

CSCI544: Homework Assignment No .1

1. Dataset Preparation

We will use the Amazon reviews dataset which contains real reviews for kitchen products sold on Amazon. The dataset is downloadable at: https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon_reviews_us_Kitchen_v1_00.tsv.gz

A: Downloaded the file from the link given above.

- a) Read the data as a Pandas frame using Pandas package and only keep the Reviews and Ratings _elds in the input data frame to generate data. Our goal is to train sentiment analysis classifiers that can predict the sentiment (positive/negative) for a given review.

Used pandas package to read the tsv file. In this dataset some of the rows were corrupted or had bad rows, therefore, we needed to mention a parameter called “error_bad_lines” to FALSE which skip all the bad rows.

Read Data

```
df = pd.read_csv("../amazon_reviews_us_Kitchen_v1_00.tsv", sep = "\t", error_bad_lines=False, warn_bad_lines=False)
```

As dataset is too large and we only need reviews and ratings, we filtered out our dataframe keeping only ratings and reviews. There might be some cells which have missing values so to get rid of that we dropped those rows so that our dataframe is correct and consistent.

Keep Reviews and Ratings

```
df = df[['star_rating', 'review_body']]
print("REVIEWS SAMPPL WITH RATING")
print(df.head(3))
print("STATISIICS OF THE REVIEWS: ")
print("rating 1: ", len(df[df.star_rating == 1]))
print("rating 2: ", len(df[df.star_rating == 2]))
print("rating 3: ", len(df[df.star_rating == 3]))
print("rating 4: ", len(df[df.star_rating == 4]))
print("rating 5: ", len(df[df.star_rating == 5]))
df = df.dropna()
```

REVIEWS SAMPPL WITH RATING

	star_rating	review_body
0	5.0	Beautiful. Looks great on counter.
1	5.0	I personally have 5 days sets and have also bo...
2	5.0	Fabulous and worth every penny. Used for clean...

STATISIICS OF THE REVIEWS:

```
rating 1: 426900
rating 2: 241948
rating 3: 349547
rating 4: 731733
rating 5: 3124759
```

To have binary class labels we added an extra column in our dataframe called “labels” The mapping of labels has been defined as follows:

If any reviews rating is more than 3 it will be classified as 1.

If any reviews rating is less than 3 it will be classified as 0.

And rest of the reviews which have rating equals 3 will be discarded.

The reviews with rating 4,5 are labelled to be 1 and 1,2 are labelled as 0. Discard the reviews with rating 3'

```
df['label'] = df.apply(lambda row: 1 if row.star_rating > 3 else 0, axis = 1)
r0 = len(df[df.label == 0])
r1 = len(df[df.label == 1])
r3 = len(df[df.star_rating == 3])
df = df[df.star_rating != 3]
print(r0, " , ", r1, " , ", r3)
```

1018348 , 3856296 , 349539

As the dataset is too large, we are going to take only 200000 rows which will have 50% positive samples and 50% negative samples. After having those samples, we are going to shuffle the dataset so that when we are splitting the model can be trained on both the samples.

We select 200000 reviews randomly with 100,000 positive and 100,000 negative reviews.

```
positive = df[df['label'] == 1]
negative = df[df['label'] == 0]
positive = positive.sample(n = 100000)
negative = negative.sample(n = 100000)
dataset = pd.concat([positive, negative])
dataset = dataset.reset_index(drop = True)
train=dataset.sample(frac=0.8,random_state=200)
test = dataset.drop(train.index)
train = train.reset_index(drop = True)
test = test.reset_index(drop = True)
```

We are splitting our dataset into training and testing dataset with 80 -20 split.

2. Data Cleaning

Implement the following steps to preprocess the dataset you created:

- convert the all reviews into the lower case.
- remove the HTML and URLs from the reviews
- remove non-alphabetical characters
- remove extra spaces
- perform contractions on the reviews, e.g., won't -> will not. Include as many contractions in English that you can think of.

You can either use Pandas package functions or any other built-in functions. Do not try to implement the above processes manually. Most of the above processes can be performed with one line of code

In your report, print the average length of the reviews in terms of character length in your dataset before and after cleaning (to be printed by .py file).

- a. Average length of reviews before data cleaning: 322.619475
- b. Average length of reviews after data cleaning: 301.630745
 - Str.lower() is used to convert all the characters to lowercase.

- To remove URL, we used regex which filters out every url embedded into the string. BeautifulSoup was used to eliminate HTML markups from the review body
- Regex was used with substitution to remove all the non-alphabetic characters.
- Same for the multiple spaces we used regex with substitution.
- Used the function contraction which substitutes any word which is in contraction dictionary to its root word.

