in the variable named `soup`.

Code Snippet A

```
soup.find('p', attrs = {'class': 'instock availability'})\n.get('icon-ok').strip()
```

Code Snippet B

```
soup.find('p', attrs = {'class': 'instock availability'}).text.strip()
```

Code Snippet C

```
soup.find('p', attrs = {'class': 'instock availability'}).find('i')\n.text.strip()
```

Code Snippet D

```
soup.find('div', attrs = {'class': 'product_price'})\n.find('p', attrs = {'class': 'instock availability'})\n.find('i').text.strip()
```

```
<HTML>\n<H1>The Book Club</H1>\n<BODY BGCOLOR="FFFFFF">\nEmail us at <a href="mailto:support@thebookclub.com">\nsupport@thebookclub.com</a>.
```

```
<div>\n  <ol class="row">\n    <li class="book_list">\n\n      <article class="product_pod">\n        <div class="image_container">\n          \n        </div>\n\n        <p class="star-rating Three"></p>\n\n        <h3>\n        <a href="cat/index.html" title="A Light in the Attic">\n        A Light in the Attic\n        </a>\n        </h3>\n\n        <div class="product_price">\n          <p class="price_color">£51.77</p>\n\n          <p class="instock availability">\n            <i class="icon-ok"></i>\n            In stock\n          </p>\n        </div>\n      </article>\n    </li>\n  </ol>\n</div>\n</BODY>\n</HTML>
```

'In Stock'

X

✓

X

||

||

<p class="instock availability">

In stock

↑

Problem 6.3

Rahul also wants to extract the number of stars that the book titled 'A Light in the Attic' received. If you look at the HTML file, you will notice that the book received a star rating of three. Which code snippet will evaluate to "Three"?

Code Snippet A

```
soup.find('article').get('class').strip()
```

Code Snippet B

```
soup.find('p').text.split(' ')
```

Code Snippet C

```
soup.find('p').get('class')[1]
```

None of the above

"Three"
'map'
'map'
.get('class')

'star-rating Three'
['star-rating', 'Three']
↑
'star-rating-Three'

```
<HTML>
<H1>The Book Club</H1>
<BODY BGCOLOR="FFFFFF">
Email us at <a href="mailto:support@thebookclub.com">
support@thebookclub.com</a>.

<div>
  <ol class="row">
    <li class="book_list">

      <article class="product_pod">
        <div class="image_container">
          
        </div>
        <p class="star-rating Three"></p>

        <h3>
        <a href="cat/index.html" title="A Light in the Attic">
A Light in the Attic
        </a>
        </h3>

        <div class="product_price">
          <p class="price_color">£51.77</p>

          <p class="instock availability">
            <i class="icon-ok"></i>
            In stock
          </p>

        </div>
      </article>
    </li>
  </ol>

</div>
</BODY>
</HTML>
```


Text as Data - Abhi

$$\cos \theta = \frac{\bar{u} \cdot \bar{v}}{\|\bar{u}\| \|\bar{v}\|}$$

$$\bar{u} \quad \bar{v}$$

Problem 7

Tahseen decides to look at reviews for the same hotel, but he modifies them so that the only terms they contain are "taco" and "sand". The bag-of-words representations of three reviews are shown as vectors below.

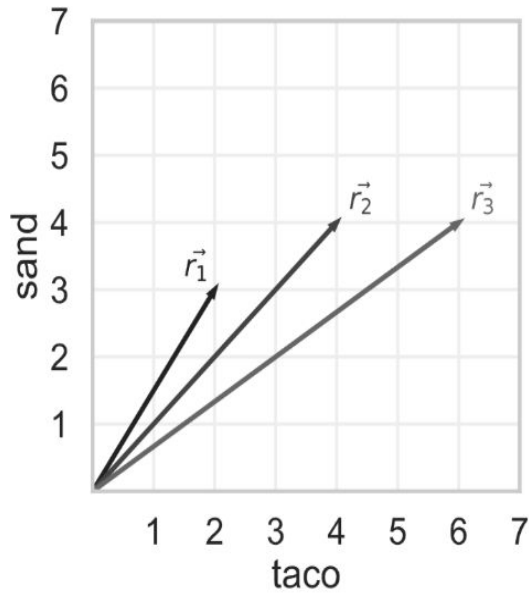
$$\vec{r}_1 = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad \vec{r}_2 = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$$

