

INLS 560 – Project: SMS Spam Detection

Date Assigned: November 7, 2019

Completion Date: November 21, 2019

Software issues:

If you feel there are mistakes in this assignment, check Piazza for corrections, and report them to us if they have not been made.

Assignment

In this assignment, you will create a python program that automatically predicts whether an email is spam or not spam.

This application consists of several parts.

Parts

1. Dataset
2. Python program to compute features
3. Weka
4. Python program to predict whether a SMS message is spam or not spam
5. Your Presentation

Your presentation will be presented at regular class time on April 25.

Dataset

The dataset is on Piazza under a note named: Project Dataset.
Please download the dataset from Piazza.

The rows with ham are not spam.
The rows with spam are spam.

Create a file named predict_sms_spam.py.

This file should contain:

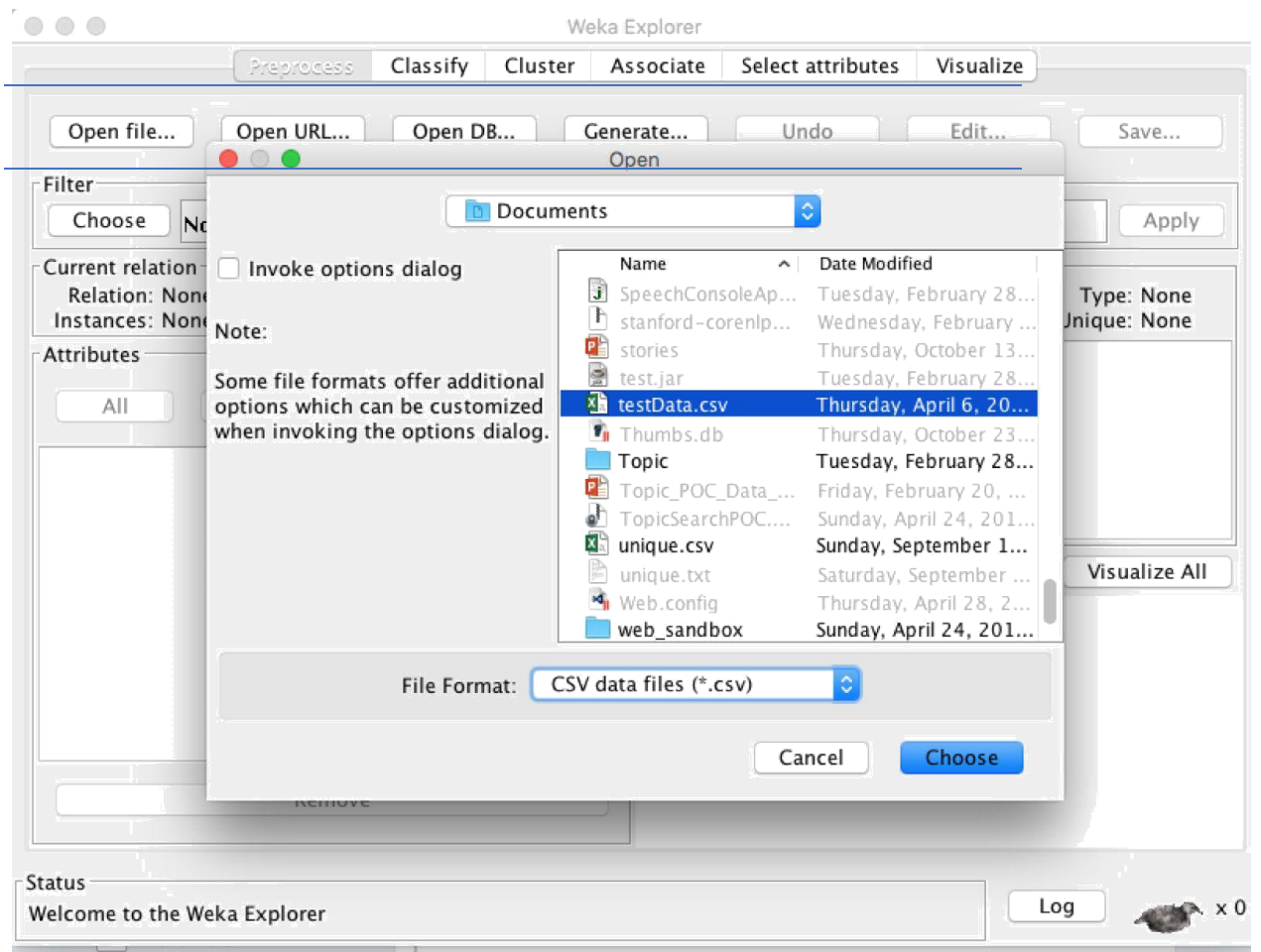
- A main function
- A function for each feature you will compute
- A function for reading in the dataset
- A function for writing the features to a csv file
 - o The format for this file should be:
 - doesHaveLinks, doesHaveSpammyWords, LengthOfText, NumberOfSymbols, class (spam or ham)
 - o Make sure that your file the headings above in the first row and the values of each column in subsequent rows

