

# (SentiMix Hindi-English)

M. Shehroz Wali Khan  
FAST School of Computing  
NUCES, Islamabad  
Islamabad, Pakistan  
i170308@nu.edu.pk

**Abstract**—Mixing languages, also known as code-mixing, is a norm in multilingual societies. Multilingual people, who are non-native English speakers, tend to code-mix using English-based phonetic typing and the insertion of anglicisms in their main language. In addition to mixing languages at the sentence level, it is fairly common to find the code-mixing behavior at the word level. This linguistic phenomenon poses a great challenge to conventional NLP systems, which currently rely on monolingual resources to handle the combination of multiple languages. The objective of this proposal is to bring the attention of the research community towards the task of sentiment analysis in code-mixed social media text. Specifically, we focus on the combination of English with Spanish (Spanglish) and Hindi (Hinglish), which are the 3rd and 4th most spoken languages in the world respectively.

## Introduction

The task was to predict the sentiment of a given code-mixed tweet. The sentiment labels are positive, negative, or neutral, and the code-mixed languages will be English-Hindi and English-Spanish. Besides the sentiment labels, we will also provide the language labels at the word level. The word-level language tags are en (English), spa (Spanish), hi (Hindi), mixed, and univ (e.g., symbols, @ mentions, hashtags). Efficiency will be measured in terms of Precision, Recall, and F-measure.

## Propose Approach

Here are the following steps were involved in doing the project:-

### 1. Importing required libraries

I imported all the required libraries and some extras in case I'd need them in the future.

```
1 import numpy as np
2 import pandas as pd
3 import re
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from scipy.stats import chi2_contingency
7
8 from sklearn.model_selection import train_test_split
9 from sklearn.model_selection import cross_val_score
10 from sklearn.model_selection import GridSearchCV
11 from sklearn.ensemble import RandomForestClassifier
12 from sklearn.metrics import accuracy_score
13 from sklearn.svm import SVC
14 from sklearn.feature_extraction.text import CountVectorizer
15
16 from subprocess import check_output
```

### 2. Mounting Google Drive to Colab

I uploaded all the data sets in the google drive. Before loading the datasets into the colab, the first thing that is needed is to mount the google drive in the colab. When you write the code written below and run it, colab provides a url and adds a field to write an authorization code. You have to put the authorization code which you will be getting by visiting the url and allowing some permission through your google account. In the end it provides you a code and you just have to paste it into this field and press enter. After that this text will appear and this will confirm that the drive is now mounted.

```
1 from google.colab import drive
2 drive.mount('/content/drive')
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
```

### 3. Defining Reading Functions

I have defined two functions. First for train sets and second for test sets.

**readandSplit():**

The main purpose of this function is to get open and read the datasets and use some splitting functions. Read each tweet and append it into the tweets list and read label i.e. sentiments (positive, negative, neutral) and append it into label list and return both lists.

**readInList():**

The main purpose of this function is to get open and read the test datasets and use some splitting functions. Read each tweet and append it into the tweets list and return this list.

For Train Datasets

```
1 def readandSplit(filename):
2     with open(filename) as f:
3         data = f.read().split('\n')
4     tweets = []
5     label = []
6     tweet = ""
7
8
9     for d in data:
10         if(len(d)==0):
11             continue
12         elif(len(d.split())==3):
13             label.append(d.split()[2])
14             if(tweet):
15                 tweets.append(tweet)
16                 tweet = ""
17             elif(len(d.split())==2):
18                 tweet += d.split()[0]+' '
19
20     tweets.append(tweet)
21
22     return label,tweets
```

For Test Datasets

```
1 def readInList(filename):
2     with open(filename) as f:
3         data = f.read().split('\n')
4     tweets = []
5     label = []
6     tweet = ""
7
8
9     for d in data:
10         if(len(d)==0):
11             continue
12         elif(d.split()[0]=='meta'):
13             if(tweet):
14                 tweets.append(tweet)
15                 tweet = ""
16             elif(len(d.split())==2):
17                 tweet += d.split()[0]+' '
18
19
20     tweets.append(tweet)
21
22     return tweets
```

## 4. Loading Datasets & Data Preprocessing

### Loading Datasets

We will be passing the train and test datasets into the readandSplit functions and assign it into lists of label and trainData (for Train Datasets)

And assign it into a list of vallebl, valData (Validation Datasets).

