

```

1 import nltk
2 nltk.download('gutenberg')
3 nltk.download('punkt')
4 import re
5 from nltk.stem import WordNetLemmatizer
6 from nltk.tokenize import word_tokenize
7 nltk.download('omw-1.4')
8 nltk.download('wordnet')
9 import pandas as pd
10 import string

[nltk_data] Downloading package gutenberg to /root/nltk_data...
[nltk_data] Unzipping corpora/gutenberg.zip.
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data] Downloading package wordnet to /root/nltk_data...

1 nltk.corpus.gutenberg.fileids()

['austen-emma.txt',
 'austen-persuasion.txt',
 'austen-sense.txt',
 'bible-kjv.txt',
 'blake-poems.txt',
 'bryant-stories.txt',
 'burgess-busterbrown.txt',
 'carroll-alice.txt',
 'chesterton-ball.txt',
 'chesterton-brown.txt',
 'chesterton-thursday.txt',
 'edgeworth-parents.txt',
 'melville-moby_dick.txt',
 'milton-paradise.txt',
 'shakespeare-caesar.txt',
 'shakespeare-hamlet.txt',
 'shakespeare-macbeth.txt',
 'whitman-leaves.txt']

1 #preprocessing
2 #grabbing the Shakespeare works
3
4 caesar = nltk.corpus.gutenberg.raw('shakespeare-caesar.txt')
5 hamlet = nltk.corpus.gutenberg.raw('shakespeare-hamlet.txt')
6 mcbeth = nltk.corpus.gutenberg.raw('shakespeare-macbeth.txt')
7
8 #lowercasing everything
9
10 c_lower = caesar.lower()
11 h_lower = hamlet.lower()
12 m_lower = mcbeth.lower()
13
14 #tokenizing the text blobs
15
16 from nltk.tokenize import sent_tokenize
17 caesar_tok = sent_tokenize(c_lower)
18 hamlet_tok = sent_tokenize(h_lower)
19 mcbeth_tok = sent_tokenize(m_lower)

1 def remove_punct(str_list):
2     no_punct = []
3
4     for sent in str_list:
5         sent = re.sub('[%s]' % re.escape(string.punctuation), '', sent)
6
7         no_punct.append(sent)
8     return no_punct

1 caesar_tok = remove_punct(caesar_tok)
2 hamlet_tok = remove_punct(hamlet_tok)
3 mcbeth_tok = remove_punct(mcbeth_tok)
4
5 #Make Shakespeare plays into dataframes
6 import pandas as pd
7 ceasar_df = pd.DataFrame(caesar_tok)

