

▼ HW 4: Spam/Ham Classification

Due Date: 5/20 (Fri), 11:59 PM

Collaboration Policy

Data science is a collaborative activity. While you may talk with others about the project, we ask that you **write your solutions individually**. If you do discuss the assignments with others please **include their names** at the top of your notebook.

Collaborators: *list collaborators here*

This Assignment

In this homework, you will use what you've learned in class to create a classifier that can distinguish spam (junk or commercial or bulk) emails from ham (non-spam) emails. In addition to providing some skeleton code to fill in, we will evaluate your work based on your model's accuracy and your written responses in this notebook.

After this homework, you should feel comfortable with the following:

- Part 1: Feature engineering with text data
- Part 2: Using sklearn libraries to process data and fit models
- Part 3: Validating the performance of your model and minimizing overfitting
- Part 3: Generating and analyzing precision-recall curves

Warning!

We've tried our best to filter the data for anything blatantly offensive as best as we can, but unfortunately there may still be some examples you may find in poor taste. If you encounter these examples and believe it is inappropriate for students, please let a TA know and we will try to remove it for future semesters. Thanks for your understanding!

Score Breakdown

Question	Points
1a	2
1b	2
1c	2
2	3
3	3
4	3

Question	Points
5	3
6a	3
6b	3
6c	3
6d	3
7	4
8	4
Total	38

▼ Part I - Initial Analysis

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5
6 import seaborn as sns
7 sns.set(style = "whitegrid",
8         color_codes = True,
9         font_scale = 1.5)
10
11 class bcolor:
12     BLACK = 'W033[40m'
13     YELLOW = 'W033[93m'
14     RED = 'W033[91m'
15     BOLD = 'W033[1m'
16     END = 'W033[0m'
17
18 def print_passed(str_in):
19     print(bcolor.BLACK + bcolor.YELLOW + bcolor.BOLD + str_in + bcolor.END)

```

▼ Mount your Google Drive

When you run a code cell, Colab executes it on a temporary cloud instance. Every time you open the notebook, you will be assigned a different machine. All compute state and files saved on the previous machine will be lost. Therefore, you may need to re-download datasets or rerun code after a reset. Here, you can mount your Google drive to the temporary cloud instance's local filesystem using the following code snippet and save files under the specified directory (note that you will have to provide permission every time you run this).

```

1 # mount Google drive
2 from google.colab import drive
3 drive.mount('/content/drive')

```

```

4
5 # now you can see files
6 !echo -e "Number of Google drive files in /content/drive/My Drive/:"
7 !ls -l "/content/drive/My Drive/" | wc -l
8 # by the way, you can run any linux command by putting a ! at the start of the line
9
10 # by default everything gets executed and saved in /content/
11 !echo -e "Current directory:"
12 !pwd

```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/c

Number of Google drive files in /content/drive/My Drive/:
167

Current directory:
/content



```

1 workspace_path = '/content/drive/MyDrive/hw4/' # Change this path!
2 print(f'Current Workspace: {workspace_path}')

```

Current Workspace: /content/drive/MyDrive/hw4/

▼ Loading in the Data

Our goal is to classify emails as spam or not spam (referred to as "ham") using features generated from the text in the email.

The dataset consists of email messages and their labels (0 for ham, 1 for spam). Your labeled training dataset contains 8348 labeled examples, and the test set contains 1000 unlabeled examples.

Run the following cells to load in the data into DataFrames.

The `train` DataFrame contains labeled data that you will use to train your model. It contains four columns:

1. `id`: An identifier for the training example
2. `subject`: The subject of the email
3. `email`: The text of the email
4. `spam`: 1 if the email is spam, 0 if the email is ham (not spam)

The `test` DataFrame contains 1000 unlabeled emails. You will predict labels for these emails and submit your predictions to Kaggle for evaluation.

```

1 original_training_data = pd.read_csv(f'{workspace_path}/train.csv')
2 test = pd.read_csv(f'{workspace_path}/test.csv')
3
4 # Convert the emails to lower case as a first step to processing the text
5 original_training_data['email'] = original_training_data['email'].str.lower()

```

```

6 test['email'] = test['email'].str.lower()
7
8 original_training_data.head()
9

```

	id	subject	email	spam
0	0	Subject: A&L Daily to be auctioned in bankrupt...	url: http://boingboing.net/#85534171Wn date: n...	0
1	1	Subject: Wired: "Stronger ties between ISPs an...	url: http://scriptingnews.userland.com/backiss...	0
2	2	Subject: It's just too small ...	<html>Wn <head>Wn </head>Wn <body>Wn <font siz...	1
-	-	-	depends on how much over spending vs.	-

▼ Question 1a

First, let's check if our data contains any missing values.

- Step1: Fill in the cell below to print the number of NaN values in each column. **Hint:** [pandas.isnull](#)
- Step2: If there are NaN values, replace them with appropriate filler values (i.e., NaN values in the `subject` or `email` columns should be replaced with empty strings).
- Step3: Print the number of NaN values in each column after this modification to verify that there are no NaN values left.

```

1 # BEGIN YOUR CODE
2 # -----
3 print('Before imputation:')
4 print(original_training_data.isnull().sum())
5 modified_original_training_data = original_training_data.fillna('')
6 print('-----')
7 print('After imputation:')
8 print(modified_original_training_data.isnull().sum())
9 original_training_data = modified_original_training_data
10 # -----
11 # END YOUR CODE

```

Before imputation:

```

id      0
subject 6
email    0
spam     0
dtype: int64

```

After imputation:

```

id      0
subject 0
email    0
spam     0
dtype: int64

```

```
1 assert original_training_data.isnull().sum().sum() == 0
2 print_passed('Q1a: Passed all unit tests!')
```

Q1a: Passed all unit tests!

Explanation: To print the total number of NaN values in a column I had to use the sum function. With the isnull hint provided to me, I used isnull() then the sum to find the number of NaN values in each column of original_training_data.

Step 2 -> to replace the NaN values with empty strings I used .fillna(""), which fills NaN values with " (an empty string).

Step 3-> After modification, I printed the modified version.

▼ Question 1b

In the cell below, print the text of the first ham (i.e. 1st row) and the first spam email in the original training set.