And to get the maximum training datasets, we are concatenating trainData with valData and label with vallebl.

```
1 label,trainData= readandSplit("/content/drive/MyDrive/Dataset/Train/train_14k_split_conll.txt")
2 vallebl,valData= readandSplit("/content/drive/MyDrive/Dataset/Validation/dev_3k_split_conll.txt")
3
4 trainData+=valData
5 label+=vallebl
```

### Vectorize the TrainData

I used a method from CountVectorizer from sklearn.feature\_extraction.text

I.e fit transform() on trainData which learn the vocabulary dictionary and return a document-term matrix and then convert it into an array using toarray() and assign it into variable X.

```
1 vectorizer = CountVectorizer(min_df=1)
2 X = vectorizer.fit_transform(trainData).toarray()
```

## 5. Model

**Note:** I also used LSTM Algorithm for it, but the accuracy was low around 58 and it took so much time to train so I just stopped using it and shift it towards Random Forest Classifiers

### Final Model Implemented: **Random Forest Classifier**

- Putting Data into the Random Forest Classifier
- Fitting the Training Data
- Predicting the Test Data

#### Brief Details:-

I called the RandomForestClassifier function and then used fit() to fit the train Data. It took around 5-10 minutes in training.

```
1 rfc1=RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
2                             criterion='gini', max_depth=None, max_features='auto',
3                             max_leaf_nodes=None, max_samples=None,
4                             min_impurity_decrease=0.0, min_impurity_split=None,
5                             min_samples_leaf=1, min_samples_split=2,
6                             min_weight_fraction_leaf=0.0, n_estimators=100,
7                             n_jobs=None, oob_score=False, random_state=42,
8                             verbose=0, warm_start=False)
9 rfc1.fit(X, label)
```

After fitting the test data, I used the predict() of RandForestClassifier to predict the testData after putting it into the list of list1 and then vectorizing it using vectorizer.transform(list1) and then convert into toarray() and assign it into p. After that we converted p into the list.

```
1 list1=[]
2 for i in range(len(testData)):
3     testData[i]
4     list1.append(testData[i])
5     list1
6
7 p=rfc1.predict(vectorizer.transform(list1).toarray())
```

## Results

### 1. Calculating Accuracy

We have the predicted labels. To get the accuracy, we need the original test label so we will open and read it and then add all the sentiment columns entries into the list and in the end we have the original Labels in the testLabels list.

```
1 #ORIGINAL TEST LABELS
2
3 with open("/content/drive/MyDrive/Dataset/Test/test_labels_hinglish.txt") as f:
4     data = f.read().split('\n')
5
6 intermediateLabels= []
7 tweet = ""
8
9 i=0
10
11 for i in range(len(data)):
12
13     for j in range(len(data[i])):
14         if(data[i][j]==' '):
15             intermediateLabels.append(data[i][j+1:])
16
17 testLabels=[]
18
19 for i in range(len(intermediateLabels)-1):
20     testLabels.append(intermediateLabels[i+1])
21
22
```

```
1 from sklearn.metrics import accuracy_score
2
3 accuracy_score(testLabels, p)

0.6456666666666667
```

Note: First I simply called out the RandomForestClassifier() and I got a 63.5 score but when I used the following parameters and put the random\_state=42, the accuracy increased to 64.8%.

shehroz_308_c	2	12/09/20	0.648 (101)
---------------	---	----------	-------------

### 2. Getting the uids from testData

As we needed answer.txt which has two columns.

#### 1. Uids

To get the uids, I loaded the test dataset and appended all the uids from the file into the list and that's how I got the list of uids.