```

```

8 author = 'Shakespeare'
9 ceasar_df['Author'] = author
10 ceasar_df = ceasar_df.rename(columns={0: 'Text'})
11
12 hamlet_df = pd.DataFrame(hamlet_tok)
13 hamlet_df['Author'] = author
14 hamlet_df = hamlet_df.rename(columns={0: 'Text'})
15
16 mcbeth_df = pd.DataFrame(mcbeth_tok)
17 mcbeth_df['Author'] = author
18 mcbeth_df = mcbeth_df.rename(columns={0: 'Text'})
19
20 #creating a non-Shakespeare corpus
21
22 bible = nltk.corpus.gutenberg.raw('bible-kjv.txt')
23 milton = nltk.corpus.gutenberg.raw('milton-paradise.txt')
24 blake = nltk.corpus.gutenberg.raw('blake-poems.txt')
25
26 b_lower = bible.lower()
27 mil_lower = milton.lower()
28 bl_lower = blake.lower()
29
30 #tokenizing the text-blob
31
32 bible_tok = sent_tokenize(b_lower)
33 milton_tok = sent_tokenize(mil_lower)
34 blake_tok = sent_tokenize(bl_lower)
35
36 bible_tok = remove_punct(bible_tok)
37 milton_tok = remove_punct(milton_tok)
38 blake_tok = remove_punct(blake_tok)
39
40 #making the non-Shakespeare all one dataframe.
41
42 author_2 = 'Other'
43
44 bible_df = pd.DataFrame(bible_tok)
45 bible_df['Author'] = author_2
46 bible_df = bible_df.rename(columns={0: 'Text'})
47
48 milton_df = pd.DataFrame(milton_tok)
49 milton_df['Author'] = author_2
50 milton_df = milton_df.rename(columns={0: 'Text'})
51
52 blake_df = pd.DataFrame(blake_tok)
53 blake_df['Author'] = author_2
54 blake_df = blake_df.rename(columns={0: 'Text'})
55
56 noshakespear = bible_df.append([milton_df, blake_df], ignore_index=True)
57
58 #subsampling No Shakespeare
59
60 import random
61 random.seed(1)
62
63 noshakes_sub = noshakespear.sample(n = 3500, random_state=1)
64
65 #composite data frames
66
67 no_ceasar = hamlet_df.append([noshakes_sub, mcbeth_df], ignore_index=True)
68 no_hamlet = ceasar_df.append([noshakes_sub, mcbeth_df], ignore_index=True)
69 no_mcbeth = hamlet_df.append([noshakes_sub, ceasar_df], ignore_index=True)

<ipython-input-5-d1a1d18a8dd7>:56: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future vers
noshakespear = bible_df.append([milton_df, blake_df], ignore_index=True)
<ipython-input-5-d1a1d18a8dd7>:67: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future vers
no_ceasar = hamlet_df.append([noshakes_sub, mcbeth_df], ignore_index=True)
<ipython-input-5-d1a1d18a8dd7>:68: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future vers
no_hamlet = ceasar_df.append([noshakes_sub, mcbeth_df], ignore_index=True)
<ipython-input-5-d1a1d18a8dd7>:69: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future vers
no_mcbeth = hamlet_df.append([noshakes_sub, ceasar_df], ignore_index=True)

1 no_ceasar['Author'] = no_ceasar['Author'].map({'Shakespeare':1, 'Other':0})
2 print(no_ceasar)

```

```

    Text Author
0 the tragedie of hamlet by william shakespeare ... 1
1 scoena prima 1
2 enter barnardo and francisco two centinels 1
3 barnardo 1
4 whos there 1
... ... ...
7268 this and what need full else\nthat calls vpon ... 1
7269 flourish 1
7270 exeunt omnes 1
7271 finis 1
7272 the tragedie of macbeth 1

[7273 rows x 2 columns]

1 #Test Train - The old Fashioned Way
2
3 from sklearn.model_selection import train_test_split
4
5 c_X = no_ceasar.Text
6 c_y = no_ceasar.Author
7 print(c_X.shape, c_y.shape)
8
9 c_X_train, c_X_test, c_y_train, c_y_test = train_test_split(c_X, c_y, random_state=1)
10 print(c_X_train.shape, c_y_train.shape)
11 print(c_X_test.shape, c_y_test.shape)

(7273,) (7273,)
(5454,) (5454,)
(1819,) (1819,)

1 no_ceasar_labels = no_ceasar['Author']
2 no_ceasar_text = no_ceasar['Text']

1 # pre processing
2
3 MAX_SEQUENCE_LENGTH = 300
4 MAX_NUM_WORDS = 20000
5 EMBEDDING_DIM = 100
6 VALIDATION_SPLIT = 0.3

1 from keras.preprocessing.text import Tokenizer

1 tokenizer = Tokenizer(num_words=MAX_NUM_WORDS)
2 tokenizer.fit_on_texts(c_X_train)
3 train_sequences = tokenizer.texts_to_sequences(c_X_train) #Converting text to a vector of word indexes
4 test_sequences = tokenizer.texts_to_sequences(c_X_test)
5 word_index = tokenizer.word_index
6 print('Found %s unique tokens.' % len(word_index))

Found 11460 unique tokens.

1 from sklearn.preprocessing import LabelEncoder
2
3 le = LabelEncoder()
4 le.fit(c_y_train)
5 train_labels = le.transform(c_y_train)
6 test_labels = le.transform(c_y_test)

1 from tensorflow.keras.preprocessing.sequence import pad_sequences
2 from keras.utils import to_categorical
3 import numpy as np
4
5 #Converting this to sequences to be fed into neural network. Max seq. len is 1000 as set earlier
6 #initial padding of 0s, until vector is of size MAX_SEQUENCE_LENGTH
7 trainvalid_data = pad_sequences(train_sequences, maxlen=MAX_SEQUENCE_LENGTH)
8 test_data = pad_sequences(test_sequences, maxlen=MAX_SEQUENCE_LENGTH)
9 trainvalid_labels = to_categorical(train_labels)
10
11 test_labels = to_categorical(np.asarray(test_labels), num_classes= trainvalid_labels.shape[1])
12
13 # split the training data into a training set and a validation set
14 indices = np.arange(trainvalid_data.shape[0])

