

SENTIMENT ANALYSIS USING LOGISTIC REGRESSION

Objective

The objective of this project is to perform sentiment analysis on a dataset of sentences from online reviews using logistic regression. Specifically, the project aims to:

1. **Preprocess** the text data by cleaning and transforming it into a suitable format for analysis.
2. **Train** a logistic regression model to classify sentences as positive or negative.
3. **Evaluate** the performance of the model and analyze its results.
4. **Identify** influential words that have the largest impact on the sentiment classification.

Summary

The dataset consists of 3,000 sentences from reviews on "imdb.com", "amazon.com", and "yelp.com", labeled as either positive or negative. The project follows these steps:

1. **Data Preprocessing:**
 - Cleaned the text by removing punctuation, digits, and converting sentences to lowercase.
 - Removed stop words and transformed the sentences into a bag-of-words representation.
2. **Model Training:**
 - Used logistic regression to train a classifier on the processed text data.
 - Evaluated the model's performance on training and test datasets.
3. **Margin Analysis:**
 - Analyzed the relationship between the margin (confidence) of predictions and their accuracy.
4. **Influential Words:**
 - Identified words with the highest positive and negative coefficients to understand their impact on sentiment classification.

Results

1. **Data Preparation:**
 - Preprocessed text data by removing unnecessary characters and stop words.
 - Converted text to a numerical format using the bag-of-words approach with a vocabulary size capped at 4,500 words.
 - Split the data into training (2,500 sentences) and test (500 sentences) sets.

2. Model Training and Evaluation:

- Trained a logistic regression model using stochastic gradient descent.
- Training Error: Computed as the fraction of misclassified training sentences.
- Test Error: Computed as the fraction of misclassified test sentences.
- Training Error: Approximately X% (to be filled in based on actual results).
- Test Error: Approximately Y% (to be filled in based on actual results).

3. Margin Analysis:

- Visualized the distribution of margins for the test set.
- Analyzed the relationship between margin and error rate, revealing how confidence impacts classification accuracy.

4. Influential Words:

- Identified words with the highest positive and negative coefficients.
- Highly Positive Words: Words that significantly influence positive sentiment.
- Highly Negative Words: Words that significantly influence negative sentiment.
- Example: [List of influential words] (to be filled in based on actual results).

Conclusion

The logistic regression model effectively classified sentences into positive or negative sentiment, with performance evaluated on both training and test datasets. The margin analysis demonstrated that higher confidence predictions are generally more accurate. The identification of influential words provides insight into what drives sentiment classification, offering potential for further refinement and understanding of the model.

Code :

```
'''
```

❑ The "sentiment" dataset consists of 3000 sentences which come from reviews on "imdb.com", "amazon.com", and "yelp.com". Each sentence is labeled according to whether it comes from a positive review or negative review.

❑ The data set consists of 3000 sentences, each labeled '1' (if it came from a positive review) or '0' (if it came from a negative review). To be consistent with our notation from lecture, we will change the negative review label to '-1'.

❑ We will use logistic regression to learn a classifier from this data.

```
'''
```

#1. Load and Preprocess data

```
import string
import numpy as np
import matplotlib
import matplotlib.pyplot as plt
matplotlib.rc('xtick', labelsize=14)
matplotlib.rc('ytick', labelsize=14)
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.linear_model import SGDClassifier
import warnings
warnings.filterwarnings("ignore")

#Read in the data set:

with open("D:\\Aditya's Notes\\Aditya's Data Science Notes\\Projects and Other Datasets\\ML
PROJECTS\\data\\full_set.txt") as f:
    content = f.readlines()

#Remove leading and trailing white space
```

```

content = [x.strip() for x in content]
#Separate the sentences from the labels:
sentences = [x.split("\t")[0] for x in content]
labels = [x.split("\t")[1] for x in content]
#Transform the labels from '0 versus 1' to '-1 versus 1':
y = np.array(labels, dtype='int8')
y = 2*y - 1

```

'''

Preprocessing the text data

☐ To transform this prediction problem into an linear classification, we will need to preprocess the text data. We will do four transformations:

- o Remove punctuation and numbers.
- o Transform all words to lower-case.
- o Remove stop words.
- o Convert the sentences into vectors, using a bag-of-words representation.

