

NLP assignment 2

Yogesh Kumar (C0852435)

In [1]:

```
import nltk
import numpy as np

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
import string

stopwords = nltk.corpus.stopwords.words('english')

path = 'C:/Users/Yogi/Downloads/'

X_train, X_test, y_train, y_test = train_test_split(
    messages, messages['label'], test_size=0.2, random_state=42)

messages.columns = ["label", "text"]
messages['label'] = np.where(messages['label'] == "spam", 1, 0)

def clean_text(text):
    text = " ".join(word.lower() for word in text if word not in string.punctuation)
    tokens = re.split('(\W+)', text)
    tokens = [token.strip() for token in tokens if token.strip() and token not in stopwords]
    return " ".join(tokens)

X_train_clean_text = [clean_text(x) for x in X_train]
X_test_clean_text = [clean_text(x) for x in X_test]

# Split data into train and test set
X_train, X_test, y_train, y_test = train_test_split(messages['clean_text'],
                                                    messages['label'], test_size=0.2,
                                                    random_state=42)

X_train.to_csv(path+'X_train.csv', index=False, header=True)
X_test.to_csv(path+'X_test.csv', index=False, header=True)
y_train.to_csv(path+'y_train.csv', index=False, header=True)
y_test.to_csv(path+'y_test.csv', index=False, header=True)
```

Implementing TF-IDF on Random Forest Classifier

In [2]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
import pandas as pd

tfidf_vect = TfidfVectorizer(analyzer='word', token_pattern=r'\b(?:\w+(?:\s+)?)+\b')
X_train_tfidf = tfidf_vect.fit_transform(X_train)
X_test_tfidf = tfidf_vect.transform(X_test)

print(X_train_tfidf.shape)
print(X_test_tfidf.shape)

from sklearn.metrics import accuracy_score, precision_score, recall_score, confusion_matrix, classification_report, accuracy_score

acc = accuracy_score(y_test, pred)
prec = precision_score(y_test, pred)

print("Accuracy Score :", acc)
print("Precision Score :", prec)

import seaborn as sns
sns.heatmap(confusion_matrix(y_test, pred), annot=True, cmap='magma', cbar=False, linewidth=3, linecolor='r',
           print_classification_report(y_test, pred))

Accuracy Score : 0.9739478053322287
Precision Score : 0.969824812030752
Recall Score : 0.969824812030752
f1-score support
0 0.98 1.00 0.99 963
1 0.98 0.98 0.99 152
accuracy
macro avg 0.9811500920439718 0.98 1115
weighted avg 0.9811500920439718 0.98 1115
```

In [3]:

```
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier(estimator='logit', criterion='entropy', random_state=0)
```

In [4]:

```
model.fit(X_train_tfidf, y_train)
pred = model.predict(X_test_tfidf)

from sklearn.metrics import precision_score, recall_score, confusion_matrix, classification_report, accuracy_score
acc = accuracy_score(y_test, pred)
prec = precision_score(y_test, pred)

print("Accuracy Score :", acc)
print("Precision Score :", prec)

import seaborn as sns
sns.heatmap(confusion_matrix(y_test, pred), annot=True, cmap='magma', cbar=False, linewidth=3, linecolor='r',
           print_classification_report(y_test, pred))

Accuracy Score : 0.9739478053322287
Precision Score : 0.969824812030752
Recall Score : 0.969824812030752
f1-score support
0 0.98 1.00 0.99 963
1 0.98 0.98 0.99 152
accuracy
macro avg 0.9811500920439718 0.98 1115
weighted avg 0.9811500920439718 0.98 1115
```

963

152

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

1115

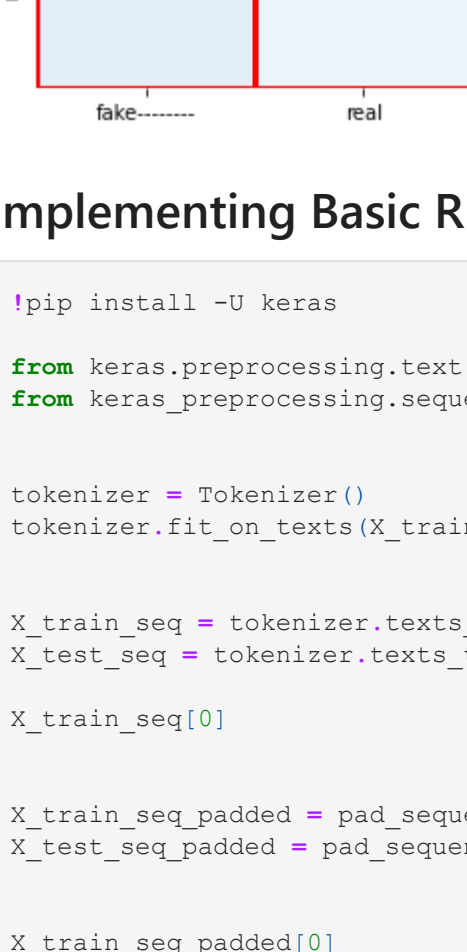
1115

1115

1115

(1115,)
(1115,)
Accuracy Score : 0.912107623318386
Precision Score : 0.964912807017544

	precision	recall	f1-score	support
0	0.91	1.00	0.95	963
1	0.96	0.36	0.53	152
accuracy			0.91	1115
macro avg	0.94	0.68	0.74	1115
weighted avg	0.92	0.91	0.89	1115



Implementing Basic RNN