```

```

15 np.random.shuffle(indices)
16 trainvalid_data = trainvalid_data[indices]
17 trainvalid_labels = trainvalid_labels[indices]
18 num_validation_samples = int(VALIDATION_SPLIT * trainvalid_data.shape[0])
19 x_train = trainvalid_data[:-num_validation_samples]
20 y_train = trainvalid_labels[:-num_validation_samples]
21 x_val = trainvalid_data[-num_validation_samples:]
22 y_val = trainvalid_labels[-num_validation_samples:]
23 #This is the data we will use for CNN training
24 print('Splitting the train data into train and valid is done')

    Splitting the train data into train and valid is done

1 import os
2
3 BASE_DIR = os.getcwd()
4 GLOVE_DIR = os.path.join(BASE_DIR, 'glove.6B')

1 from keras.layers import Conv1D, MaxPooling1D, Embedding
2 from keras.initializers import Constant
3
4 #embedding matrix
5
6 print('Preparing embedding matrix.')
7 # first, build index mapping words in the embeddings set
8 # to their embedding vector
9 embeddings_index = {}
10
11 with open( 'glove.6B.100d.txt',encoding="utf8") as f:
12     for line in f:
13         values = line.split()
14         word = values[0]
15         coefs = np.asarray(values[1:], dtype='float32')
16         embeddings_index[word] = coefs
17
18 print('Found %s word vectors in Glove embeddings.' % len(embeddings_index))
19 #print(embeddings_index["google"])
20
21 # prepare embedding matrix - rows are the words from word_index, columns are the embeddings of that word from glove.
22 num_words = min(MAX_NUM_WORDS, len(word_index)) + 1
23 embedding_matrix = np.zeros((num_words, EMBEDDING_DIM))
24 for word, i in word_index.items():
25     if i > MAX_NUM_WORDS:
26         continue
27     embedding_vector = embeddings_index.get(word)
28     if embedding_vector is not None:
29         # words not found in embedding index will be all-zeros.
30         embedding_matrix[i] = embedding_vector
31
32 # load these pre-trained word embeddings into an Embedding layer
33 # note that we set trainable = False so as to keep the embeddings fixed
34 embedding_layer = Embedding(num_words,
35                             EMBEDDING_DIM,
36                             embeddings_initializer=Constant(embedding_matrix),
37                             input_length=MAX_SEQUENCE_LENGTH,
38                             trainable=False)
39 print("Preparing of embedding matrix is done")

    Preparing embedding matrix.
    Found 2449 word vectors in Glove embeddings.
    Preparing of embedding matrix is done

1 from keras.models import Model, Sequential
2 from keras.layers import Dense, Input, GlobalMaxPooling1D
3
4 cnnmodel = Sequential()
5 cnnmodel.add(embedding_layer)
6 cnnmodel.add(Conv1D(128, 5, activation='relu'))
7 cnnmodel.add(MaxPooling1D(5))
8 cnnmodel.add(Conv1D(128, 5, activation='relu'))
9 cnnmodel.add(MaxPooling1D(5))
10 cnnmodel.add(Conv1D(128, 5, activation='relu'))
11 cnnmodel.add(GlobalMaxPooling1D())
12 cnnmodel.add(Dense(128, activation='relu'))
13 cnnmodel.add(Dense(len(trainvalid_labels[0]), activation='softmax'))