```
1 # BEGIN YOUR CODE
2 # -----
3 first_ham = original_training_data.loc[(original_training_data['spam'] == 0) , 'email'].iloc[0]
4 first_spam = original_training_data.loc[(original_training_data['spam'] == 1) , 'email'].iloc[0]
5 # -----
6 # END YOUR CODE
7
8 print('The text of the first Ham:')
9 print('-----')
10 print(first_ham)
11
12 print('The text of the first Spam:')
13 print('-----')
14 print(first_spam)
```

The text of the first Ham:

url: <http://boingboing.net/#85534171>
date: not supplied

arts and letters daily, a wonderful and dense blog, has folded up its tent due to the bankruptcy of its parent company. a&l daily will be auctioned off by the receivers. link[1] discuss[2] (_thanks, misha!_)

[1] <http://www.aldaily.com/>

[2] <http://www.quicktopic.com/boing/h/zlfterjnd6jf>

The text of the first Spam:

```
<html>
<head>
```

```

</head>
<body>
<font size=3d"4"><b> a man endowed with a 7-8" hammer is simply<br>
  better equipped than a man with a 5-6"hammer. <br>
<br>would you rather have<br>more than enough to get the job done or fall =
short. it's totally up<br>to you. our methods are guaranteed to increase y=
our size by 1-3"<br> <a href=3d"http://209.163.187.47/cgi-bin/index.php?10=
004">come in here and see how</a>
</body>
</html>

```

```

1 assert len(first_ham) == 359 and len(first_spam) == 444
2 print_passed('Q1b: Passed all unit tests!')

```

Q1b: Passed all unit tests!

Explanation: Before finding the first ham or spam in the training set, I first had to create two series where one contained all the ham emails while the other had all the spam emails. Hence, I used loc Boolean Arrays to get my desired values. For ham emails the spam column consisted of 0 (*original_training_data['spam'] == 0*) while spam emails are 1 (*original_training_data['spam'] == 1*). Also, for both, I needed the email column so I had to have email after a comma.

Next, to print the first row of the series, all I had to do was `iloc[0]` since that prints the first row.

▼ Question 1c

Discuss one thing you notice that is different between the two emails that might relate to the identification of spam.

Answer: First off, the first spam seemed to be much longer than the first ham email however I believe this doesn't really matter. So, the second observation I made was that the first spam was HTML format while the ham emails seemed to have a normal text format. So I believe the weird format of emails can give us a clue that it may be a spam email.

▼ Training Validation Split

The training data is all the data we have available for both training models and **validating** the models that we train. We therefore need to split the training data into separate training and validation datasets. You will need this **validation data** to assess the performance of your classifier once you are finished training.

Note that we set the seed (`random_state`) to 42. This will produce a pseudo-random sequence of random numbers that is the same for every student. **Do not modify this in the following questions, as our tests depend on this random seed.**

```
1 from sklearn.model_selection import train_test_split
2
3 train, val = train_test_split(
4     original_training_data, test_size=0.1, random_state=42)

1 print(train.shape, val.shape)    # 더해서 8342 맞음

(7513, 4) (835, 4)
```

▼ Basic Feature Engineering

We would like to take the text of an email and predict whether the email is **ham** or **spam**. This is a *classification* problem, and here we use logistic regression to train a classifier.

Recall that to train an logistic regression model we need:

- a numeric feature matrix X
- a vector of corresponding binary labels y .

Unfortunately, our data are text, not numbers. To address this, we can create numeric features derived from the email text and use those features for logistic regression:

- Each row of X is an email.
- Each column of X contains one feature for all the emails.

We'll guide you through creating a simple feature, and you'll create more interesting ones when you are trying to increase your accuracy.

▼ Question 2

Create a function called `words_in_texts` that takes in a list of `words` and a pandas Series of email `texts`. It should output a 2-dimensional NumPy array containing one row for each email text. The row should contain either a 0 or a 1 for each word in the list: 0 if the word doesn't appear in the text and 1 if the word does. For example:

```
>>> words_in_texts(['hello', 'bye', 'world'],
                    pd.Series(['hello', 'hello worldhello']))

array([[1, 0, 0],
       [1, 0, 1]])
```

Hint: [pandas.Series.str.contains](#)

The provided tests make sure that your function works correctly, so that you can use it for future questions.

```

1 def words_in_texts(words, texts):
2     '''
3     Args:
4         words (list-like): words to find
5         texts (Series): strings to search in
6
7     Returns:
8         NumPy array of 0s and 1s with shape (n, p) where n is the
9         number of texts and p is the number of words.
10    '''
11    # BEGIN YOUR CODE
12    # -----
13    indicator_array = (np.asarray([texts.str.contains(words) for words in words]).T ) * 1
14    # -----
15    # END YOUR CODE
16
17    return indicator_array

```