$$\vec{r}_3 = \begin{bmatrix} 6 \\ 4 \end{bmatrix}$$

$$\frac{2 \cdot 4 + 3 \cdot 4}{\sqrt{13} \cdot \sqrt{32}} = \frac{20}{4\sqrt{26}}$$

$$\sqrt{13} \cdot 4\sqrt{2}$$

$$= \frac{5}{\sqrt{26}}$$



\vec{r}_1 and \vec{r}_2

~~\vec{r}_1 and \vec{r}_3~~

\vec{r}_2 and \vec{r}_3

Using cosine similarity to measure similarity, which pair of reviews are the most similar? If there are multiple pairs of reviews that are most similar, select them all.

$$\frac{4 \cdot 6 + 4 \cdot 4}{\sqrt{32} \cdot \sqrt{32}} = \frac{40}{8\sqrt{26}} = \frac{5}{\sqrt{26}}$$

You create a table called `gums` that only contains the chewing gum purchases of `df`, then you create a bag-of-words matrix called `bow` from the `name` column of `gums`. The `bow` matrix is stored as a DataFrame shown below:

412 · #013

→

	pur	gum	...	paperboard	80
0	0	1	...	0	1
1	0	1	...	1	1
...
38	0	0	...	0	0
39	0	0	...	0	1

]

You also have the following outputs:

```
>>> bow_df.sum(axis=0)
pur          5
gum         41
sugar        2
..          ..
90           4
paperboard   22
80           20
Length: 139
```

```
>>> bow_df.sum(axis=1)
0    21
1    22
...  ..
37   22
38   10
39   17
Length: 40
```

```
>>> bow_df.loc[0, 'pur']
0
```

```
>>> (bow_df['paperboard'] > 0).sum()
20
```

```
>>> bow_df['gum'].sum()
41
```

For each question below, write your answer as an unsimplified math expression (no need to simplify fractions or logarithms) in the space provided, or write "Need more information" if there is not enough information provided to answer the question.

Problem 8.1

What is the TF-IDF for the word "pur" in document 0?

$$TF-IDF = TF \cdot IDF$$

$$TF = 0$$

$$TF-IDF = 0$$

	pur	gum	...	paperboard	80
0	0	1	...	0	1
1	0	1	...	1	1
...
38	0	0	...	0	0
39	0	0	...	0	1

```
>>> bow_df.sum(axis=0)
```

```
pur      5
gum     41
sugar     2
..
90        4
paperboard 22
80       20
Length: 139
```

```
>>> bow_df.sum(axis=1)
```

```
pur      0  21
gum      1  22
sugar    2  22
..
90       37  22
paperboard 38  10
80       39  17
Length: 40
```

```
>>> bow_df.loc[0, 'pur']
```

```
0
```

```
>>> (bow_df['paperboard'] > 0).sum()
```

```
20
```

```
>>> bow_df['gum'].sum()
```

```
41
```

Problem 8.2

What is the TF-IDF for the word "gum" in document 0?

$$IDF = \log \left(\frac{\# \text{ of docs.}}{\# \text{ of docs w/ term}} \right)$$

$$TF = \frac{1}{21}$$

$$IDF = \log \left(\frac{40}{??} \right)$$

need more information

	pur	gum	...	paperboard	80
0	0	1	...	0	1
1	0	1	...	1	1
...
38	0	0	...	0	0
39	0	0	...	0	1

```
>>> bow_df.sum(axis=0)
```

```
pur      5
gum     41
sugar     2
..
90        4
paperboard 22
80       20
Length: 139
```

```
>>> bow_df.sum(axis=1)
```

```
0      21
1      22
2      22
..
37     22
38     10
39     17
Length: 40
```

```
>>> bow_df.loc[0, 'pur']
```

```
0
```

```
>>> (bow_df['paperboard'] > 0).sum()
```

```
20
```

```
>>> bow_df['gum'].sum()
```

```
41
```

Problem 8.3

What is the TF-IDF for the word "paperboard" in document 1?