```
In [12]: !pip install -U keras

from keras.preprocessing.text import Tokenizer # simple pre-process function from genesis celan and tokenize on
from keras.preprocessing.sequence import pad_sequences

tokenizer = Tokenizer()
tokenizer.fit_on_texts(X_train)

X_train_seq = tokenizer.texts_to_sequences(X_train)
X_test_seq = tokenizer.texts_to_sequences(X_test)

X_train_seq[0]

X_train_seq_padded = pad_sequences(X_train_seq, 50)
X_test_seq_padded = pad_sequences(X_test_seq, 50)

X_train_seq_padded[0]

import keras.backend as K
from keras.layers import Dense, Embedding, LSTM
from keras.models import Sequential

def recall_m(y_true, y_pred):
    true_positives = K.sum(K.round(K.clip(y_true * y_pred, 0, 1)))
    possible_positives = K.sum(K.round(K.clip(y_true, 0, 1)))
    recall = true_positives / (possible_positives + K.epsilon())
    return recall

def precision_m(y_true, y_pred):
    true_positives = K.sum(K.round(K.clip(y_true * y_pred, 0, 1)))
    predicted_positives = K.sum(K.round(K.clip(y_pred, 0, 1)))
    precision = true_positives / (predicted_positives + K.epsilon())
    return precision

# Construct a simple RNN model
model = Sequential()
model.add(Embedding(len(tokenizer.index_word*), 32))
model.add(LSTM(32, dropout=0, recurrent_dropout=0))
model.add(Dense(32, activation='relu'))
model.add(Dense(1, activation='sigmoid'))
model.summary()

model.compile(optimizer='adam',
              loss='binary_crossentropy',
              metrics=['accuracy', precision_m, recall_m])

11.
# Fit the RNN model
history = model.fit(X_train_seq_padded, y_train,
                   batch_size=32, epochs=10,
                   validation_data=(X_test_seq_padded, y_test))

12.

# Plot the evaluation metrics by each epoch for the model to see if we are over or underfitting
import matplotlib.pyplot as plt

for i in ['accuracy', 'precision_m', 'recall_m']:
    acc = history.history[i]
    val_acc = history.history['val_'+i].format(i)
    epochs = range(1, len(acc) + 1)

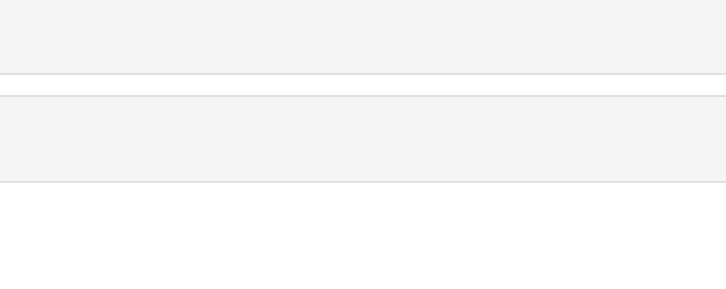
    plt.figure()
    plt.plot(epochs, acc, label='Training Accuracy')
    plt.plot(epochs, val_acc, label='Validation Accuracy')
    plt.title('Results for %s' % i)
    plt.legend()
    plt.show()

