# CODE: Sentiment Classification - Impact of Pretrained Word Embedding

```
1   import numpy as np
2
3   import tensorflow_datasets as tfds
4   import tensorflow as tf
5
6   data, info = tfds.load('imdb_reviews', with_info=True, as_supervised=True)
7   train, test = data['train'], data['test']
```

**Downloading and preparing dataset imdb_reviews/plain_text/1.0.0 (download: 80.23**

Dl Completed...: 100%      1/1 [00:05<00:00, 5.32s/ url]

Dl Size...: 100%      80/80 [00:05<00:00, 23.93 MiB/s]


Shuffling and writing examples to /root/tensorflow_datasets/imdb_reviews/plain_t
100%                                                24999/25000 [00:00<00:00, 116024.35 examples/s]
Shuffling and writing examples to /root/tensorflow_datasets/imdb_reviews/plain_t
100%                                                24999/25000 [00:00<00:00, 118347.06 examples/s]
Shuffling and writing examples to /root/tensorflow_datasets/imdb_reviews/plain_t
100%                                                49999/50000 [00:00<00:00, 133471.40 examples/s]
WARNING:absl:Dataset is using deprecated text encoder API which will be removed
**Dataset imdb_reviews downloaded and prepared to /root/tensorflow_datasets/imdb_r**

```
1 VOCAB_SIZE = 10000
2 BATCH_SIZE = 64
3 embedding_dim=100
4 max_length = 120
5 trunc_type= 'post'
6 UNK="<UNK>"
```

```
1 ## Processing the data
2
3 train_input = []  #X_train
4 test_input = []   #X_test
5 train_labels = [] #y_train
6 test_labels = []  #y_test
7
8 for i,j in train:
9     train_input.append(str(i.numpy()))
```

```
10      train_labels.append(j.numpy())
11
12 for i,j in test:
13      test_input.append(str(i.numpy()))
14      test_labels.append(j.numpy())
15
16
17 train_labels = np.array(train_labels)
18 test_labels = np.array(test_labels)
```

```
1 print(test_input[1])
```

```
  e they concern his moral integrity and we are never quite sure whether it remains
```

```
1 from tensorflow.keras.preprocessing.text import Tokenizer
2 from tensorflow.keras.preprocessing.sequence import pad_sequences
3
4 tokenizer = Tokenizer(num_words = VOCAB_SIZE, oov_token=UNK)
5 tokenizer.fit_on_texts(train_input)
6 word_index = tokenizer.word_index
```

```
1 seq = tokenizer.texts_to_sequences(train_input)
2 padding = pad_sequences(seq, maxlen=max_length, truncating = trunc_type)
3 test_seq = tokenizer.texts_to_sequences(test_input)
4 test_padded = pad_sequences(test_seq, maxlen=max_length)
```

```
1 print("train : ", len(seq[0]))
2 print("test : ", len(test_seq[0]))
```

```
    train :  118
    test :  173
```

```
1 import nltk
2 nltk.download('punkt')
3 nltk.download('stopwords')
```

```
    [nltk_data] Downloading package punkt to /root/nltk_data...
    [nltk_data]   Package punkt is already up-to-date!
    [nltk_data] Downloading package stopwords to /root/nltk_data...
    [nltk_data]   Unzipping corpora/stopwords.zip.
    True
```

```
1 import string
2 from nltk.tokenize import word_tokenize
3 from nltk.corpus import stopwords
4
5 output = []
6
```

```
 7 for line in train_input:
 8   tokens = word_tokenize(line)
 9   tokens = [word.lower() for word in tokens]
10   processed = str.maketrans('','', string.punctuation)
11   translated = [word.translate(processed) for word in tokens]
12   words = [word for word in translated if word.isalpha()]
13   stopWords = set(stopwords.words('english'))
14   words = [word for word in words if not word in stopWords]
15   output.append(words)
```

```
1 print(output[0])
2 print(len(output))
```

```
['b', 'absolutely', 'terrible', 'movie', 'nt', 'lured', 'christopher', 'walken',
25000
```

```
1   import gensim
2
3   model_1 = gensim.models.Word2Vec(sentences = output, size= embedding_dim, window
4   words = list(model_1.wv.vocab)
```

```
1 model_1.wv.save_word2vec_format("w2v_embedding.txt", binary = False)
```

```
1 import os
2 path_to_w2v_file = "w2v_embedding.txt"
3
4 embeddings_index = {}
5 with open(path_to_w2v_file) as f:
6     for line in f:
7         word, coefs = line.split(maxsplit=1)
8         coefs = np.fromstring(coefs, "f", sep=" ")
9         embeddings_index[word] = coefs
10
11 print("Found %s word vectors." % len(embeddings_index))
```

```
    Found 93966 word vectors.
```

```
 1 num_tokens = len(word_index) + 1
 2 hits = 0
 3 misses = 0
 4
 5 embedding_matrix = np.zeros((num_tokens, embedding_dim))
 6
 7 for word, i in word_index.items():
 8     embedding_vector = embeddings_index.get(word)
 9     if embedding_vector is not None:
10         # Words not found in embedding index will be all-zeros.
11         # This includes the representation for "padding" and "OOV"
```

```
12          embedding_matrix[i] = embedding_vector
13          hits += 1
14      else:
15          misses += 1
16  print("Converted %d words (%d misses)" % (hits, misses))

     Converted 70670 words (15869 misses)
```

```
1 from keras.layers import Embedding
2 from keras.initializers import Constant
3
4 embedding_layer = Embedding(num_tokens, embedding_dim, embeddings_initializer= Co
```

```
1 ## LSTM Model
2 from keras.models import Sequential
3 from keras.layers import Dense, LSTM, GRU
4 from keras.layers.embeddings import Embedding
5
6 model_1 = Sequential()
7 model_1.add(embedding_layer)
8 model_1.add(LSTM(units=64, dropout=0.2))
9 model_1.add(Dense(32, activation='relu'))
10 model_1.add(Dense(1, activation='sigmoid'))
11
12 model_1.summary()
```

```
Model: "sequential"

_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_1 (Embedding)     (None, 120, 100)          8654000

 lstm (LSTM)                 (None, 64)                42240

 dense (Dense)               (None, 32)                2080

 dense_1 (Dense)             (None, 1)                 33

=================================================================
Total params: 8,698,353
Trainable params: 44,353
Non-trainable params: 8,654,000
_____
```

```
1 model_1.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'
2 history = model_1.fit(padding, train_labels, epochs=20, validation_data = (test_p
```

```
Epoch 1/20
782/782 [==============================] - 37s 38ms/step - loss: 0.5588 - accura
Epoch 2/20
```

```
782/782 [==============================] – 29s 37ms/step – loss: 0.4431 – accura
Epoch 3/20
782/782 [==============================] – 29s 37ms/step – loss: 0.4215 – accura
Epoch 4/20
782/782 [==============================] – 29s 37ms/step – loss: 0.4035 – accura
Epoch 5/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3999 – accura
Epoch 6/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3898 – accura
Epoch 7/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3851 – accura
Epoch 8/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3772 – accura
Epoch 9/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3726 – accura
Epoch 10/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3632 – accura
Epoch 11/20
782/782 [==============================] – 31s 40ms/step – loss: 0.3573 – accura
Epoch 12/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3478 – accura
Epoch 13/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3412 – accura
Epoch 14/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3316 – accura
Epoch 15/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3266 – accura
Epoch 16/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3170 – accura
Epoch 17/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3071 – accura
Epoch 18/20
782/782 [==============================] – 28s 36ms/step – loss: 0.3018 – accura
Epoch 19/20
782/782 [==============================] – 29s 37ms/step – loss: 0.2941 – accura
Epoch 20/20
782/782 [==============================] – 28s 36ms/step – loss: 0.2828 – accura
```

```
 1 from sklearn import metrics
 2 from sklearn.metrics import precision_score, recall_score
 3
 4 predicted_1 = model_1.predict(test_padded)
 5
 6 pred_1 = np.zeros(len(predicted_1))
 7 for i, score in enumerate(predicted_1):
 8     if score > 0.5:
 9       pred_1[i] = 1
10     else:
11       pred_1[i] = 0
12
13 print(metrics.classification_report(test_labels, pred_1))
14 print("Precision = ",precision_score(test_labels, pred_1))
15 print("Recall = ", recall_score(test_labels, pred_1))
```