You should compute the following features:

- **doesHaveLinks:** This feature is True if a sms message has links and False if a sms message doesn't have links
- **doesHaveSpammyWords:** This feature is True if the sms contains spammy words and False if a sms message does not have spammy words
 - o To determine what words are spammy, look through the dataset and pay close attention to curse words and any words that are in the spam category but not in the ham category.
- **LengthOfText**
 - o The number of characters including spaces in the text message
- **NumberOfSymbols**
 - o The number of symbols in the text message

Create two additional features of your choice.

Select the csv file you created in the first python program. Make sure to select CSV data files under File Format.



After you press select, your screen should look similar to the one below (**your attributes will be different**):

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Generate... | Undo | Edit... | Save...

Filter: Choose **None** Apply

Current relation
Relation: testData-weka.filt... | Attributes: 3
Instances: 1000 | Sum of weights: 1000

Attributes
All | None | Invert | Pattern

No.	Name
1	<input checked="" type="checkbox"/> id
2	<input type="checkbox"/> value
3	<input type="checkbox"/> textLength

Remove

Selected attribute
Name: id | Type: Numeric
Missing: 0 (0%) | Distinct: 1000 | Unique: 1000 (100%)

Statistic	Value
Minimum	2795
Maximum	3794
Mean	3294.5
StdDev	288.819

Class: textLength (Num) Visualize All

100 100 100 100 100 100 100 100 100 100

2795 100.0 (3094.7, 3194.6) 3794

Status: OK Log x 0

Next select, the Classify tab and **make sure the field (Nom) value is selected**. Then press the choose button, select trees and choose J48. J48 is the decision tree classifier. Finally press Start.

The screenshot shows the Weka Explorer application window. The 'Classify' tab is active. The classifier is set to 'J48 -C 0.25 -M 2'. The test options are set to 'Cross-validation' with 10 folds. The '(Nom) value' dropdown is selected. The 'Start' button has been pressed, and the results are displayed in the 'Classifier output' pane.

Classifier output

Correctly Classified Instances	633	63.3	%
Incorrectly Classified Instances	367	36.7	%
Kappa statistic	0.2666		
Mean absolute error	0.4568		
Potential mean squared error	0.4787		
Relative absolute error	91.357	%	
Potential relative squared error	95.742	%	
Coverage of cases (0.95 level)	100	%	
Mean rel. region size (0.95 level)	100	%	
Total Number of Instances	1000		

Detailed Accuracy By Class

	TP Rate	FP Rate	Precision	Recall	F-Measure	MC
Weighted Avg.	0.429	0.162	0.726	0.429	0.540	0.
	0.838	0.571	0.594	0.838	0.695	0.