Requirement already satisfied: keras in c:\Users\yogiv\anaconda3\lib\site-packages (2.9.0)
Model: "sequential_1"
```

Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, None, 32)	265888
lstm (LSTM)	(None, 32)	8320
dense (Dense)	(None, 32)	1056
dense_1 (Dense)	(None, 1)	33

Total params: 275,297
Trainable params: 275,297
Non-trainable params: 0

Epoch 1/10
140/140 [=====] - 5s 23ms/step - loss: 0.2648 - accuracy: 0.9172 - precision_m: 0.5200
- recall_m: 0.3856 - val_loss: 0.0959 - val_accuracy: 0.9776 - val_precision_m: 0.9845 - val_recall_m: 0.8613
Epoch 2/10
140/140 [=====] - 3s 21ms/step - loss: 0.0369 - accuracy: 0.9910 - precision_m: 0.9670
- recall_m: 0.9376 - val_loss: 0.0598 - val_accuracy: 0.9830 - val_precision_m: 1.0000 - val_recall_m: 0.8801
Epoch 3/10
140/140 [=====] - 3s 22ms/step - loss: 0.0113 - accuracy: 0.9973 - precision_m: 0.9914
- recall_m: 0.9777 - val_loss: 0.0677 - val_accuracy: 0.9874 - val_precision_m: 1.0000 - val_recall_m: 0.9227
Epoch 4/10
140/140 [=====] - 3s 21ms/step - loss: 0.0034 - accuracy: 0.9996 - precision_m: 0.9857
- recall_m: 0.9825 - val_loss: 0.0746 - val_accuracy: 0.9839 - val_precision_m: 1.0000 - val_recall_m: 0.8970
Epoch 5/10
140/140 [=====] - 3s 25ms/step - loss: 0.0017 - accuracy: 0.9996 - precision_m: 1.0000
- recall_m: 0.9968 - val_loss: 0.0528 - val_accuracy: 0.9874 - val_precision_m: 0.9636 - val_recall_m: 0.9442
Epoch 6/10
140/140 [=====] - 3s 24ms/step - loss: 4.3904e-04 - accuracy: 1.0000 - precision_m: 1.0000
- recall_m: 1.0000 - val_loss: 0.0695 - val_accuracy: 0.9874 - val_precision_m: 0.9845 - val_recall_m: 0.9251
Epoch 7/10
140/140 [=====] - 3s 23ms/step - loss: 1.5941e-04 - accuracy: 1.0000 - precision_m: 0.9885
- recall_m: 0.9857 - val_loss: 0.0645 - val_accuracy: 0.9883 - val_precision_m: 0.9752 - val_recall_m: 0.9401
Epoch 8/10
140/140 [=====] - 3s 20ms/step - loss: 9.9409e-05 - accuracy: 1.0000 - precision_m: 0.9929
- recall_m: 0.9929 - val_loss: 0.0676 - val_accuracy: 0.9883 - val_precision_m: 0.9752 - val_recall_m: 0.9401
Epoch 9/10
140/140 [=====] - 3s 19ms/step - loss: 6.9337e-05 - accuracy: 1.0000 - precision_m: 0.9786
- recall_m: 0.9786 - val_loss: 0.0705 - val_accuracy: 0.9883 - val_precision_m: 0.9752 - val_recall_m: 0.9401
Epoch 10/10
140/140 [=====] - 3s 21ms/step - loss: 5.1763e-05 - accuracy: 1.0000 - precision_m: 0.9714
- recall_m: 0.9714 - val_loss: 0.0726 - val_accuracy: 0.9883 - val_precision_m: 0.9752 - val_recall_m: 0.9401



Summary:

```
In [13]: print(f"The accuracy and precision score of TF_IDF on Random forest are {acc1} and {pre1} respectively.\n The accuracy and precision score of word2vec on Random forest is {acc2} and {pre2} respectively.\n The accuracy and precision score of doc2vec on Random forest is {acc3} and {pre3} respectively.\n The accuracy and precision score of basic RNN is {acc4} and {pre4} respectively.")
```

The accuracy and precision score of TF_IDF on Random forest are 0.9757847533632287 and 0.9699248120300752 respectively.
The accuracy and precision score of word2vec on Random forest is 0.9730941704035875 and 0.952238805970149 respectively.
The accuracy and precision score of doc2vec on Random forest is 0.9112107623318386 and 0.964912807017544 respectively.
The accuracy and precision score of basic RNN is 0.9208761274814605 and 0.9786 respectively.

Conclusion: TF_IDF is best technique

we can conclude that TF_IDF is best algorithm as we can see accuracy is more than other models. The advantage of TF-IDF is its simplicity and ease of use. It is easy to compute, has low computational cost, and is an easy starting point for similarity computation. TF-IDF is a algorithm that uses the frequency of words to determine how relevant those words are to a particular document. This is a relatively simple than other algorithms.

```
In [ ]: 
```

```
In [ ]: 
```

```
In [ ]: 
```

```
In [ ]: 
```

```
In [ ]: 
```