log₂

$$TF = \frac{1}{22}$$

$$IDF = \log_2 \left(\frac{40}{20} \right)$$

$$TF-IDF = \frac{1}{22} \cdot 1 = \frac{1}{22}$$

← final answer

	pur	gum	...	paperboard	80
0	0	1	...	0	1
1	0	1	...	1	1
...
38	0	0	...	0	0
39	0	0	...	0	1

```
>>> bow_df.sum(axis=0)
```

```
pur      5
gum     41
sugar     2
..
90        4
paperboard 22
80       20
Length: 139
```

```
>>> bow_df.sum(axis=1)
```

```
0      21
1      22
2      22
..
37     22
38     10
39     17
Length: 40
```

```
>>> bow_df.loc[0, 'pur']
```

```
0
```

```
>>> (bow_df['paperboard'] > 0).sum()
```

```
20
```

```
>>> bow_df['gum'].sum()
```

```
41
```

→ 0,
↳ 1
← 0

Constant Model - Angela

constant model

Problem 9.1

Which of the following is closest to the constant prediction h^* that minimizes:

Choose a constant prediction h
same number h for all data points

$$L_{0,1}(h, y) = \frac{1}{n} \sum_{i=1}^n \begin{cases} 0 & y_i = h \\ 1 & y_i \neq h \end{cases}$$

minimize zero-one loss \uparrow single prediction

- 1
- 5
- 6
- 7
- 11
- 15
- 30

find the integer h that makes the fewest mistakes compared to the observed values of y_i

30

from histogram, is the mode of this distribution

Consider a dataset of n integers, y_1, y_2, \dots, y_n , whose histogram is given below:

multiple data points y_1, y_2, \dots, y_n

\rightarrow the empirical risk (average loss)

$$R(h) = \frac{1}{n} \sum_{i=1}^n L_{0,1}(h; y_i)$$

This counts the fraction of points for which h is not equal to y_i

\uparrow

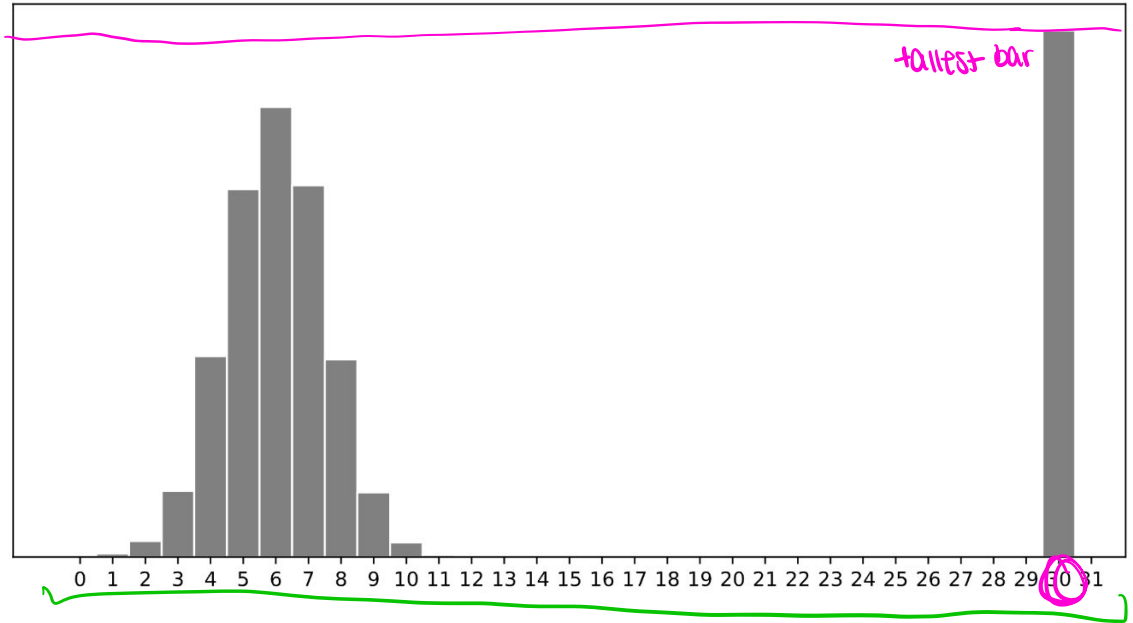
To minimize this count, choose h that appears most frequently

i.e. the mode

intuition: if h is the most common value (mode), you will

have the fewest points that differ from h , hence

the fewest mistakes.