☐ We begin with first two steps

'''

#"full_remove" takes a string x and a list of characters

#"removal_list" and returns with all the characters in removal_list replaced by ' '.

def full_remove(x, removal_list):

 for w in removal_list:

 x = x.replace(w, ' ')

 return x

digits = [str(x) for x in range(10)] # Remove digits

digit_less = [full_remove(x, digits) for x in sentences]

punc_less = [full_remove(x, list(string.punctuation)) for x in digit_less] # Remove punctuation

sents_lower = [x.lower() for x in punc_less] # Make everything lowercase

#Stop words - Stop words are the words that are filtered out because

```
# they are believed to contain no useful information for the task at hand.
```

```
# These usually include articles such as 'a' and 'the', pronouns such
```

```
# as 'i' and 'they', and prepositions such 'to' and 'from'.
```

```
# We have put together a very small list of stop words,
```

```
# but these are by no means comprehensive.
```

```
# Define our stop words
```

```
stop_set = set(['the', 'a', 'an', 'i', 'he', 'she', 'they', 'to', 'of', 'it', 'from'])
```

```
# Remove stop words
```

```
sents_split = [x.split() for x in sents_lower]
```

```
sents_processed = [" ".join(list(filter(lambda a: a not in stop_set, x))) for x in sents_split]
```

```
#Let us look at the sentences:
```

```
sents_processed[0:10]
```

```
'''
```

Bag of words

☐ In order to use linear classifiers on our data set, we need to transform our textual data into numeric

data. The classical way to do this is known as the bag of words representation.

☐ In this representation, each word is thought of as corresponding to a number in "{1, 2, ..., V}" where

"V" is the size of our vocabulary. And each sentence is represented as a V-dimensional vector x , where x_i is the number of times that word occurs in the sentence.

☐ To do this transformation, we will make use of the "CountVectorizer" class in "scikit-learn". We will cap the number of features at 4500, meaning a word will make it into our vocabulary only if it is one of the 4500 most common words in the corpus. This is often a useful step as it can weed out spelling mistakes and words which occur too infrequently to be useful.

☐ Finally, we will also append a '1' to the end of each vector to allow our linear classifier to learn a bias

term.

```
'''
```

```

# Transform to bag of words representation.

vectorizer = CountVectorizer(analyzer = "word", tokenizer = None,
preprocessor = None, stop_words = None, max_features = 4500)
data_features = vectorizer.fit_transform(sents_processed)

# Append '1' to the end of each vector.
data_mat = data_features.toarray()

#Training / test split - We split the data into a training set of 2500 sentences
# and a test set of 500 sentences (of which 250 are positive and 250 negative).
# Split the data into testing and training sets
np.random.seed(0)

test_inds = np.append(np.random.choice((np.where(y==-1))[0], 250,
replace=False), np.random.choice((np.where(y==1))[0], 250,
replace=False))

train_inds = list(set(range(len(labels))) - set(test_inds))

train_data = data_mat[train_inds,]
train_labels = y[train_inds]
test_data = data_mat[test_inds,]
test_labels = y[test_inds]

print("train data: ", train_data.shape)
print("test data: ", test_data.shape)

```

'''

2. Fitting a logistic regression model to the training data

❑ We could implement our own logistic regression solver using stochastic gradient descent, but fortunately, there is already one built into "scikit-learn".

❑ Due to the randomness in the SGD procedure, different runs can yield slightly different solutions (and thus different error values).

'''

```

# Fit logistic classifier on training data

```

```

clf = SGDClassifier(loss="log_loss", penalty=None) # Corrected the loss and penalty parameters

```

```

clf.fit(train_data, train_labels)

# Pull out the parameters (w,b) of the logistic regression model
w = clf.coef_[0,:]
b = clf.intercept_

# Get predictions on training and test data
preds_train = clf.predict(train_data)
preds_test = clf.predict(test_data)

# Compute errors
errs_train = np.sum((preds_train > 0.0) != (train_labels > 0.0))
errs_test = np.sum((preds_test > 0.0) != (test_labels > 0.0))
print ("Training error: ", float(errs_train)/len(train_labels))
print ("Test error: ", float(errs_test)/len(test_labels))

```

'''

3. Analyzing the margin

□ The logistic regression model produces not just classifications but also conditional probability estimates.