```

```

14
15 cnnmodel.compile(loss='categorical_crossentropy',
16                 optimizer='rmsprop',
17                 metrics=['acc'])
18
19 cnnmodel.summary()
20
21 #you can replace softmax by hyperbolic

```

Model: "sequential"

| Layer (type) | Output Shape | Param # |
|---|------------------|---------|
| embedding (Embedding) | (None, 300, 100) | 1146100 |
| conv1d (Conv1D) | (None, 296, 128) | 64128 |
| max_pooling1d (MaxPooling1D) | (None, 59, 128) | 0 |
| conv1d_1 (Conv1D) | (None, 55, 128) | 82048 |
| max_pooling1d_1 (MaxPooling1D) | (None, 11, 128) | 0 |
| conv1d_2 (Conv1D) | (None, 7, 128) | 82048 |
| global_max_pooling1d (GlobalMaxPooling1D) | (None, 128) | 0 |
| dense (Dense) | (None, 128) | 16512 |
| dense_1 (Dense) | (None, 2) | 258 |

=====
 Total params: 1,391,094
 Trainable params: 244,994
 Non-trainable params: 1,146,100
 =====

```

1 #Train the model. Tune to validation set.
2 cnnmodel.fit(x_train, y_train,
3             batch_size=128,
4             epochs=3, validation_data=(x_val, y_val))
5 #Evaluate on test set:
6 score, acc = cnnmodel.evaluate(test_data, test_labels)
7 print('Test accuracy with CNN:', acc) #consider decreasing batch size and increasing epoch size

Epoch 1/3
30/30 [=====] - 25s 807ms/step - loss: 0.5859 - acc: 0.7391 - val_loss: 0.4406 - val_acc: 0.8081
Epoch 2/3
30/30 [=====] - 20s 663ms/step - loss: 0.4406 - acc: 0.8070 - val_loss: 0.3986 - val_acc: 0.8417
Epoch 3/3
30/30 [=====] - 15s 500ms/step - loss: 0.3783 - acc: 0.8363 - val_loss: 0.3378 - val_acc: 0.8570
57/57 [=====] - 4s 62ms/step - loss: 0.3408 - acc: 0.8505
Test accuracy with CNN: 0.8504672646522522

```

Confusion Matrix

```

1 #x_val, y_val, test_data, test_labels
2
3 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
4
5 class_labels = [0, 1]
6
7 Y_test_preds = cnnmodel.predict(test_data)
8 rounded_labels = np.argmax(test_labels, axis=1)
9 Y_test_preds = np.argmax(Y_test_preds, axis=1)
10
11 print("Test Accuracy : {}".format(accuracy_score(rounded_labels, Y_test_preds)))
12 print("\nConfusion Matrix : ")
13 print(confusion_matrix(rounded_labels, Y_test_preds, labels=[0, 1]))
14 print("\nClassification Report :")
15 #print(classification_report(rounded_labels, Y_test_preds, target_names=class_labels))
16
17 #https://stackoverflow.com/questions/54589669/confusion-matrix-error-classification-metrics-cant-handle-a-mix-of-multilabel
18 #https://coderrzcolumn.com/tutorials/artificial-intelligence/lime-explain-keras-image-classification-network-predictions

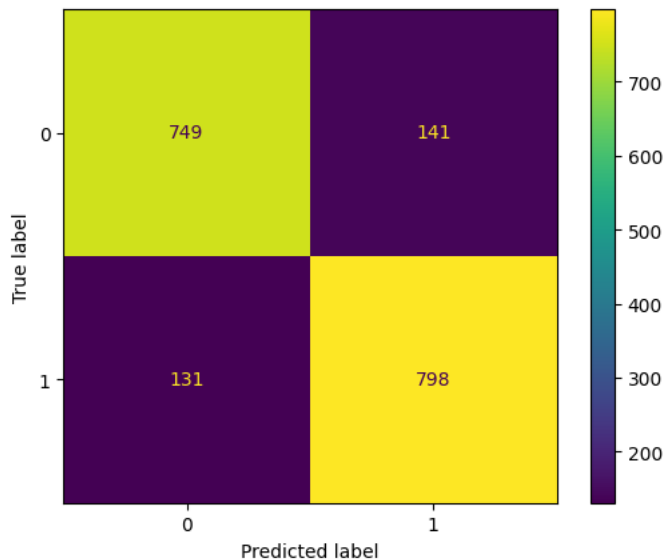
```

```
57/57 [=====] - 2s 37ms/step
Test Accuracy : 0.8504672897196262
```

```
Confusion Matrix :
[[749 141]
 [131 798]]
```

```
Classification Report :
```

```
1 from sklearn.metrics import ConfusionMatrixDisplay
2 import matplotlib.pyplot as plt
3
4 cm = confusion_matrix(rounded_labels, Y_test_preds, labels=class_labels)
5 disp = ConfusionMatrixDisplay(confusion_matrix=cm,
6                               display_labels=class_labels)
7 disp.plot()
8 plt.show()
9
10 #https://stackoverflow.com/questions/54589669/confusion-matrix-error-classification-metrics-cant-handle-a-mix-of-multilabel
```



