

HomeWork_1_CSE517

1 Text Classification – Eisenstein 4.6 (p. 89)

After you run `tar -xzf A1.tgz`, in the directory `review polarity`, you will find a dataset of positively and negatively classified reviews that was used by Pang and Lee [2], a seminal paper about sentiment classification. Consult the `readme` file for more information. Hold out a randomly selected 400 reviews as a test set.

```
In [38]: 1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 from math import sqrt
```

In this step below the text files from the `pos` folder are being read and has been transformed into dataframe (`df_pos`). This dataframe contains two columns namely "Reviews" and "Sentiments (with values 'Positive')".

```
In [39]: 1 # Import Module
2 import os
3
4 # Folder Path
5 path = r'/Users/nehakardam/Documents/UWclasses /CSE NLP/A1/A1_data/revi
6
7 # Change the directory
8 os.chdir(path)
9
10 # Read text File
11 def read_text_file(file_path):
12     with open(file_path, 'r') as f:
13         return f.read().splitlines()
14
15 # iterate through all file
16 collection_of_data_pos = []
17 for file in os.listdir():
18     # Check whether file is in text format or not
19     if file.endswith(".txt"):
20         file_path = f"{path}/{file}"
21         # call read text file function
22         collection_of_data_pos = collection_of_data_pos + read_text_fil
```

```
In [40]: 1 df_pos = pd.DataFrame(collection_of_data_pos, columns = ['Reviews'])
2 df_pos['Sentiments'] = 'Positive'
```

In this step below the text files from the `neg` folder are being read and has been transformed into dataframe (`df_pos`). This dataframe contains two columns namely "Reviews" and "Sentiments (with values 'Negative')".

```
In [41]: 1 # Folder Path
2 path = r'/Users/nehakardam/Documents/UWclasses /CSE NLP/A1/A1_data/revi
3
4 # Change the directory
5 os.chdir(path)
6
7 # iterate through all file
8 collection_of_data_neg = []
9 for file in os.listdir():
10     # Check whether file is in text format or not
11     if file.endswith(".txt"):
12         file_path = f"{path}/{file}"
13         # call read text file function
14         collection_of_data_neg = collection_of_data_neg + read_text_fil
```

```
In [42]: 1 df_neg = pd.DataFrame(collection_of_data_neg, columns = ['Reviews'])
2 df_neg['Sentiments'] = 'Negative'
```

```
In [43]: 1 #both the dataframe of positive and negative reviews are concatenated t
2 df_pos_neg = [df_pos, df_neg]
3 Collection_of_data = pd.concat(df_pos_neg)
4 #Shuffle the columns of the data
5 All_data= Collection_of_data.sample(frac=1).reset_index(drop=True)
6 All_data
```

Out[43]:

	Reviews	Sentiments
0	whew . okay .	Negative
1	warning : this review contains some spoilers f...	Positive
2	its nice to see a pg rated movie out at christ...	Positive
3	he gives each one a different series of cranky...	Negative
4	but like everything else , tarantino takes thi...	Positive
...
64715	everything about this movie was just completel...	Negative
64716	basically , he says , it boils down to a lack ...	Negative
64717	y'know , that just sends chills down my spine .	Negative
64718	bryan singer's " the usual suspects " for one ...	Negative
64719	i knew it was going to be a bad print with onl...	Negative

64720 rows × 2 columns

```
In [66]: 1 All_data.to_csv('/Users/nehakardam/Documents/UWclasses /CSE NLP/A1/A1_d
2 lex_neg = pd.read_csv('/Users/nehakardam/Documents/UWclasses /CSE NLP/A
3 lex_pos = pd.read_csv('/Users/nehakardam/Documents/UWclasses /CSE NLP/A
```

```
In [67]: 1 dict_neg_pos = {}
2 for index, value in lex_neg[0].items():
3     dict_neg_pos[value] = "negative"
4 for index, value in lex_pos[0].items():
5     dict_neg_pos[value] = "positive"
```

```
In [216]: 1 import re
2 def tokenizer(theText):
3     theTokens = re.findall(r'\b\w[\w-]*\b', theText.lower())
4     return theTokens
5
6 reviews_tokenized = []
7 reviews = All_data['Reviews']
8
9 for index, value in reviews.items():
10     reviews_tokenized.append(tokenizer(value))
11 All_data['Reviews_Tokenized'] = reviews_tokenized
12
```

```
In [69]: 1 #Creating test data set
2 Testset_data = All_data.sample(n=400, random_state=1)
```

Part A: Sentiment lexicon-based classifier.