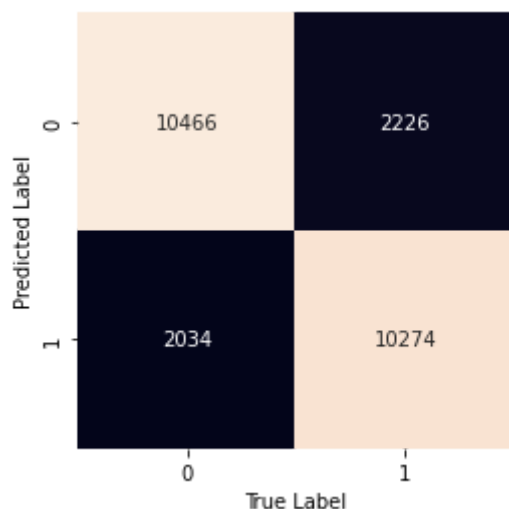|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.82      | 0.84   | 0.83     | 12500   |
| 1            | 0.83      | 0.82   | 0.83     | 12500   |
|              |           |        |          |         |
| accuracy     |           |        | 0.83     | 25000   |
| macro avg    | 0.83      | 0.83   | 0.83     | 25000   |
| weighted avg | 0.83      | 0.83   | 0.83     | 25000   |

```
Precision =  0.8347416314592135
Recall =  0.82192
```

```
1 from sklearn.metrics import confusion_matrix
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 mat = confusion_matrix(test_labels, pred_1)
6 sns.heatmap(mat.T, square = True, annot = True, fmt = 'd', cbar = False)
7 plt.xlabel("True Label")
8 plt.ylabel("Predicted Label")
```

```
Text(91.68, 0.5, 'Predicted Label')
```