Confusion Matrix

```

a  b  <-- classified as
215 286 | a = spam
81  418 | b = not spam
  
```

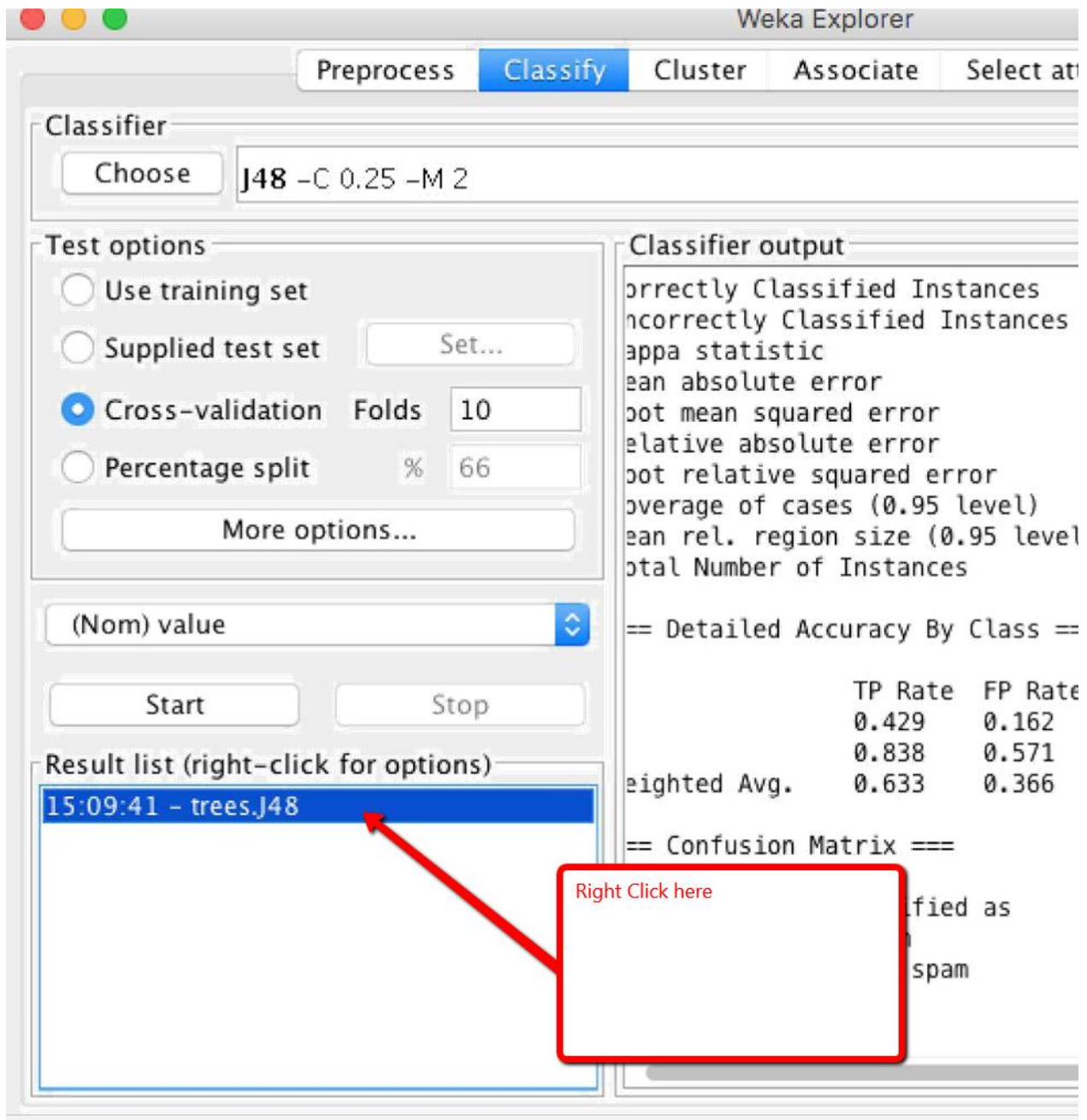
The 'Result list' on the left shows a single entry: '15:09:41 - trees.J48'.