Problem 9.2

Which of the following is closest to the constant prediction h^* that minimizes:

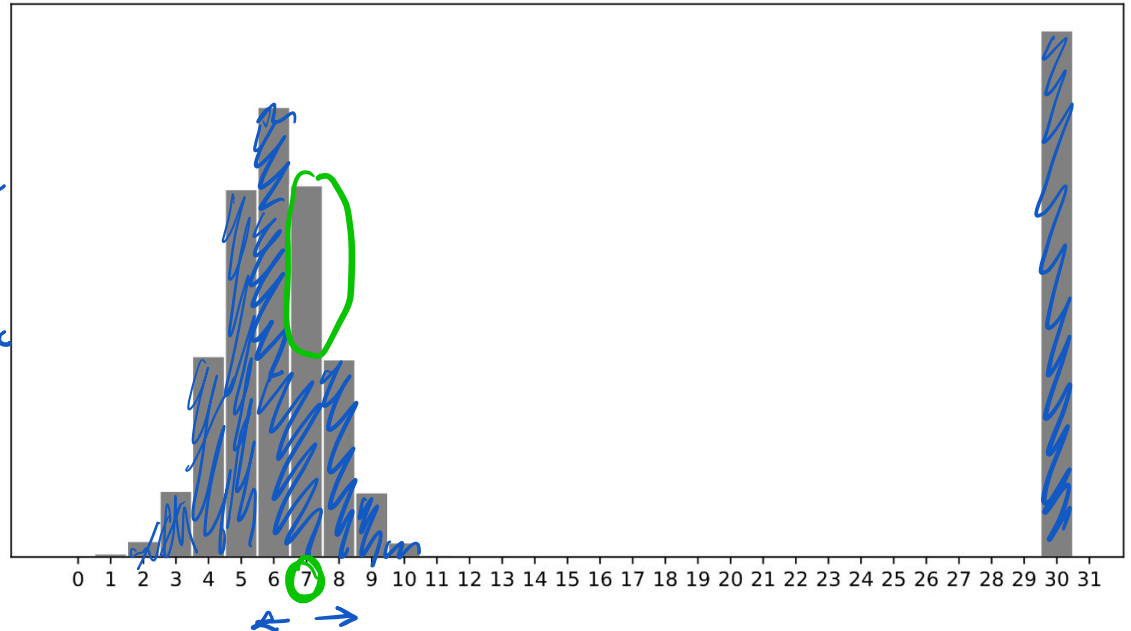
7

$$\frac{1}{n} \sum_{i=1}^n |y_i - h|$$

minimize absolute loss $|y-h|$

- 1
 - 5
 - 6
 - 7
 - 11
 - 15
 - 30
- minimizer for empirical risk for the constant model when using absolute loss is the median.

Consider a dataset of n integers, y_1, y_2, \dots, y_n , whose histogram is given below:



intuition: if you move h left of the median, you increase the absolute distance to all the data points on the right

→ if you move h right of the median, you decrease the absolute distance to all the data on the left.

“best balance” point is where half the data is on each side
→ median

Problem 9.3

Which of the following is closest to the constant prediction h^* that minimizes:

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

minimizes squared loss

- 1
- 5
- 6
- 7
- 11
- 15
- 30

minimizer for empirical risk for the constant model when using squared loss is the mean
mean > 7
11 next largest option

