

# DSC 80 Discussion 6 Worksheet

Name: \_\_\_\_\_

## 1 FA23 Final Exam Problem 7

Alan set up a web page for his DSC 80 note with the following HTML:

```
<html><body>
<div id = "hero">DSC 80 NOTES</div>
  <div class="notes">
    <div class="notes">
      <p>Lecture 1: 5/5 stars!</p>
    </div>
    <div class="lecture_notes">
      <p>Lecture 2: 6/5 stars!!</p>
    </div>
  </div>
  <div class="lecture">
    <p>Lecture 3: 10/5 stars!!!!</p>
  </div>
</body></html>
```

Assume that the web page is parsed into a BeautifulSoup object named `soup`. Fill in each of the expressions below to evaluate to the desired string. Pay careful attention to the indexes after each call to `find_all`!

Desired string: "Lecture 1: 5/5 stars!"

`soup.find_all(_____)[0].text`

Desired string: "Lecture 2: 6/5 stars!!"

`soup.find_all(_____)[3].text`

Desired string: "Lecture 3: 10/5 stars!!!!"

`soup.find_all(_____)[1].text`

## 2 Spring 2023 Final Exam Problem 1

Consider the following line. Choose the regex pattern that, when filled in the blank, will return the desired matches.

`re.findall(r'_____', 'my cat is hungry, concatenate!, catastrophe! What a cat!')`

Output: ['my', 'a']

- [ ] `\b([a-z]*)cat\b`
- [ ] `\b[a-z]*\scat\b`
- [ ] `([a-z]*)\scat\b`
- [ ] `\b([a-z]*\scat)\b`

Output: ['concatenate']

- [ ] `\b.*cat.*\b`
- [ ] `[a-z]*cat[a-z]*`
- [ ] `[a-z]+cat[a-z]+`
- [ ] `\b[a-z]*cat[a-z]*\b`

Output: ['cat', 'concatenate', 'catastrophe', 'cat']

- [ ] `.*cat.*`
- [ ] `\b.*cat.*\b`
- [ ] `\b[a-z]*cat[a-z]*\b`
- [ ] `\b[a-z]+cat[a-z]+\b`

### 3 WI23 Final Exam Problem 4

Nicole decides to read research papers on data science. While reading these papers, she notices that there are many citations interspersed that refer to other research papers, and she'd like to read the cited papers as well. An excerpt from one such paper is stored in the string `s` below.

```
s = '''In DSC 10 [3], you learned about babypandas, a strict subset
of pandas [15][4]. It was designed [5] to provide programming
beginners [3][91] just enough syntax to be able to perform
meaningful tabular data analysis [8] without getting lost in
100s of details.'''
```

For each expression below, select the list it evaluates to, or select "None of the above."

```
list1 = ['10', '100']
list2 = ['3', '15', '4', '5', '3', '91', '8']
list3 = ['10', '3', '15', '4', '5', '3', '91', '8', '100']
list4 = ['[3]', '[15]', '[4]', '[5]', '[3]', '[91]', '[8]']
list5 = ['1', '0', '3', '1', '5', '4', '5', '3',
        '9', '1', '8', '1', '0', '0']
```

<code>re.findall(r'\d+', s)</code>	<code>re.findall(r'[\d+]', s)</code>	<code>re.findall(r'\[(\d+)\]', s)</code>
<input type="checkbox"/> list1	<input type="checkbox"/> list1	<input type="checkbox"/> list1
<input type="checkbox"/> list2	<input type="checkbox"/> list2	<input type="checkbox"/> list2
<input type="checkbox"/> list3	<input type="checkbox"/> list3	<input type="checkbox"/> list3
<input type="checkbox"/> list4	<input type="checkbox"/> list4	<input type="checkbox"/> list4
<input type="checkbox"/> list5	<input type="checkbox"/> list5	<input type="checkbox"/> list5

### 4 WI23 Final Exam Problem 5

Suppose `soup` is a BeautifulSoup object instantiated using the following HTML document.

```
<college>Your score is ready!</college>
<sat verbal="ready" math="ready">Your percentiles are as follows:
  <scorelist listtype="percentiles">
    <scorerow kind="verbal" subkind="per">Verbal: <scorenum>84</scorenum></
    scorerow>
    <scorerow kind="math" subkind="per">Math: <scorenum>99</scorenum></scorerow>
  </scorelist>
  And your actual scores are as follows:
  <scorelist listtype="scores">
    <scorerow kind="verbal"> Verbal: <scorenum>680</scorenum> </scorerow>
    <scorerow kind="math"> Math: <scorenum>800</scorenum> </scorerow>
  </scorelist></sat>
```

Which expression evaluates to "verbal"?

- `soup.find("scorerow").get("kind")`
- `soup.find("sat").get("ready")`
- `soup.find("scorerow").text.split(":")[0].lower()`
- `[s.get("kind") for s in soup.find_all("scorerow")][-2]`
- `soup.find("scorelist", attrs="listtype":"scores").get("kind")`
- None of the above

# DSC 80 Discussion 7 Worksheet

## 1 FA23 Final Exam Problem 8

Consider the following corpus:

Document number	Content
1	yesterday rainy today sunny
2	yesterday sunny today sunny
3	today rainy yesterday today
4	yesterday yesterday today today

Using a bag-of-words representation, which two documents have the largest dot product?

Using a bag-of-words representation, what is the cosine similarity between documents 2 and 3?