## GLOVE

```
1 !wget http://nlp.stanford.edu/data/glove.6B.zip
2 !unzip -q glove.6B.zip
```

```
--2022-03-07 05:16:00--  http://nlp.stanford.edu/data/glove.6B.zip
Resolving nlp.stanford.edu (nlp.stanford.edu)... 171.64.67.140
Connecting to nlp.stanford.edu (nlp.stanford.edu)|171.64.67.140|:80... connected
HTTP request sent, awaiting response... 302 Found
Location: https://nlp.stanford.edu/data/glove.6B.zip [following]
--2022-03-07 05:16:00--  https://nlp.stanford.edu/data/glove.6B.zip
Connecting to nlp.stanford.edu (nlp.stanford.edu)|171.64.67.140|:443... connecte
HTTP request sent, awaiting response... 301 Moved Permanently
Location: http://downloads.cs.stanford.edu/nlp/data/glove.6B.zip [following]
```

```
--2022-03-07 05:16:00--  http://downloads.cs.stanford.edu/nlp/data/glove.6B.zip
Resolving downloads.cs.stanford.edu (downloads.cs.stanford.edu)... 171.64.64.22
Connecting to downloads.cs.stanford.edu (downloads.cs.stanford.edu)|171.64.64.22
HTTP request sent, awaiting response... 200 OK
Length: 862182613 (822M) [application/zip]
Saving to: 'glove.6B.zip'

glove.6B.zip          100%[===================>] 822.24M  5.11MB/s    in 2m 40s

2022-03-07 05:18:41 (5.13 MB/s) - 'glove.6B.zip' saved [862182613/862182613]
```

```python
1 import os
2 path_to_glove_file = "glove.6B.100d.txt"
3
4 embeddings_index = {}
5 with open(path_to_glove_file) as f:
6     for line in f:
7         word, coefs = line.split(maxsplit=1)
8         coefs = np.fromstring(coefs, "f", sep=" ")
9         embeddings_index[word] = coefs
10
11 print("Found %s word vectors." % len(embeddings_index))
```

```
Found 400000 word vectors.
```

```python
1 num_tokens = len(word_index) + 1
2 embedding_dim = 100 ## 100 dimensions
3 hits = 0 ## number of words that were found in the pretrained model
4 misses = 0 ## number of words that were missing in the pretrained model
5
6 # Prepare embedding matrix for our word list
7 embedding_matrix = np.zeros((num_tokens, embedding_dim))
8 for word, i in word_index.items():
9     embedding_vector = embeddings_index.get(word)
10     if embedding_vector is not None:
11         # Words not found in embedding index will be all-zeros.
12         # This includes the representation for "padding" and "OOV"
13         embedding_matrix[i] = embedding_vector
14         hits += 1
15     else:
16         misses += 1
17 print("Converted %d words (%d misses)" % (hits, misses))
```

```
Converted 60197 words (26342 misses)
```

```python
1 from keras.layers import Embedding
2 from keras.initializers import Constant
3
4 embedding_layer = Embedding(num_tokens, embedding_dim, embeddings_initializer= Co
```

```
 1 ## LSTM Model
 2 from keras.models import Sequential
 3 from keras.layers import Dense, LSTM, GRU
 4 from keras.layers.embeddings import Embedding
 5
 6 model_2 = Sequential()
 7 model_2.add(embedding_layer)
 8 model_2.add(LSTM(units=64, dropout=0.2))
 9 model_2.add(Dense(32, activation='relu'))
10 model_2.add(Dense(1, activation='sigmoid'))
11
12 model_2.summary()
```

```
Model: "sequential_2"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_2 (Embedding)     (None, 120, 100)          8654000

 lstm_2 (LSTM)               (None, 64)                42240

 dense_4 (Dense)             (None, 32)                2080

 dense_5 (Dense)             (None, 1)                 33

=================================================================
Total params: 8,698,353
Trainable params: 44,353