```

1 assert np.allclose(
2     words_in_texts(
3         ['hello', 'bye', 'world'],
4         pd.Series(['hello', 'hello worldhello'])),
5     np.array([[1, 0, 0], [1, 0, 1]]))
6
7 assert np.allclose(
8     words_in_texts(
9         ['a', 'b', 'c', 'd', 'e', 'f', 'g'],
10        pd.Series(['a b c d e f g', 'a', 'b', 'c', 'd e f g', 'h', 'a h'])),
11    np.array(
12        [[1,1,1,1,1,1,1],
13         [1,0,0,0,0,0,0],
14         [0,1,0,0,0,0,0],
15         [0,0,1,0,0,0,0],
16         [0,0,0,1,1,1,1],
17         [0,0,0,0,0,0,0],
18         [1,0,0,0,0,0,0]]))
19 print_passed('Q2: Passed all unit tests!')

```

Q2: Passed all unit tests!

Explanation: Using the numpy as np, I was able to create an array with just np.asarray(). Inside that np.asarray argument I used Series.str.contains(pattern) which returns a boolean series based on whether the pattern is included in the Series.

The text (the strings to search in) is the Series and 'words' is the pattern (words to find). So if words inside 'words' was in the text(Series) then it would be 1 or else 0. Whenever there is more

than one word in "words" then we have to cover all of it, hence I used the for-in loop to do that.

Next, because it was a boolean series the value of whether the words appeared in the text will be displayed in columns. In other words, the return value is (p,n). So to change this to a desired (n,p) format I had to transpose the function, which can be done by adding ".T"

Finally, I had to multiply the entire array by 1 to convert True/False to 1s and 0s

▼ Basic EDA

We need to identify some features that allow us to distinguish spam emails from ham emails.

One idea is to compare the distribution of a single feature in spam emails to the distribution of the same feature in ham emails.

If the feature is itself a binary indicator (such as whether a certain word occurs in the text), this amounts to comparing the proportion of spam emails with the word to the proportion of ham emails with the word.

```
1 from IPython.display import display, Markdown
2 df = pd.DataFrame({
3     'word_1': [1, 0, 1, 0],
4     'word_2': [0, 1, 0, 1],
5     'type': ['spam', 'ham', 'ham', 'ham']
6 })
7 display(Markdown("> Our Original DataFrame has some words column and a type column. You can thin
8 display(df);
9 display(Markdown("> `melt` will turn columns into variable, notice how `word_1` and `word_2` beco
10 display(df.melt("type"))
```

Our Original DataFrame has some words column and a type column. You can think of each row is a sentence, and the value of 1 or 0 indicates the number of occurrences of the word in this sentence.

word_1	word_2	type	
0	1	0	spam

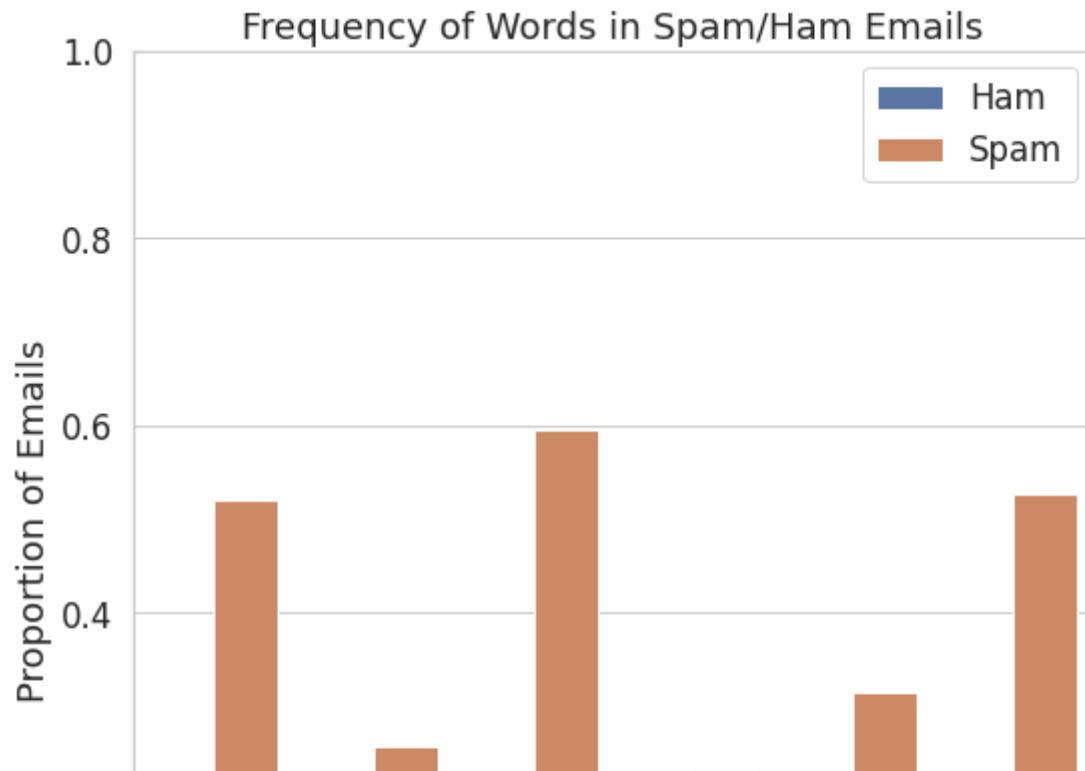
We can create a bar chart like the one above comparing the proportion of spam and ham emails containing certain words. Choose a set of words that are different from the ones above, but also have different proportions for the two classes. Make sure that we only consider emails from train.

their values are stored in the value column

```

1 # We must do this in order to preserve the ordering of emails to labels for words_in_texts
2 train=train.reset_index(drop=True)
3
4 some_words = ['body', 'html', 'please', 'money', 'business', 'offer']
5 Phi_train = words_in_texts(some_words, train['email'])
6
7 df = pd.DataFrame(data = Phi_train, columns = some_words)
8 df['label'] = train['spam']
9
10 plt.figure(figsize=(8,8))
11 sns.barplot(x = "variable",
12             y = "value",
13             hue = "label",
14             data = (df
15                     .replace({'label':
16                               {0 : 'Ham',
17                                1 : 'Spam'}})
18                     .melt('label')
19                     .groupby(['label', 'variable'])
20                     .mean()
21                     .reset_index()))
22
23 plt.ylim([0, 1])
24 plt.xlabel('Words')
25 plt.ylabel('Proportion of Emails')
26 plt.legend(title = "")
27 plt.title("Frequency of Words in Spam/Ham Emails")
28 plt.tight_layout()
29 plt.show()

```



▼ Question 3

When the feature is binary, it makes sense to compare its proportions across classes (as in the previous question). Otherwise, if the feature can take on numeric values, we can compare the distributions of these values for different classes.

Create a *class conditional density plot* like the one above (using `sns.distplot`), comparing the distribution of the length of spam emails to the distribution of the length of ham emails in the training set. Set the x-axis limit from 0 to 50000.