3. PREPROCESSING

Use NLTK package to process your dataset:

- Remove the stopwords
- Perform lemmatization

Print three sample reviews before and after data cleaning + preprocessing. In the .py file, print the average length of the reviews in terms of character length in before and after processing.

A:

3 Sample Reviews before data cleaning and data preprocessing

- I put these on a south facing window and parts of the shades disintegrated over time due to the intense sun where I live - Colorado. Hunter Douglas would not cover them with their warranty as it states that exposure to the elements voids the warranty. They consider the sun an element...I guess you need to use these shades where the sun doesn't shine. These are way too expensive to be replacing them every several years. Just a heads up if you are considering buying these. I won't spend the money for these a gain. As a follow up- 8/29/07 - I wrote to the corporate office and Hunter Douglas is replacing the Sillouettes under warranty.
- The recipient loved her gift and that makes for a very satisfied customer. Thanks.
- Exactly what I expected when I ordered the product.

```
print("Sample Reviews before data cleaning \n")
for ele in review_sample:
    print(ele)
```

I put these on a south facing window and parts of the shades disintegrated over time due to the intense sun where I live - Colorado. Hunter Douglas would not cover them with their warranty as it states that exposure to the elements voids the warranty. They consider the sun an element...I guess you need to use these shades where the sun doesn't shine. These are way too expensive to be replacing them every several years. Just a heads up if you are considering buying these. I won't spend the money for these again. As a follow up- 8/29/07 - I wrote to the corporate office and Hunter Douglas is replacing the Sillouettes under warranty.

The recipient loved her gift and that makes for a very satisfied customer. Thanks.

Exactly what I expected when I ordered the product.

Sample Reviews after data cleaning and pre-processing:

3 Sample Reviews after data cleaning and data preprocessing

- put south facing window part shade disintegrated time due intense sun live colorado hunter douglas would cover warranty state exposure element void warranty consider sun shine way expensive replacing every several year head considering buying spend money follow wrote corporate office hunter douglas replacing sillouettes warranty
- recipient loved gift make satisfied customer thanks
- exactly expected ordered product

```
print("Sample Reviews after data cleaning and pre-processing:\n")
for ele in review_sample_2:
    print(ele)
```

put south facing window part shade disintegrated time due intense sun live colorado hunter douglas would cover warranty state e
xposure element void warranty consider sun shine way expensive replacing every several year head considering buying spend money
follow wrote corporate office hunter douglas replacing sillouettes warranty
recipient loved gift make satisfied customer thanks
exactly expected ordered product

```
print("Average length of reviews before data cleaning: ", aclbdc)
print("Average length of reviews after data cleaning: ", acladc)

print("Average length of reviews before data pre-processing: ", aclbdc)
print("Average length of reviews after data pre-processing: ", acladc)

Average length of reviews before data cleaning: 322.619475
Average length of reviews after data cleaning: 301.630745
Average length of reviews before data pre-processing: 301.630745
Average length of reviews after data pre-processing: 184.227195
```

4. Feature Extraction

Use sklearn to extract TF-IDF features. At this point, you should have created a dataset which consists of features and binary labels for the reviews you selected.

A:

TF-IDF Vectorizer library was used to extract features from the reviews. Fit and transform was applied on the training dataset and only transform on testing dataset.

TF-IDF Feature Extraction

```
tfidf_vect = TfidfVectorizer(min_df = 0.001)
X_train = tfidf_vect.fit_transform(train['review_body'])
X_train = pd.DataFrame(X_train.toarray(), columns = tfidf_vect.get_feature_names())
X_test = tfidf_vect.transform(test['review_body'])
X_test = pd.DataFrame(X_test.toarray(), columns = tfidf_vect.get_feature_names())
Y_train = train['label']
Y_test = test['label']
```

5. Perceptron

Train a Perceptron model on your training dataset using the sklearn built-in implementation. Report Accuracy, Precision, Recall, and f1-score on both the training and testing split of your dataset. These 8 values should be printed in separate lines by the .py file.

A:

Perceptron was imported from the sklearn library with max iteration/epochs to 80, stopping criterion to 0.00001, random_state to 200 for consistency and learning rate at 0.01.