Which words have a TF-IDF score of 0 for all four documents? Select all words that apply.

yesterday  rainy  today  sunny

## 2 WI23 Final Problem 7

We decide to build a classifier that takes in a state's demographic information and predicts whether, in a given year, a state's mean math score was greater than its mean verbal score (1), or a state's mean math score was less than or equal to its mean verbal score (0). The simplest possible classifier we could build is one that predicts the same label (1 or 0) every time, independent of all other features.

*If  $a > b$ , then the constant classifier that maximizes training accuracy predicts 1 every time; otherwise, it predicts 0 every time.*

For which combination of a and b is the above statement not guaranteed to be true? **Select one.**

- $a = (\text{sat}[\text{'Math'}] > \text{sat}[\text{'Verbal'}]).\text{mean}(); b = 0.5$
- $a = (\text{sat}[\text{'Math'}] - \text{sat}[\text{'Verbal'}]).\text{mean}(); b = 0$
- $a = (\text{sat}[\text{'Math'}] - \text{sat}[\text{'Verbal'}] > 0).\text{mean}(); b = 0.5$
- $a = ((\text{sat}[\text{'Math'}] / \text{sat}[\text{'Verbal'}]) > 1).\text{mean}() - 0.5; b = 0$

Suppose we train a classifier that achieves an accuracy of 5/9 on our training set. Typically, RMSE is used as a performance metric for regression models, but mathematically, nothing is stopping us from using it for classification models as well. What is the RMSE of our classifier on our training set?

### 3 SP23 Final Problem 5.3

Chen downloaded 4 reviews of a new vacuum cleaner from Amazon (as shown in the 4 sentences below).

Sentence 1: 'if i could give this vacuum zero stars i would'

Sentence 2: 'i will not order again this vacuum is garbage'

Sentence 3: 'Love Love Love i love this product'

Sentence 4: 'this little vacuum is so much fun to use i love it'

$X$  is the Term frequency-Inverse Document Frequency (TF-IDF) of the word `vacuum` in sentence 1. Chen replaces sentence 3 with the following new sentence/review.

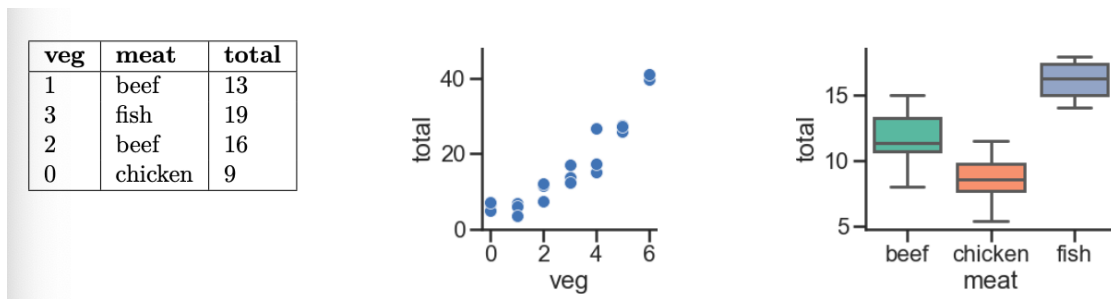
New Sentence 3: 'Love Love Love i love this vacuum'

$Y$  is the TF-IDF of the word `vacuum` in sentence 1 after the sentence 3 is replaced by the new sentence 3. Given the above information, which of the following statements is true? **Select one.**

- $X = Y$
- $X > Y > 0$
- $X > 0$  and  $Y = 0$
- $X > 0$  and  $Y > X$

### 4 FA23 Final Problem 9

Every week, Lauren goes to her local grocery store and buys a varying amount of vegetable but always exactly one pound of meat. We use a linear regression model to predict her total grocery bill. We've collected a dataset containing the pounds of vegetables bought, the type of meat bought, and the total bill. Below we display the first few rows of the dataset and two plots generated using the entire training set.



Suppose we fit the following linear regression models to predict `total`. Based on the data and visualizations above, for each model, determine whether **each fitted model coefficient**  $w^*$  is positive, negative, or exactly 0. The notation "meat=beef", for example, refers to the one-hot encoded `meat` column with value 1 if the original value in the `meat` column was `beef` and 0 otherwise.

1.  $H(x) = w_0 + w_1 \cdot \text{veg}$      $w_0$  \_\_\_\_\_  $w_1$  \_\_\_\_\_
2.  $H(x) = w_0 + w_1 \cdot \text{veg} + w_2 \cdot (\text{meat} = \text{beef}) + w_3 \cdot (\text{meat} = \text{chicken})$      $w_0$  \_\_\_\_\_  $w_1$  \_\_\_\_\_  $w_2$  \_\_\_\_\_  $w_3$  \_\_\_\_\_
3.  $H(x) = w_0 + w_1 \cdot \text{meat} = \text{chicken}$      $w_0$  \_\_\_\_\_  $w_1$  \_\_\_\_\_
4.  $H(x) = w_0 + w_1 \cdot (\text{meat} = \text{beef}) + w_2 \cdot (\text{meat} = \text{chicken}) + w_3 \cdot (\text{meat} = \text{fish})$      $w_0$  \_\_\_\_\_  $w_1$  \_\_\_\_\_  $w_2$  \_\_\_\_\_  $w_3$  \_\_\_\_\_

Suppose we fit  $H(x) = w_0 + w_1 \cdot \text{veg} + w_2 \cdot (\text{meat} = \text{beef}) + w_3 \cdot (\text{meat} = \text{fish})$ , and find that  $\vec{w}^* = [-3, 5, 8, 12]$ .

What is the prediction of this model on the **first** point in our dataset? \_\_\_\_\_