```

1 # BEGIN SOLUTION
2
3 spam_emails = train[train['spam'] == 0]['email'].str.len()
4 ham_emails = train[train['spam'] == 1]['email'].str.len()
5
6 sns.distplot(spam_emails, hist = False, label = 'ham')
7 sns.distplot(ham_emails, hist = False, label = 'spam')
8
9 plt.xlim(0, 50000)
10 plt.legend();
11 plt.title("Distribution of length of spam emails and ham emails in training set")
12 plt.show
13 # END SOLUTION

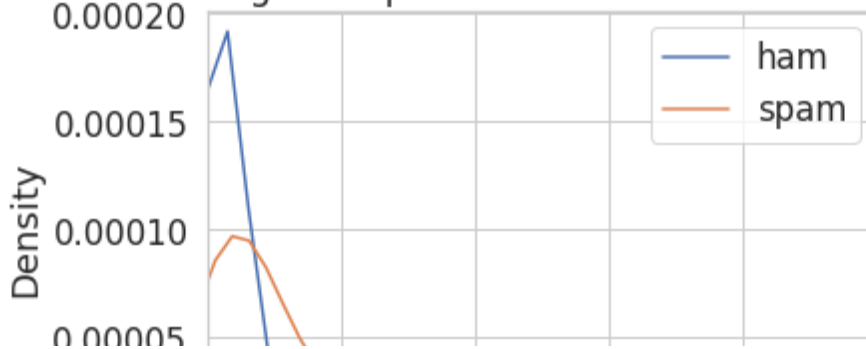
```

```

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `dist
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `dist
warnings.warn(msg, FutureWarning)
<function matplotlib.pyplot.show>

```

Distribution of length of spam emails and ham emails in training set



Explanation: Imported from seaborn, the distplot had to have the 'sns' to properly form the conditional density plot.

Next, to compare the length of each email I used `['email'].str.len()` which computes the string length of all the elements in the email column and labels the x-axis as 'email'. But before that, since we are talking about the spam and ham emails, I had to use the source field to properly retrieve the appropriate data, so if spam = 0 then it's ham email and vice versa.

To get rid of the bar, I used `hist = False`, and finally used the label parameter for the legend.

Because the instructions were to set the x limit from 0 to 50000, I used `xlim(0,50000)` to do that and later completed the plot with a title and a legend.

▼ Basic Classification

Notice that the output of `words_in_texts(words, train['email'])` is a numeric matrix containing features for each email. This means we can use it directly to train a classifier!

▼ Question 4

We've given you 5 words that might be useful as features to distinguish spam/ham emails. Use these words as well as the `train` DataFrame to create two NumPy arrays: `X_train` and `Y_train`.

- `X_train` should be a matrix of 0s and 1s created by using your `words_in_texts` function on all the emails in the training set.
- `Y_train` should be a vector of the correct labels for each email in the training set.

The provided tests check that the dimensions of your feature matrix (X) are correct, and that your features and labels are binary (i.e. consists of 0 and 1, no other values). It does not check that your function is correct; that was verified in a previous question.

```
1 some_words = ['drug', 'bank', 'prescription', 'memo', 'private']
```

```

2
3 # BEGIN YOUR CODE
4 # -----
5 X_train = words_in_texts(some_words, train['email'])
6 Y_train = np.array(train['spam'])
7 # -----
8 # END YOUR CODE
9
10 X_train[:5], Y_train[:5]

(array([[0, 0, 0, 0, 0],
        [0, 0, 0, 0, 0],
        [0, 0, 0, 0, 0],
        [0, 0, 0, 0, 0],
        [0, 0, 0, 1, 0]]), array([0, 0, 0, 0, 0]))

```

Explanation: X_train is for testing the words_in_texts function so for input I placed some_words, and the text to search the words in is train['email'].

Y_train is the vector of the **correct labels** for each email, so it should be an array of spam emails in the training data set.

```

1 assert X_train.shape == (7513, 5)
2 assert len(np.unique(X_train)) == 2
3 assert len(np.unique(Y_train)) == 2
4
5 print_passed('Q4: Passed all unit tests!')

```

Q4: Passed all unit tests!

▼ Question 5

Now we have matrices we can give to scikit-learn!

- Using the [LogisticRegression](#) classifier, train a logistic regression model using X_train and Y_train.
- Then, output the accuracy of the model (on the training data) in the cell below. You should get an accuracy around 75%.

```

1 from sklearn.linear_model import LogisticRegression
2
3 # BEGIN YOUR CODE
4 # -----
5 model = LogisticRegression().fit(X_train, Y_train)
6 training_accuracy = model.score(X_train, Y_train)
7 Y_predict = model.predict(X_train)
8 # -----
9 # END YOUR CODE
10
11 print("Training Accuracy: ", training_accuracy)

```

Training Accuracy: 0.7576201251164648

Explanation: After using the Logistic Regression classifier, I used `.fit` (which is shown in the example of the scikit-learn website) to fit and train the model according to x training data(`X_train`) and y training data (`Y_train`).

The accuracy should also be based on the two training data sets, as a result I used the `.score method` to return the mean accuracy of the data.

Next, I added the `.predict(X_train)`, which predicts values(or labels) for samples in `X_train`.

The accuracy was 75.76% which is around 75%.