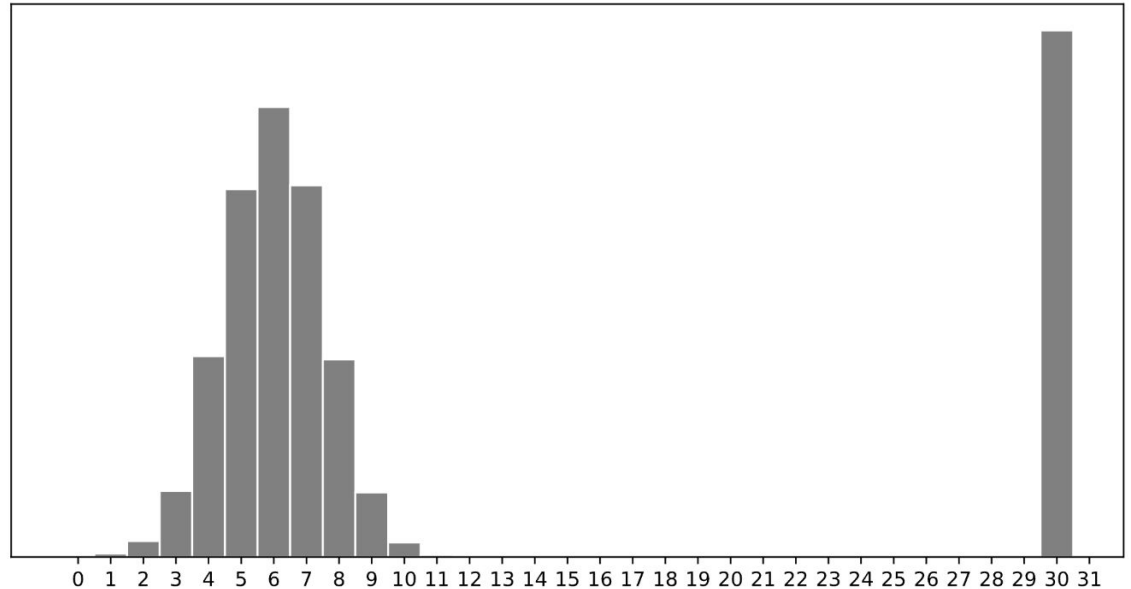
Consider a dataset of n integers, y_1, y_2, \dots, y_n , whose histogram is given below:

intuition: this function is minimized when

$$h = \frac{1}{n} \sum_{i=1}^n y_i \quad (\text{mean})$$

if you take the derivative of the sum of squared errors and set it to zero, you get this result.

median < mean since there is a large outlier skewing the distribution to the right.



Problem 9.4

Which of the following is closest to the constant prediction h^* that minimizes:

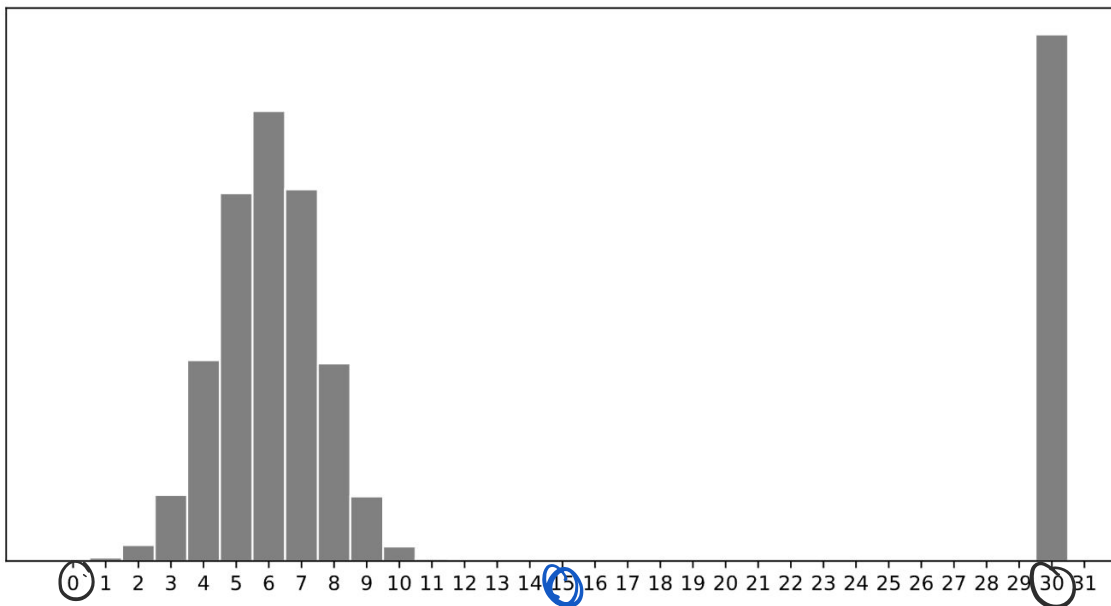
$$\frac{30-0}{2} = \boxed{15}$$

$$\lim_{p \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n |y_i - h|^p$$

minimize infinity loss

- 1
 - 5
 - 6
 - 7
 - 11
 - 15
 - 30
- minimizer of empirical risk for the constant model when using infinity loss is the midrange, i.e. halfway between the min and max

Consider a dataset of n integers, y_1, y_2, \dots, y_n , whose histogram is given below:



intuition: $\max_{i \in \mathbb{N}} |y_i - h|$ essentially means max absolute error

→ choose constant h that minimizes this worst-case (maximum) error.