Sentiment lexicon-based classifier. Create a classifier using a sentiment lexicon. A lexicon from Hu and Liu [1] is provided in the directory opinion lexicon English, but you are welcome to find and use (with attribution, of course) another. Tokenize the data, and classify each document as positive if and only if it has more positive sentiment words than negative sentiment words. Compute and report the accuracy and F1 score (on detecting positive reviews) on the test set, using this lexicon-based classifier.

```
In [210]: 1 def result_of_sentiment_analysis(review):
2     pos_count = 0
3     neg_count = 0
4     for token in review:
5         if token in dict_neg_pos:
6             s = dict_neg_pos.get(token)
7             if s == "positive":
8                 pos_count = pos_count + 1
9             else:
10                neg_count = neg_count + 1
11     if pos_count >= neg_count:
12         return "Positive"
13     else:
14         return "Negative"
15
16 sentiments_y = []
17 reviews_tokens = Testset_data['Reviews_Tokenized']
18 for review in reviews_tokens:
19     sentiments_y.append(result_of_sentiment_analysis(review))
```

```
In [211]: 1 Testset_data['Sentiments_Y'] = sentiments_y
          2 Testset_data
```

Out[211]:

	Reviews	Sentiments	Reviews_Tokenized	Sentiments_Y
54756	you're so . . .	Negative	[you, re, so]	Positive
43964	no one would deny that spillane's writing has ...	Negative	[no, one, would, deny, that, spillane, s, writ...	Negative
21339	i can ? t recall seeing a film recently that w...	Positive	[i, can, t, recall, seeing, a, film, recently,...	Positive
56309	the times are changing , though , and the gear...	Positive	[the, times, are, changing, though, and, the, ...	Positive
54493	i certainly wish that was a new trend in holly...	Negative	[i, certainly, wish, that, was, a, new, trend,...	Positive
...
9945	details essential to the story are so improbab...	Negative	[details, essential, to, the, story, are, so, ...	Negative
62137	as many great movies that got cult status thro...	Positive	[as, many, great, movies, that, got, cult, sta...	Positive
10376	if we encourage terrorism abroad (from bening...	Positive	[if, we, encourage, terrorism, abroad, from, b...	Negative
50983	alas , the times are a-changin' .	Positive	[alas, the, times, are, a-changin]	Positive
6022	the chase is also shot and edited in that " je...	Positive	[the, chase, is, also, shot, and, edited, in, ...	Positive

400 rows × 4 columns

Accuracy and F1 Score calculation

```
In [72]: 1 actual = Testset_data['Sentiments']
          2 predicted = Testset_data['Sentiments_Y']
          3
          4 tp = np.sum((actual=='Positive') & (predicted=='Positive'))
          5 fp = np.sum((actual!='Positive') & (predicted=='Positive'))
          6 tn = np.sum((actual=='Negative') & (predicted=='Negative'))
          7 fn = np.sum((actual!='Negative') & (predicted=='Negative'))
          8
          9 accuracy = (tp+tn)/(tp+tn+fp+fn)
         10 precision = tp/(tp+fp)
         11 recall = tp/(tp+fn)
         12 f1 = 2 * (precision * recall) / (precision + recall)
```

```
In [73]: 1 print (accuracy, precision, recall, f1)
```

0.5375 0.4854014598540146 0.751412429378531 0.5898004434589801

Part B: Logistic regression classifier.