Status: OK

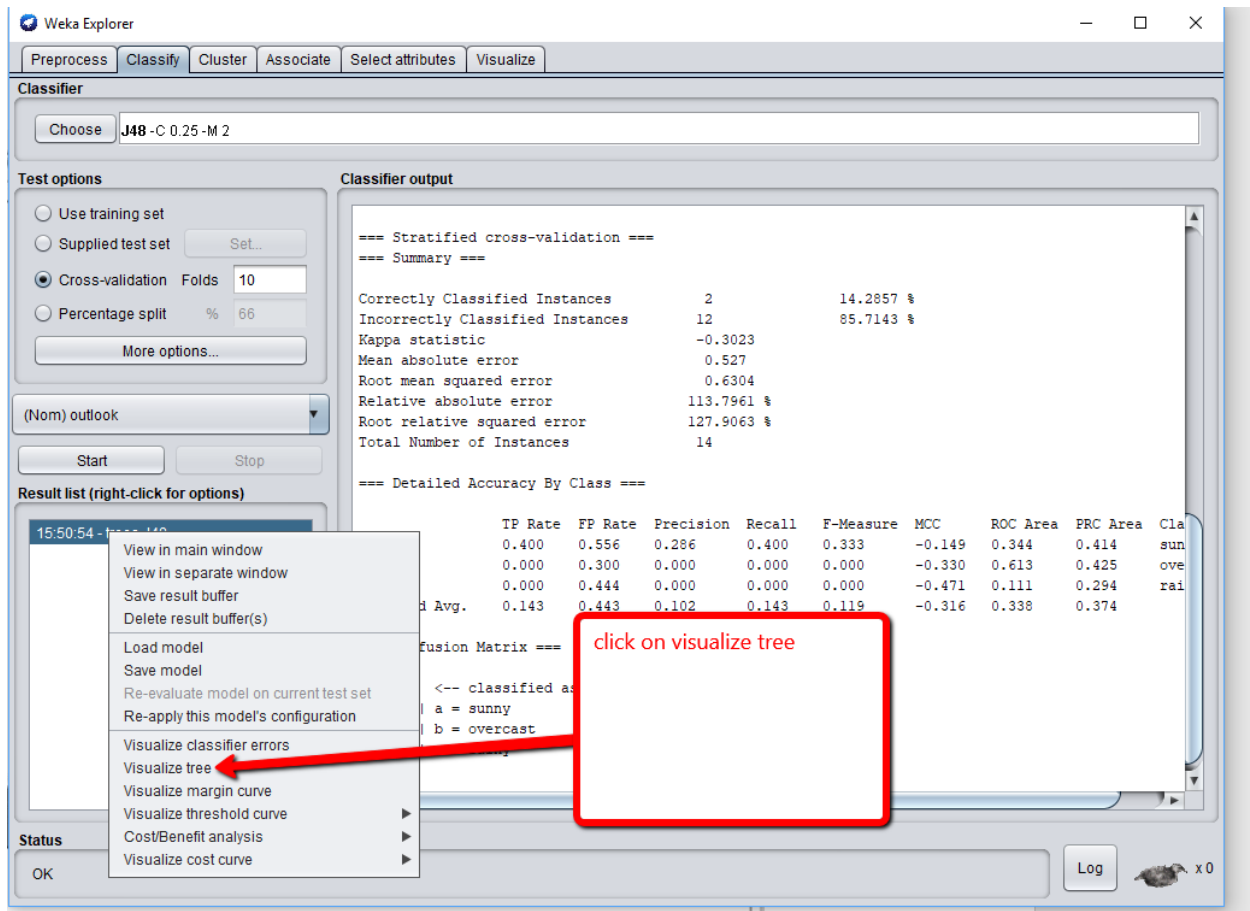
Log: x 0

Python program to predict whether a SMS message is spam or not spam (add to your previous program)