```
1 assert training_accuracy > 0.72
2 print_passed('Q5: Passed all unit tests!')
```

Q5: Passed all unit tests!

▼ Evaluating Classifiers

That doesn't seem too shabby! But the classifier you made above isn't as good as this might lead us to believe. First, we are evaluating accuracy on the training set, which may lead to a misleading accuracy measure, especially if we used the training set to identify discriminative features. In future parts of this analysis, it will be safer to hold out some of our data for model validation and comparison.

Presumably, our classifier will be used for **filtering**, i.e. preventing messages labeled `spam` from reaching someone's inbox. There are two kinds of errors we can make:

- False positive (FP): a ham email gets flagged as spam and filtered out of the inbox.
- False negative (FN): a spam email gets mislabeled as ham and ends up in the inbox.

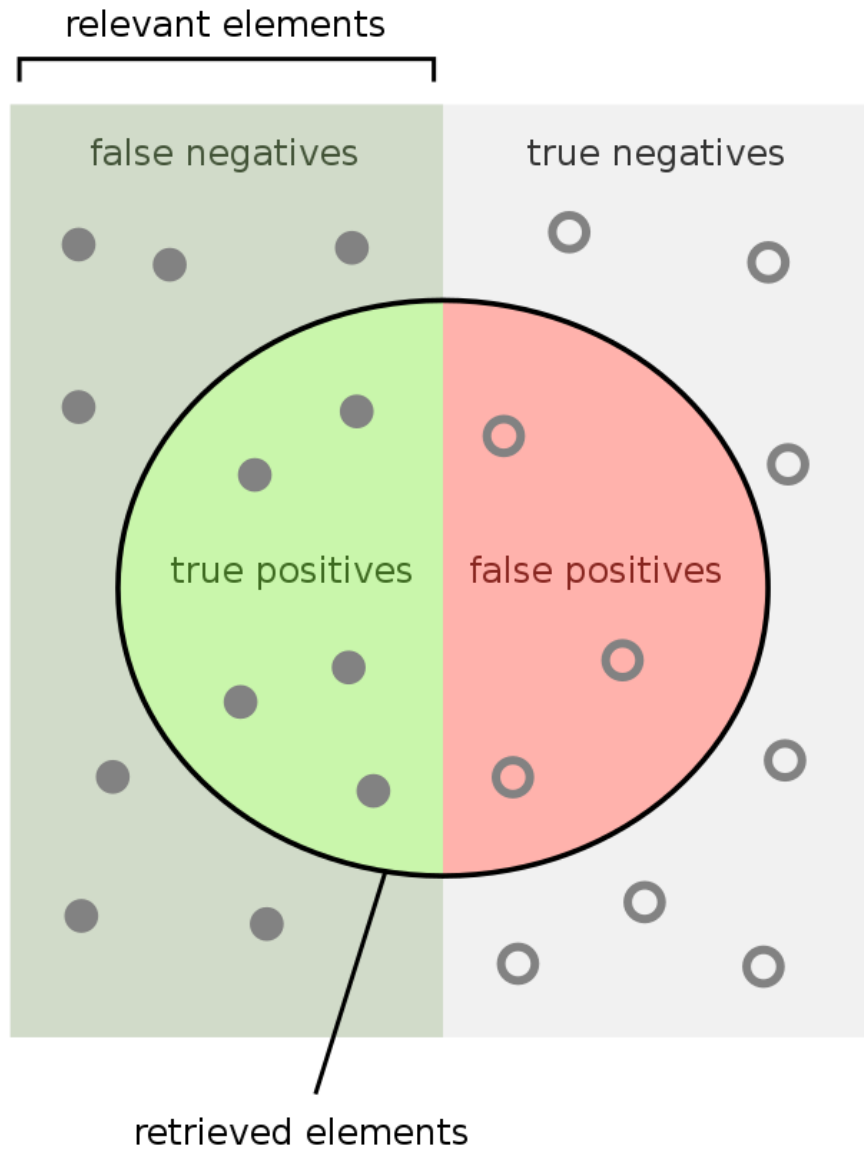
These definitions depend both on the true labels and the predicted labels. False positives and false negatives may be of differing importance, leading us to consider more ways of evaluating a classifier, in addition to overall accuracy:

Precision measures the proportion $\frac{TP}{TP+FP}$ of emails flagged as spam that are actually spam.

Recall measures the proportion $\frac{TP}{TP+FN}$ of spam emails that were correctly flagged as spam.

False-alarm rate measures the proportion $\frac{FP}{FP+TN}$ of ham emails that were incorrectly flagged as spam.

The following image might help:



How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

Note that a true positive (TP) is a spam email that is classified as spam, and a true negative (TN) is a ham email that is classified as ham.

▼ Question 6a

Suppose we have a classifier `zero_predictor` that always predicts 0 (never predicts positive). How many false positives and false negatives would this classifier have if it were evaluated on the training set and its results were compared to `Y_train`? Fill in the variables below (answers can be hard-coded):

```
1 # BEGIN YOUR CODE
2 # -----
3 zero_predictor_fp = 0
4 zero_predictor_fn = sum(Y_train == 1)
5 # -----
6 # END YOUR CODE
```

Explanation: Since the `zero_predictor` will never predict a positive(spam), there is 0 possibility that the classifier would have a false positive.

Since False Negative means the email is spam but is actually labeled/predicted as ham, the possibility of FN is equal to the sum of when `Y_train` data is correctly labeled as spam, which is 1.

```
1 assert zero_predictor_fp + zero_predictor_fn == 1918
2 print_passed('Q6a: Passed all unit tests!')
```

Q6a: Passed all unit tests!

▼ Question 6b

What are the accuracy and recall of `zero_predictor` (classifies every email as ham) on the training set? Do NOT use any `sklearn` functions.