→ center our prediction so that the worst error on the low side is exactly balanced by the worst error on the right side.

→ choose h as the midrange

left error: $\left| \frac{y_{\min} + y_{\max}}{2} - y_{\min} \right| = y_{\max} - y_{\min}$

right error: $\left| y_{\max} - \frac{y_{\min} + y_{\max}}{2} \right| = y_{\max} - y_{\min}$

both errors are equal, max error is as small as it can be given the spread of the data.

Regression - Angela

Problem 10

derivative of risk function

set to zero, then solve for optimal h^*

Consider a dataset that consists of y_1, \dots, y_n . In class, we used calculus to minimize mean squared error, $R_{sq}(h) = \frac{1}{n} \sum_{i=1}^n (h - y_i)^2$. In this problem, we want you to apply the same approach to a slightly different loss function defined below:

Loss: $L_{\text{midterm}}(y, h) = (\alpha y - h)^2 + \lambda h$

for a single data point y_i

Problem 10.1

Write down the empirical risk $R_{\text{midterm}}(\hat{h})$ by using the above loss function.

average loss across all data points y_1, \dots, y_n

$$R_{\text{midterm}}(h) = \frac{1}{n} \sum_{i=1}^n L_{\text{midterm}}(y_i, h)$$

$$\Rightarrow R_{\text{midterm}}(h) = \frac{1}{n} \sum_{i=1}^n [(\alpha y_i - h)^2 + \lambda h]$$

depends on i

$$\Rightarrow R_{\text{midterm}}(h) = \left[\frac{1}{n} \sum_{i=1}^n (\alpha y_i - h)^2 \right] + \lambda h$$

does not depend on i

Problem 10.2

The mean of dataset is \bar{y} , i.e. $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$. Find h^* that minimizes $R_{\text{midterm}}(h)$ using calculus. Your result should be in terms of \bar{y} , α and λ .

1.) take derivative of risk function
in respect to h

$$\frac{d}{dh} R_{\text{midterm}}(h) = \frac{d}{dh} \left[\frac{1}{n} \sum_{i=1}^n (\alpha y_i - h)^2 \right] + \lambda h$$

$$\frac{d}{dh} R_{\text{midterm}}(h) = \frac{1}{n} \sum_{i=1}^n 2(\alpha y_i - h)(-1) + \lambda$$

$$= -\frac{2}{n} \sum_{i=1}^n (\alpha y_i - h) + \lambda$$

$$= -\frac{2}{n} \cdot \left[\sum_{i=1}^n \alpha y_i - \sum_{i=1}^n h \right] + \lambda$$

$$= -\frac{2}{n} \cdot \left[\alpha \sum_{i=1}^n y_i - nh \right] + \lambda$$

$$= -\frac{2}{n} \cdot [\alpha n \bar{y} - nh] + \lambda$$

$$= -2\alpha \bar{y} + 2h + \lambda$$

$$= 2h - 2\alpha \bar{y} + \lambda = 0 \quad \text{set to 0}$$

$$L_{\text{midterm}}(y, h) = (\alpha y - h)^2 + \lambda h$$

$$\text{chain rule: } \frac{d}{dh} [f(g(h))] = f'(g(h)) \cdot g'(h)$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$n\bar{y} = \sum_{i=1}^n y_i$$

$$2h^* - 2\alpha \bar{y} + \lambda = 0$$

$$2h^* = 2\alpha \bar{y} - \lambda$$

$$h^* = \alpha \bar{y} - \frac{\lambda}{2}$$

intuition: term $\alpha \bar{y}$ would be the best constant if the loss were purely $\sum (\alpha y_i - h)^2$
→ extra $-\frac{\lambda}{2}$ shift comes from the linear penalty λh .