#### 2. Predicted labels/sentiment.

We already got the list i.e. p.

```
1 p=list(p)
2
3 #Taking out Uids
4 with open("/content/drive/MyDrive/Dataset/Test/Hindi_test_unalbelled_conll_updated.txt") as f:
5     data = f.read().split('\n')
6
7 uids = []
8 tweet = ""
9
10
11 for d in data:
12     if(len(d)-->0):
13         continue
14     elif(d.split()[0]=='meta'):
15         uids.append(d.split()[1])
16
```

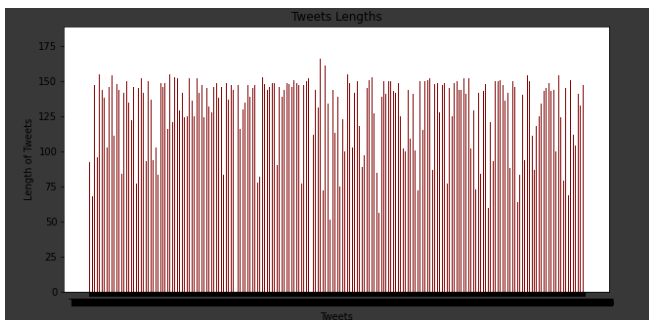
### 3. Saving the results in answer.txt

Save the sentiments along with uids in the answer.txt file

```
3 print(len(p))
4 count=0
5 store1 = []
6 for i in range(0, len(uids)):
7     save = ""
8     res = p[i]
9
10    if (res == "positive"):
11        save += uids[i]
12        save += ",positive\n"
13        count+=1
14
15    elif (res == "neutral"):
16        save += uids[i]
17        save += ",neutral\n"
18        count+=1
19
20    elif (res == "negative"):
21        save += uids[i]
22        save += ",negative\n"
23        count+=1
24    else:
25        print(uids[i])
26
27    store1.append(save)
28
29
30
31 with open('answer.txt','w',encoding='utf-8') as file:
32     i=0
33     while (i<len(store1)):
34         file.write(store1[i])
35         i+=1
```

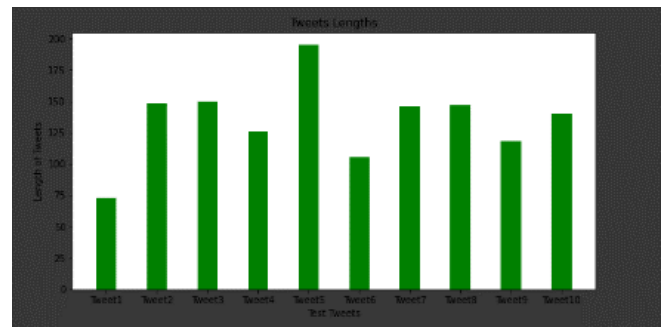
## Graphs

### 1. Length of tweets of Train Data



```
1 sizeTweets=[]
2 tweetNumber=[]
3 for i in range(len(trainData)):
4     sizeTweets.append(len(trainData[i]))
5     tweetNumber.append("Tweet"+str(i+1))
6
7 fig = plt.figure(figsize = (10, 5))
8
9 # creating the bar plot
10 plt.bar(tweetNumber, sizeTweets, color = 'maroon',
11         width = 0.4)
12
13 plt.xlabel("Tweets")
14 plt.ylabel("Length of Tweets")
15 plt.title("Tweets Lengths")
16 plt.show()
```

### 2. Length of first 10 tweets of Test Data



```
1 sizeTweets=[]
2 tweetNumber=[]
3 for i in range(len(testData[:10])):
4     sizeTweets.append(len(testData[i]))
5     tweetNumber.append("Tweet"+str(i+1))
6
7 fig = plt.figure(figsize = (10, 5))
8
9 # creating the bar plot
10 plt.bar(tweetNumber, sizeTweets, color = 'green',
11         width = 0.4)
12
13 plt.xlabel("Test Tweets")
14 plt.ylabel("Length of Tweets")
15 plt.title("Tweets Lengths")
16 plt.show()
```

## ACKNOWLEDGMENT

I would like to express my special thanks of gratitude to my Artificial Intelligence Teacher \*Mr. Umair Arshad\* for his able guidance and support in completing my Project.

## REFERENCES