```
1 # BEGIN YOUR CODE
2 # -----
3 zero_predictor_acc = sum(Y_train == 0) / len(Y_train)
4 zero_predictor_recall = 0 / ( 0 + sum(Y_train == 1) )
5
6 print(zero_predictor_acc)
7 # -----
8 # END YOUR CODE
```

0.7447091707706642

Explanation: First off, the *accuracy's formula* is $(\text{True Positive} + \text{True Negative}) / (n)$. Since the `zero_predictor` only classifies negative(never assumes positive), the True Positive is an automatic 0. So only the True Negative and 'n' remain. True Negative means the predicted

negative is actually correct, so the correct label should equal 0, so it is replaced by `sum(Y_train = 0)`. On the other hand, n is the total data so is then replaced by `len(Y_train)`.

Second, recall formula is $\text{True Positive} / (\text{True Positive} + \text{False Negative})$. **This one was quite simple because earlier we said True Positive is 1 and False Negative is `*sum(Y_train == 1)`*

```
1 assert np.isclose(zero_predictor_acc + zero_predictor_recall, 0.7447091707706642)
2 print_passed('Q6b: Passed all unit tests!')
```

Q6b: Passed all unit tests!

▼ Question 6c

Compute the precision, recall, and false-alarm rate of the `LogisticRegression` classifier created and trained in Question 5. **Note: Do NOT use any sklearn built-in functions.**

```
1 # BEGIN YOUR CODE
2 # -----
3 TP = sum((Y_predict == Y_train) & (Y_train == 1))
4 FP = sum((Y_predict != Y_train) & (Y_train == 0))
5 TN = sum((Y_predict == Y_train) & (Y_train == 0))
6 FN = sum((Y_predict != Y_train) & (Y_train == 1))
7
8 logistic_predictor_precision = TP / (TP + FP)
9 logistic_predictor_recall = TP / (TP + FN)
10 logistic_predictor_far = FP / (FP + TN)
11 # -----
12 # END YOUR CODE
```

Explanation: For this, I had to clearly define each one of the cases.

True Positive: predicted as Positive and is actually Positive. So, the predicted value is equal to the correct(actual) value -> `(Y_predict == Y_train)` while the correct(actual) value must be Positive (1), so `(& (Y_train == 1))`. Since we are going for the total number of cases, it is important to add the sum, or else the result would be an array.

False Positive: predicted as Positive but is actually Negative. So the predicted value should not equal the correct(actual) value -> `(Y_predict != Y_train)` while the correct(actual) value is Negative (0), so `(& (Y_train == 0))`.

True Negative: predicted as Negative and is actually a Negative. So predicted should equal the correct value, and the correct value(or labels) should be 0

False Negative: predicted as Negative but is actually a Positive. So predicted should *not equal* the correct value, and the correct value(or label) should be 1

After setting all of the cases, I inserted the appropriate cases inside the formulas of precision, recall, and false alarm rate.