□ We will say that "x" has margin "gamma" if (according to the logistic regression model)

$\Pr(y=1|x) > (1/2)+\text{gamma}$ or $\Pr(y=1|x) < (1/2)-\text{gamma}$ ". The following function

margin_counts takes as input as the classifier ("clf", computed earlier), the test set

("test_data"), and a value of "gamma", and computes how many points in the test set have margin of at least "gamma".

'''

Return number of test points for which $\Pr(y=1)$ lies in $[0, 0.5 - \text{gamma})$ or $(0.5 + \text{gamma}, 1]$

```
def margin_counts(clf, test_data, gamma):
```

```
    # Compute probability on each test point
```

```
    preds = clf.predict_proba(test_data)[:,-1]
```

```
    # Find data points for which prediction is at least gamma away from 0.5
```

```
    margin_inds = np.where((preds > (0.5+gamma)) | (preds < (0.5- gamma)))[0]
```

```
    return float(len(margin_inds))
```

#Let us visualize the test set's distribution of margin values.

```
gammas = np.arange(0, 0.5, 0.01)
```

```
f = np.vectorize(lambda g: margin_counts(clf, test_data, g))
```

```
plt.plot(gammas, f(gammas) / 500.0, linewidth=2, color='green')
```

```
plt.xlabel('Margin', fontsize=14)
```

```
plt.ylabel('Fraction of points above margin', fontsize=14)
```

```
plt.show()
```

#We investigate a natural question: "Are points "x" with larger margin more likely to

be classified correctly? To address this, we define a function margin_errors

that computes the fraction of points with margin at least "gamma" that are misclassified.

Return error of predictions that lie in intervals $[0, 0.5 - \gamma)$ and $(0.5 + \gamma, 1]$

```
def margin_errors(clf, test_data, test_labels, gamma):
```

```
    # Compute probability on each test point
```

```
    preds = clf.predict_proba(test_data)[:,-1]
```

```
    # Find data points for which prediction is at least gamma away from 0.5
```

```
    margin_inds = np.where((preds > (0.5+gamma)) | (preds < (0.5-gamma)))[0]
```

```
    # Compute error on those data points.
```

```
    num_errors = np.sum((preds[margin_inds] > 0.5) !=
```

```
    (test_labels[margin_inds] > 0.0))
```

```
    return float(num_errors)/len(margin_inds)
```

#Let us visualize the relationship between margin and error rate.

Create grid of gamma values

```
gammas = np.arange(0, 0.5, 0.01)
```

Compute margin_errors on test data for each value of g

```
f = np.vectorize(lambda g: margin_errors(clf, test_data, test_labels,  
g))
```

Plot the result

```
plt.plot(gammas, f(gammas), linewidth=2)
```



```
plt.ylabel('Error rate', fontsize=14)
plt.xlabel('Margin', fontsize=14)
plt.show()
```

'''

4. Words with large influence

¶ Finally, we attempt to partially interpret the logistic regression model.

¶ Which words are most important in deciding whether a sentence is positive? As a first approximation

to this, we simply take the words whose coefficients in "w" have the largest positive values.

¶ Likewise, we look at the words whose coefficients in "w" have the most negative values, and we think of these as influential in negative predictions.

'''

Convert vocabulary into a list:

```
vocab = np.array([z[0] for z in sorted(vectorizer.vocabulary_.items(),
key=lambda x:x[1])])
```

Get indices of sorting w

```
inds = np.argsort(w)
```

Words with large negative values

```
neg_inds = inds[0:50]
```

```
print("Highly negative words: ")
```

```
print([str(x) for x in list(vocab[neg_inds])])
```

Words with large positive values

```
pos_inds = inds[-49:-1]
```

```
print("Highly positive words: ")
```

```
print([str(x) for x in list(vocab[pos_inds])])
```