Non-trainable params: 8,654,000
_____
```

```
 1 model_2.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'
 2 history = model_2.fit(padding, train_labels, epochs=20, validation_data = (test_p
```

```
Epoch 1/20
782/782 [==============================] - 31s 37ms/step - loss: 0.6761 - accura
Epoch 2/20
782/782 [==============================] - 28s 36ms/step - loss: 0.4945 - accura
Epoch 3/20
782/782 [==============================] - 28s 36ms/step - loss: 0.4385 - accura
Epoch 4/20
782/782 [==============================] - 28s 36ms/step - loss: 0.4172 - accura
Epoch 5/20
782/782 [==============================] - 28s 36ms/step - loss: 0.4015 - accura
Epoch 6/20
782/782 [==============================] - 28s 36ms/step - loss: 0.3898 - accura
Epoch 7/20
782/782 [==============================] - 28s 36ms/step - loss: 0.3760 - accura
Epoch 8/20
782/782 [==============================] - 29s 37ms/step - loss: 0.3621 - accura
Epoch 9/20
782/782 [==============================] - 28s 36ms/step - loss: 0.3513 - accura
Epoch 10/20
```

```
782/782 [==============================] - 28s 36ms/step - loss: 0.3393 - accura
Epoch 11/20
782/782 [==============================] - 28s 36ms/step - loss: 0.3271 - accura
Epoch 12/20
782/782 [==============================] - 28s 36ms/step - loss: 0.3163 - accura
Epoch 13/20
782/782 [==============================] - 29s 37ms/step - loss: 0.3029 - accura
Epoch 14/20
782/782 [==============================] - 28s 36ms/step - loss: 0.2969 - accura
Epoch 15/20
782/782 [==============================] - 29s 36ms/step - loss: 0.2831 - accura
Epoch 16/20
782/782 [==============================] - 29s 37ms/step - loss: 0.2760 - accura
Epoch 17/20
782/782 [==============================] - 28s 36ms/step - loss: 0.2648 - accura
Epoch 18/20
782/782 [==============================] - 28s 36ms/step - loss: 0.2575 - accura
Epoch 19/20
782/782 [==============================] - 28s 36ms/step - loss: 0.2511 - accura
Epoch 20/20
782/782 [==============================] - 28s 36ms/step - loss: 0.2413 - accura
```

```python
1 from sklearn import metrics
2 from sklearn.metrics import precision_score, recall_score
3
4 predicted_2 = model_2.predict(test_padded)
5
6 pred_2 = np.zeros(len(predicted_2))
7 for i, score in enumerate(predicted_2):
8     if score > 0.5:
9         pred_2[i] = 1
10     else:
11         pred_2[i] = 0
12
13 print(metrics.classification_report(test_labels, pred_2))
14 print("Precision = ",precision_score(test_labels, pred_2))
15 print("Recall = ", recall_score(test_labels, pred_2))
```

```
              precision    recall  f1-score   support

           0       0.86      0.82      0.84     12500
           1       0.83      0.87      0.85     12500

    accuracy                           0.84     25000
   macro avg       0.85      0.84      0.84     25000
weighted avg       0.85      0.84      0.84     25000

Precision =  0.8279381521821921
Recall =  0.8696
```
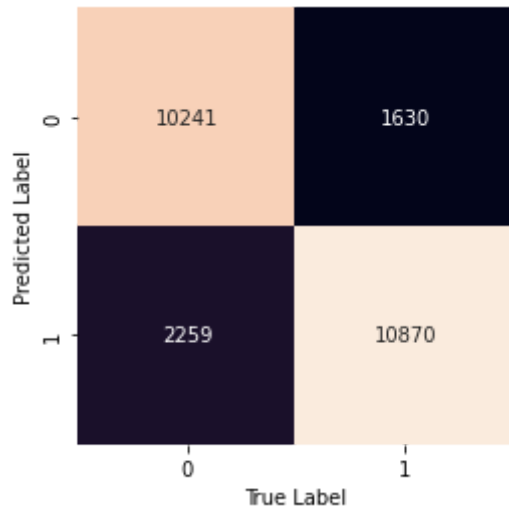
```python
1 from sklearn.metrics import confusion_matrix
2 import seaborn as sns
```

```
3 import matplotlib.pyplot as plt
4
5 mat = confusion_matrix(test_labels, pred_2)
6 sns.heatmap(mat.T, square = True, annot = True, fmt = 'd', cbar = False)
7 plt.xlabel("True Label")
8 plt.ylabel("Predicted Label")
```

Text(91.68, 0.5, 'Predicted Label')



● 0s    completed at 7:50 AM                                                    ● ✕