```
1 assert np.isclose(logistic_predictor_precision, 0.6422287390029325)
2 assert np.isclose(logistic_predictor_recall, 0.11418143899895725)
3 assert np.isclose(logistic_predictor_far, 0.021805183199285077)
4 print_passed('Q6c: Passed all unit tests!')
```

Q6c: Passed all unit tests!

▼ Question 6d

1. Our logistic regression classifier got 75.6% prediction accuracy (number of correct predictions / total). How does this compare with predicting 0 for every email?
2. Given the word features we gave you above, name one reason this classifier is performing poorly. Hint: Think about how prevalent these words are in the email set.
3. Which of these two classifiers would you prefer for a spam filter and why? Describe your reasoning and relate it to at least one of the evaluation metrics you have computed so far.

Answer:

1. The zero predictor accuracy was 0.7447091707706642 or 74.47%. This can prove that our logistic regression classifier is more accurate than the zero predictors.
2. The number of words in the word features is not enough. We can't really distinguish spam and ham emails just by drug, bank, prescription, memo, and private. Both spam and ham emails might use the word memo a lot, which ruins the classifiers.
3. Even though the zero predictor classifier has a lower accuracy than that of the logistic regression classifier, the zero predictors will never wrongly classify a ham email as a spam email. I believe that knowing none of my important ham emails will get misclassified to spam is more important than the 1~2% accuracy increase of the logistic regression classifier.

▼ Part II - Moving Forward

With this in mind, it is now your task to make the spam filter more accurate. In order to get full credit on the accuracy part of this assignment, you must get at least **88%** accuracy on the validation set.

Here are some ideas for improving your model:

1. Finding better features based on the email text. Some example features are:
 1. Number of characters in the subject / body
 2. Number of words in the subject / body
 3. Use of punctuation (e.g., how many '!' were there?)
 4. Number / percentage of capital letters

5. Whether the email is a reply to an earlier email or a forwarded email
2. Finding better words to use as features. Which words are the best at distinguishing emails? This requires digging into the email text itself.
3. Better data processing. For example, many emails contain HTML as well as text. You can consider extracting out the text from the HTML to help you find better words. Or, you can match HTML tags themselves, or even some combination of the two.
4. Model selection. You can adjust parameters of your model (e.g. the regularization parameter) to achieve higher accuracy. Recall that you should use cross-validation to do feature and model selection properly! Otherwise, you will likely overfit to your training data.

You may use whatever method you prefer in order to create features, but **you are not allowed to import any external feature extraction libraries**. In addition, **you are only allowed to train logistic regression models**. No random forests, k-nearest-neighbors, neural nets, etc.

▼ Question 7: EDA

In the cell below, show a visualization that you used to select features for your model. Include both

1. A plot showing something meaningful about the data that helped you during feature / model selection.
2. 2-3 sentences describing what you plotted and what its implications are for your features.

Feel free to create as many plots as you want in your process of feature selection, but select one for the response cell below.

You should not just produce an identical visualization to question 3. Specifically, don't show us a bar chart of proportions, or a one-dimensional class-conditional density plot. Any other plot is acceptable, as long as it comes with thoughtful commentary. Here are some ideas:

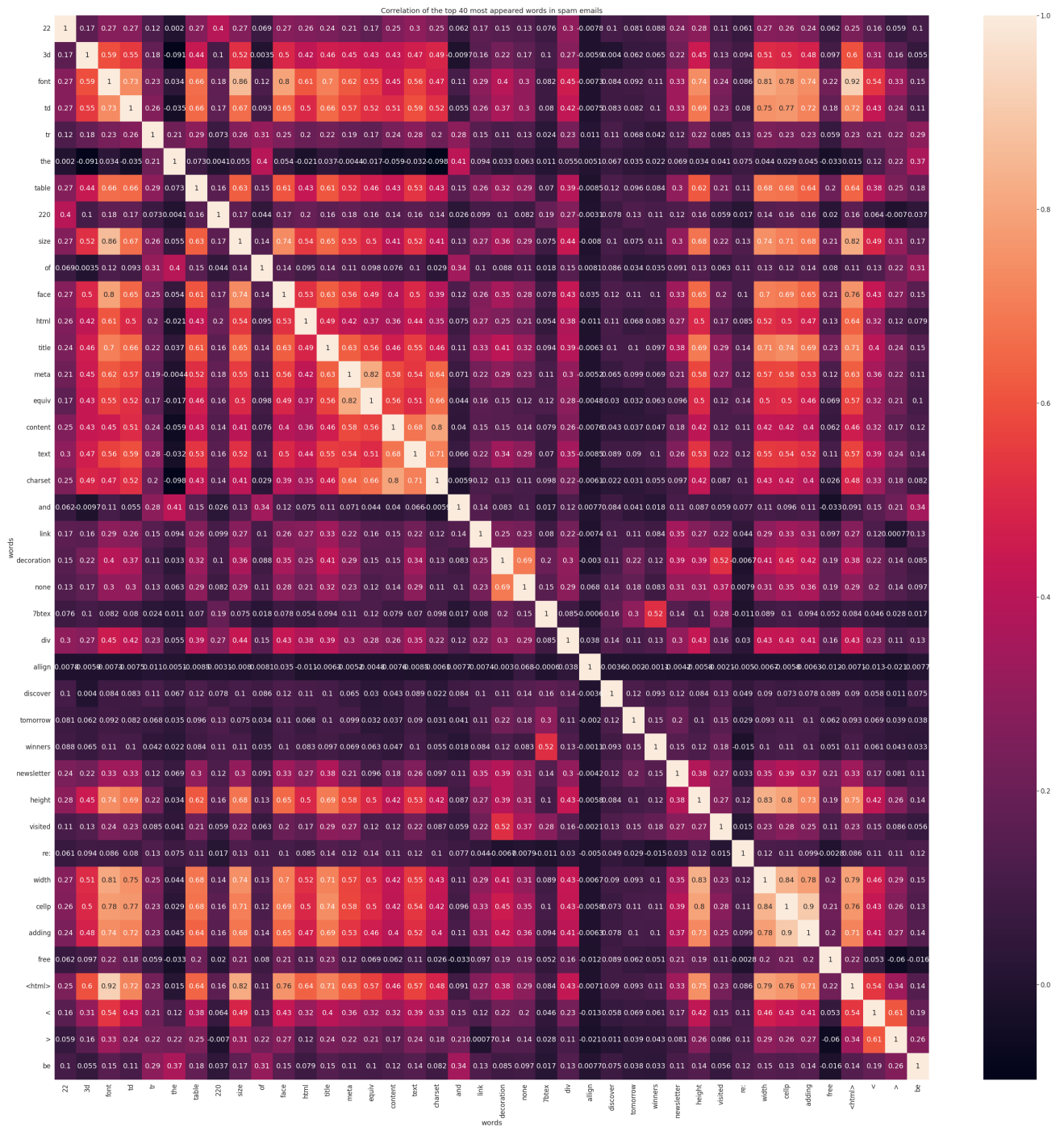
1. Consider the correlation between multiple features (look up correlation plots and `sns.heatmap`).
2. Try to show redundancy in a group of features (e.g. `body` and `html` might co-occur relatively frequently, or you might be able to design a feature that captures all html tags and compare it to these).
3. Visualize which words have high or low values for some useful statistic.
4. Visually depict whether spam emails tend to be wordier (in some sense) than ham emails.

Generate your visualization in the cell below and provide your description in a comment.

```
1 # Write your description (2-3 sentences) as a comment here:
2 # First, I found out the top 50 most appeared words in the spam emails(using value_counts or mod
3 # From that 50 words,
4 # I chose the top 40 (not including 'b') as the forty_most_words.
```