Train a (binary) logistic regression classifier on your training set using features of your own choosing, and report its accuracy and F1 score (as above) on the test set. In your write-up, describe the features you have chose and explain the reasoning behind your choice. Do not use pretrained word vectors or any features implemented or constructed by anyone else. Do not use an existing implementation of logistic regression, stochastic gradient descent, or automatic differentiation.

```
In [364]: 1 Training_data = All_data.copy()
          2 Testset_data = All_data.sample(n=4000, random_state=1)
          3 Training_data.drop(Testset_data.index)
          4 Training_data = Training_data.sample(n=20000, random_state=1)
```

Bag of word Model

```
In [365]: 1 word_dict = {}
          2
          3 def add_all_workds_to_dict(words):
          4     for word in words:
          5         if word in word_dict:
          6             word_dict[word] = word_dict.get(word) + 1
          7         else:
          8             word_dict[word] = 1
          9
          10 for words in Training_data["Reviews_Tokenized"]:
          11     add_all_workds_to_dict(words)
```

```
In [366]: 1 frequent_word_set = set()
          2 for key, value in word_dict.items():
          3     if value > 10:
          4         frequent_word_set.add(key)
```

```
In [367]: 1 len(frequent_word_set)
```

Out[367]: 3654

In [368]: 1 Training_data

Out[368]:

	Reviews	Sentiments	Reviews_Tokenized
54756	you're so . . .	Negative	[you, re, so]
43964	no one would deny that spillane's writing has ...	Negative	[no, one, would, deny, that, spillane, s, writ...
21339	i can ? t recall seeing a film recently that w...	Positive	[i, can, t, recall, seeing, a, film, recently,...
56309	the times are changing , though , and the gear...	Positive	[the, times, are, changing, though, and, the, ...
54493	i certainly wish that was a new trend in holly...	Negative	[i, certainly, wish, that, was, a, new, trend,...
...
44239	in pay it forward leder continues in this vein...	Negative	[in, pay, it, forward, leder, continues, in, t...
29593	the humor especially works well in this case ,...	Positive	[the, humor, especially, works, well, in, this...
22651	here's a film that actually gives away most of...	Positive	[here, s, a, film, that, actually, gives, away...
7397	a bad guy that's not really bad .	Positive	[a, bad, guy, that, s, not, really, bad]
10320	this vacuous picture throws in a standard down...	Negative	[this, vacuous, picture, throws, in, a, standa...

20000 rows × 3 columns

```
In [369]: 1 Training_data_with_features = pd.DataFrame(columns = frequent_word_set)
2 for review in Training_data["Reviews_Tokenized"]:
3     i = len(Training_data_with_features)
4     Training_data_with_features.loc[i] = 0
5     for word in review:
6         if word in frequent_word_set:
7             count = Training_data_with_features.loc[i, word]
8             Training_data_with_features.loc[i, word] = count + 1
9 train_X = Training_data_with_features.astype(float)
```

```
In [370]: 1 train_Y = []
2 for sentiment in Training_data["Sentiments"]:
3     if sentiment == "Positive":
4         train_Y.append(1)
5     else:
6         train_Y.append(0)
```

Logistic Regression Model Classifier

```

In [371]: 1 class Logistic_Regression() :
2         def __init__(self, learning_rate=0.001, n_iters=1000):
3             self.lr = learning_rate
4             self.n_iters = n_iters
5             self.weights = None
6             self.bias = None
7
8         def fit(self, X, y):
9             n_samples, n_features = X.shape
10            self.weights = np.zeros(n_features)
11            self.bias = 0
12
13            for i in range(self.n_iters):
14                model_linear = np.dot(X, self.weights) + self.bias
15                predict_y = self._sigmoid(model_linear)
16                dw = (1 / n_samples) * np.dot(X.T, (predict_y - y))
17                db = (1 / n_samples) * np.sum(predict_y - y)
18                self.weights -= self.lr * dw
19                self.bias -= self.lr * db
20
21            def predict(self, X):
22                model_linear = np.dot(X, self.weights) + self.bias
23                predict_y = self._sigmoid(model_linear)
24                predict_y_cls = [1 if i > 0.5 else 0 for i in predict_y]
25                return np.array(predict_y_cls)
26
27            def _sigmoid(self, x):
28                return 1 / (1 + np.exp(-x))

```

```

In [372]: 1 regression = Logistic_Regression()
2         regression.fit(train_X, train_Y)