Accuracy, Precision, Recall and F1-score

Train: 0.87, 0.84, 0.91, 0.87

Test: 0.86, 0.84, 0.90, 0.87

Perceptron ¶

```
def metrics(true, pred):
    tn, fp, fn, tp = cm(true, pred).ravel()
    acc = (tp + tn)/(tn + fp + fn + tp)
    prec = tp/(tp + fp)
    rec = tp / (tp + fn)
    f1 = 2*(rec * prec) / (rec + prec)
    return [acc, prec, rec, f1]

def print_seq(score_list):
    print("%.2f" % score_list[0], "\n%.2f" % score_list[1], "\n%.2f" % score_list[2], "\n%.2f" % score_list[3])

percept = Perceptron(max_iter = 80, tol = 1e-5, random_state = 200, eta0 = 0.01)
percept.fit(X_train, Y_train)
Y_train_pred = percept.predict(X_train)
train_score = metrics(Y_train, Y_train_pred)

# Predicting on test dataset
Y_test_pred = percept.predict(X_test)
test_score = metrics(Y_test, Y_test_pred)
# print(train_score)
# print(test_score)
# print("TRAIN_SCORES:")
print_seq(train_score)
# print("TEST_SCORES:")
print_seq(test_score)
```

0.87
0.84
0.91
0.88
0.86
0.84
0.90
0.87

6. SVM

Train an SVM model on your training dataset using the sklearn built-in implementation. Report Accuracy, Precision, Recall, and f1-score on both the training and testing split of your dataset. These 8 values should be printed in separate lines by the .py file.

A:

Linear SVM was used with max iteration at 5000. Fit method was used to train the SVM model on the training dataset and predict on testing dataset.

Accuracy, Precision, Recall and F1-score:

Train: 0.89, 0.90, 0.89, 0.89

Test: 0.89, 0.89, 0.88, 0.89

SVM

```
: svc = SVC(max_iter = 5000)
model_svc = svc.fit(X_train, Y_train)
pred_train_svc = model_svc.predict(X_train)
pred_test_svc = model_svc.predict(X_test)
train_svc_score = metrics(Y_train, pred_train_svc)
test_svc_score = metrics(Y_test, pred_test_svc)

# print(train_svc_score)
# print(test_svc_score)
# print("TRAIN_SCORES:")
print_seq(train_svc_score)
# print("TEST_SCORES:")
print_seq(test_svc_score)
```

0.90
0.90
0.89
0.90
0.89
0.90
0.88
0.89

7. Logistic Regression

Train a Logistic Regression model on your training dataset using the sklearn built-in implementation. Report Accuracy, Precision, Recall, and f1-score on both the training and testing split of your dataset. These 8 values should be printed in separate lines by the .py file.

A:

In Logistic Regression we used max iteration to 1000, penalty to l2 and solver to liblinear.

Accuracy, Recall, Precision and F1-score:

Train: 0.89, 0.90, 0.89, 0.89

Test: 0.89, 0.89, 0.88, 0.89

Logistic Regression

```
log_reg = LR(solver='liblinear', random_state=0, C=5, penalty='l2', max_iter=1000)
model = log_reg.fit(X_train, Y_train)
pred_train_LR = model.predict(X_train)
pred_test_LR = model.predict(X_test)
train_LR_score = metrics(Y_train, pred_train_LR)
test_LR_score = metrics(Y_test, pred_test_LR)

# print(train_LR_score)
# print(test_LR_score)
# print("TRAIN_SCORES:")
print_seq(train_LR_score)
# print("TEST_SCORES:")
print_seq(test_LR_score)

0.90
0.90
0.89
0.90
0.89
0.90
0.88
0.89
```

8. Naïve Bayes

Train a Multinomial Naive Bayes model on your training dataset using the sklearn built-in implementation. Report Accuracy, Precision, Recall, and f1-score on both the training and testing split of your dataset. These 8 values should be printed in separate lines by the .py file.

A:

Accuracy, Precision, Recall and F1-score

Train: 0.86, 0.86, 0.86, 0.86

Test: 0.86, 0.86, 0.86, 0.86

Naive Bayes

```
MNB = MultinomialNB()
MNB.fit(X_train, Y_train)
pred_train_MNB = MNB.predict(X_train)
pred_test_MNB = MNB.predict(X_test)
train_MNB_score = metrics(Y_train, pred_train_MNB)
test_MNB_score = metrics(Y_test, pred_test_MNB)

# print(train_MNB_score)
# print(test_MNB_score)
# print("TRAIN_SCORES:")
print_seq(train_MNB_score)
# print("TEST_SCORES:")
print_seq(test_MNB_score)
```

```
0.87
0.86
0.87
0.87
0.86
0.86
0.86
0.86
```