```
5 # With those 40 words, I wanted to create a correlation matrix.
6 # In other words, I wanted to find the relationship between those words.
7 # The results show various results, for example, the word 'the' seems to have
8 # a negative relationship with '3d', which means that the
9 # two have a negative linear correlation, but a very weak one.
10 # So one might not say something about the other.
11 # Another example is 'face' and 'size' has a 0.74,
12 # which indicates a strong positive linear correlation.
13 # This can mean that whenever 'face' appears, there is a high chance
14 # that 'size' appears in the text as well.
15 # So with these values, we can use a combination of
16 # words to predict whether the email is spam or ham.
17
18 # Write the code to generate your visualization here:
19 forty_most_words = ['22', '3d', 'font', 'td', 'tr', 'the', 'table', '220', 'size', 'of', 'face',
20                     'html', 'title', 'meta', 'equiv', 'content', 'text', 'charset', 'and',
21                     'link', 'decoration', 'none', '7btex', 'div', 'align', 'discover',
22                     'tomorrow', 'winners', 'newsletter', 'height', 'visited', 're:',
23                     'width', 'cellp', 'adding', 'free', '<html>', ' <', '>', 'be']
24 for word in forty_most_words:
25     train[word] = words_in_texts([word], train['email'])
26 plt.figure(figsize = (50,50))
27 sns.heatmap(train[forty_most_words].corr(method = 'pearson'), annot = True)
28
29 plt.xlabel("words")
30 plt.ylabel("words")
31 plt.title("Correlation of the top 40 most appeared words in spam emails")
32 plt.show()
33
34
35
36
```





▼ Question 8: Precision-Recall Curve

We can trade off between precision and recall. In most cases we won't be able to get both perfect precision (i.e. no false positives) and recall (i.e. no false negatives), so we have to compromise.

Recall that logistic regression calculates the probability that an example belongs to a certain class.

- Then, to classify an example we say that an email is spam if our classifier gives it ≥ 0.5 probability of being spam.
- However, *we can adjust that cutoff*: we can say that an email is spam only if our classifier gives it ≥ 0.7 probability of being spam.

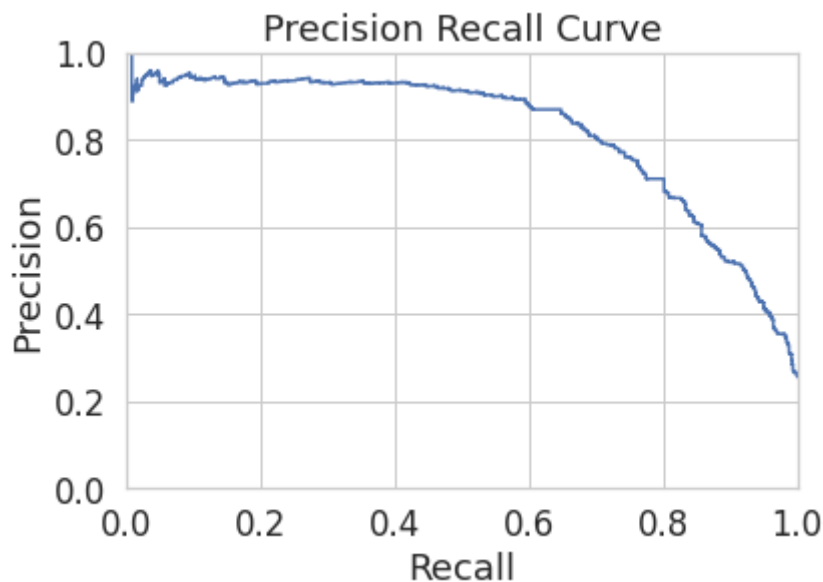
This is how we can trade off false positives and false negatives. The precision-recall curve shows this trade off for each possible cutoff probability. In the cell below, [plot a precision-recall curve](#) for your final classifier.

```
1 from sklearn.metrics import precision_recall_curve
2
3 # Note that you'll want to use the .predict_proba(...) method for your classifier
4 # instead of .predict(...) so you get probabilities, not classes
5
6
7
8
9
```

```

10
11 # BEGIN YOUR CODE
12 # -----
13 forty_most_words = ['22', '3d', 'font', 'td', 'tr', 'the', 'table', '220', 'size', 'of', 'face',
14                     'html', 'title', 'meta', 'equiv', 'content', 'text', 'charset', 'and',
15                     'link', 'decoration', 'none', '7btex', 'div', 'align', 'discover',
16                     'tomorrow', 'winners', 'newsletter', 'height', 'visited', 're:',
17                     'width', 'cellp', 'adding', 'free', '<html>', ' <', '>', 'be']
18
19 X_train = words_in_texts(forty_most_words, train['email'])
20 Y_train = np.array(train['spam'])
21 model = LogisticRegression().fit(X_train, Y_train)
22 y_predict_scores = model.predict_proba(X_train)[ :,1]
23 precision, recall, thresholds = precision_recall_curve(np.array(train['spam']), y_predict_scores)
24 plt.step(recall, precision)
25
26 plt.xlim([0, 1])
27 plt.ylim([0, 1])
28 plt.xlabel('Recall')
29 plt.ylabel('Precision')
30 plt.title('Precision Recall Curve');
31 plt.show()
32
33 model = LogisticRegression().fit(X_train, Y_train)
34 training_accuracy = model.score(X_train, Y_train)
35 Y_predict = model.predict(X_train)
36 print("accuracy:", training_accuracy)
37 # -----
38 # END YOUR CODE

```



accuracy: 0.8836683082656729

Explanation: As I was plotting a Precision-Recall Curve from my earlier classifier, I used the same set of forty_most_words. Since I am using my logistic regression classifier, I retrieved my earlier *X_train*, *Y_train*, and *model formula*, from the earlier questions.

A new array called `y_predict_scores` is created for the target scores in the Precision-Recall Curve parameter and is considered the probability estimates of the positive class.

Following the examples in https://scikit-learn.org/stable/modules/generated/sklearn.metrics.precision_recall_curve.html#sklearn.metrics.precision_recall_curve, I start the `precision_recall_curve`, by precision, recall, thresholds with the necessary parameters of the precision-recall curve. The first parameter is the True labels, which are the labels that are labeled/predicted correctly, hence it is the `Y_train` or `(np.array(train['spam']))` and the second parameter is the target or predicted scores, which was explained above in `y_predict_scores`.

After that, it is straightforward as the curve of precision and recall are both limited from 0 to 1. Next, I have to draw the curve by `plt.step`, which makes a step plot with the x-axis as recall and the y-axis as precision. Finally, to complete the Precision-Recall Curve there must be labels and a title.

Finally, I got the `model.score` and `training_accuracy` from Question 5 to grade my final accuracy percentage

Congratulations! You have completed HW 4.

Make sure you have run all cells in your notebook in order before running the cell below, so that all images/graphs appear in the output.,

Please save before submitting!

Please generate pdf as follows and submit it to Gradescope.

File > Print Preview > Print > Save as pdf