```

```

In [373]: 1 Testset_data_with_features = pd.DataFrame(columns = frequent_word_set)
2         for review in Testset_data["Reviews_Tokenized"]:
3             i = len(Testset_data_with_features)
4             Testset_data_with_features.loc[i] = 0
5             for word in review:
6                 if word in frequent_word_set:
7                     count = Testset_data_with_features.loc[i, word]
8                     Testset_data_with_features.loc[i, word] = count + 1
9         test_X = Testset_data_with_features.astype(float)

```

```

In [374]: 1 test_Y = []
2         for sentiment in Testset_data["Sentiments"]:
3             if sentiment == "Positive":
4                 test_Y.append(1)
5             else:
6                 test_Y.append(0)

```

```

In [375]: 1 predictions = regression.predict(test_X)
2         type(predictions)

```

Out[375]: numpy.ndarray

```
In [376]: 1 actual = np.array(test_Y)
2 predicted = predictions
3
4 tp = np.sum((actual==1) & (predicted==1))
5 fp = np.sum((actual!=1) & (predicted==1))
6 tn = np.sum((actual==0) & (predicted==0))
7 fn = np.sum((actual!=0) & (predicted==0))
8
9 accuracy = (tp+tn)/(tp+tn+fp+fn)
10 precision = tp/(tp+fp)
11 recall = tp/(tp+fn)
12 f1 = 2 * (precision * recall) / (precision + recall)
```

```
In [377]: 1 print (accuracy, precision, recall, f1)

0.53325 0.5237562884292901 0.9199803632793323 0.6674977738201247
```

Breaking good.

Case 1: Positive by human English speakers, and your lexicon classifier predicts it as positive, whereas your logistic regression classifier predicts it as negative.

```
In [340]: 1 text1 = tokenizer("I've recently noticed something I had not even consi
2 result_of_sentiment_analysis(text1)
```

Out[340]: 'Positive'

```
In [341]: 1 breaking_bad_text = text1
2 breaking_bad_text_with_features = pd.DataFrame(columns = frequent_word_
3 breaking_bad_text_with_features.loc[0] = 0
4 for word in breaking_bad_text:
5     if word in frequent_word_set:
6         count = breaking_bad_text_with_features.loc[0, word]
7         breaking_bad_text_with_features.loc[0, word] = count + 1
8 breaking_bad_text_X = breaking_bad_text_with_features.astype(float)
9
10 regression.predict(breaking_bad_text_X) # 0 means negative and 1 means
```

Out[341]: array([0])

Case 2: Human positive and your lexicon classifier predicts it as negative, whereas your logistic regression classifier predicts it as positive.

```
In [380]: 1 text2 = tokenizer("You can bet that anybody who has ever seen Karan's m
2 result_of_sentiment_analysis(text2)
```

Out[380]: 'Negative'


```
In [381]: 1 breaking_bad_text = text2
2 breaking_bad_text_with_features = pd.DataFrame(columns = frequent_word_
3 breaking_bad_text_with_features.loc[0] = 0
4 for word in breaking_bad_text:
5     if word in frequent_word_set:
6         count = breaking_bad_text_with_features.loc[0, word]
7         breaking_bad_text_with_features.loc[0, word] = count + 1
8 breaking_bad_text_X = breaking_bad_text_with_features.astype(float)
9
10 regression.predict(breaking_bad_text_X) # 0 means negative and 1 means
```

Out[381]: array([1])

Case 3: Both of your classifiers predict it as negative.

```
In [409]: 1 text3 = tokenizer("They're all unhappy or marginally so about noodles i
2 result_of_sentiment_analysis(text3)
```

Out[409]: 'Negative'

```
In [410]: 1 breaking_bad_text = text3
2 breaking_bad_text_with_features = pd.DataFrame(columns = frequent_word_
3 breaking_bad_text_with_features.loc[0] = 0
4 for word in breaking_bad_text:
5     if word in frequent_word_set:
6         count = breaking_bad_text_with_features.loc[0, word]
7         breaking_bad_text_with_features.loc[0, word] = count + 1
8 breaking_bad_text_X = breaking_bad_text_with_features.astype(float)
9
10 regression.predict(breaking_bad_text_X) # 0 means negative and 1 means
```

Out[410]: array([0])

```
In [ ]: 1
```