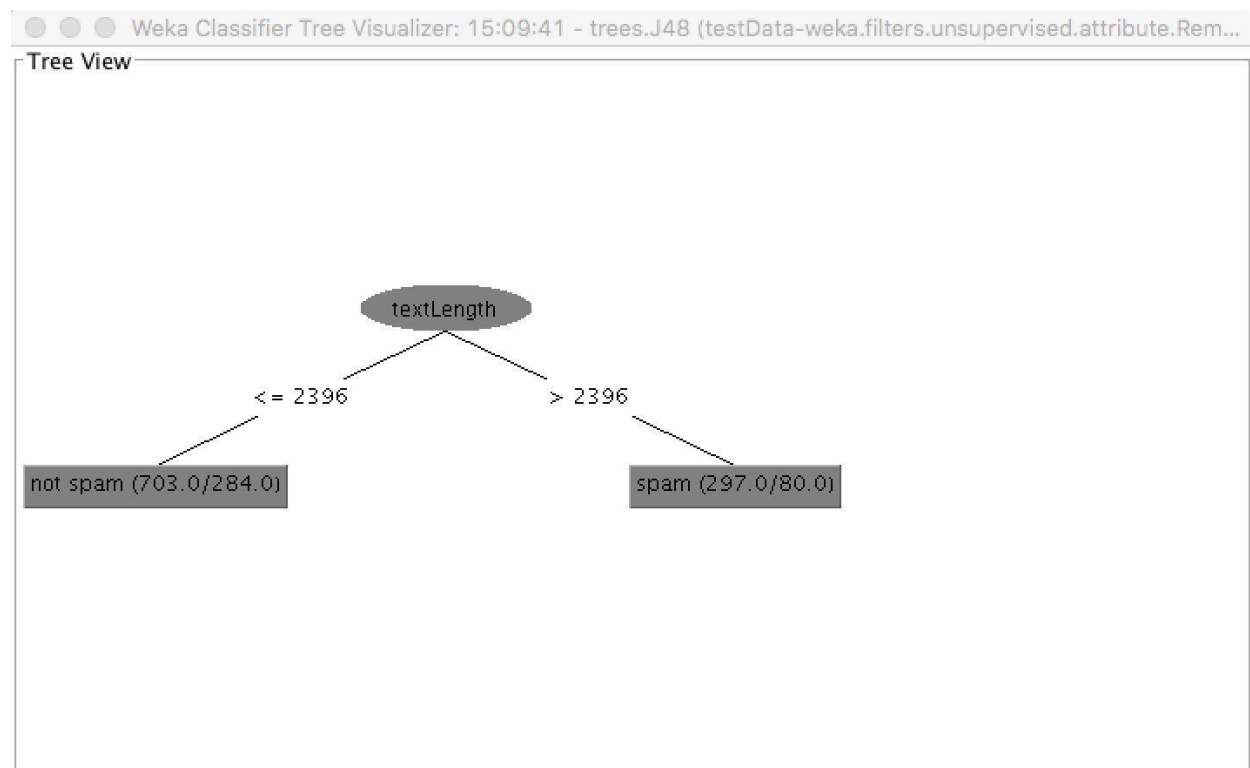
Follow the instructions in red in the image below.



Follow the instructions in red in the image below.



You should see a visualization similar to the one below.



Write down the if/else rules.

For example,

- If the text length is greater than 2396, then the message is spam
- Else if the text length is less than or equal to 2396 then the message is ham (not spam)

You will convert these rules into a function in Python.

Next, create a function named: `predict_spam` in your file `predict_sms_spam.py` **(DO NOT CREATE A NEW FILE)**

- This function should take as arguments, each feature
- This function should return “spam” or “ham”

Your Presentation

You will have a maximum of 10 minutes to present your program.

Discuss how you designed your program.

- a. Walk us through your program
- b. What functions you used and why
- c. What data structures did you use to compute your features and why did you choose those data structures?
 - i. Did you use lists, dictionaries, arrays?
 - ii. What made you choose these data structures?
- d. What features did you choose and why?
- e. Is there anything you wish you could have done better?
- f. What did you learn from this project?