Predicting with CNN

```
1 #Predicting with a CNN
2
3 #experimental set
4 #caesar_experiment = vect.transform(caesar_tok)
5
6 caesar_tok = tokenizer.texts_to_sequences(caesar_tok)
7 caesar = pad_sequences(caesar_tok, maxlen=MAX_SEQUENCE_LENGTH)
8
9
10 #https://medium.com/analytics-vidhya/multi-class-classification-using-cnn-for-custom-dataset-7759865bd19
11 #https://towardsdatascience.com/how-to-predict-an-image-with-keras-ca97d9cd4817
12 #https://www.analyticsvidhya.com/blog/2021/12/intent-classification-with-convolutional-neural-networks/
```

```
1 is_ceasar = cnnmodel.predict(caesar)
```

```
49/49 [=====] - 2s 36ms/step
```

```
1 print(is_ceasar)
```

```
[[0.6988547 0.30114526]
 [0.11961149 0.88038856]
 [0.17559522 0.8244048 ]
 ...
 [0.11961149 0.88038856]
 [0.11961149 0.88038856]
 [0.21503216 0.7849678  ]]
```

1 #DF Analysis

```
1 ceasar_df['Author'] = author
2 ceasar_df = ceasar_df.rename(columns={0: 'Text'})
3
4 is_ceasar_df = pd.DataFrame(is_ceasar)
5 is_ceasar_df['Text'] = ceasar_df['Text']
6
7 is_ceasar_df = is_ceasar_df.rename(columns={0: 'Other'})
8 is_ceasar_df = is_ceasar_df.rename(columns={1: 'Shakespeare'})
```

```
1 def count_if_1(series):
2     count = 0
3
4     for cell in series:
5         if cell >= 0.5:
6             count += 1
7     return count
```

```
1 count_if_1(is_ceasar_df['Shakespeare'])

1380
```

1 is_ceasar_df

| | Other | Shakespeare | Text |
|------|----------|-------------|---|
| 0 | 0.698855 | 0.301145 | the tragedie of julius caesar by william shake... |
| 1 | 0.119611 | 0.880389 | scoena prima |
| 2 | 0.175595 | 0.824405 | enter flavius murellus and certaine commoners ... |
| 3 | 0.119611 | 0.880389 | flavius |
| 4 | 0.057073 | 0.942927 | hence home you idle creatures get you home\nis... |
| ... | ... | ... | ... |
| 1547 | 0.177149 | 0.822851 | within my tent his bones to night shall ly\nmo... |
| 1548 | 0.119611 | 0.880389 | exeunt |
| 1549 | 0.119611 | 0.880389 | omnes |
| 1550 | 0.119611 | 0.880389 | finis |
| 1551 | 0.215032 | 0.784968 | the tragedie of ivlivs caesar |

1552 rows × 3 columns