



65070501037

Paweeorn Soratyathorn

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
```

```
1 df = pd.read_csv('IMDB Dataset.csv')
2 df.head()
```





|   | review  | sentiment |   |
|---|---|-----------|---|
| 0 | One of the other reviewers has mentioned that ... | positive  |  |
| 1 | A wonderful little production. <br /><br />The... | positive  |   |
| 2 | I thought this was a wonderful way to spend ti... | positive  |   |
| 3 | Basically there's a family where a little boy ... | negative  |   |
| 4 | Petter Mattei's "Love in the Time of Money" is... | positive  |   |

Next steps:

[Generate code with df](#)[View recommended plots](#)[New interactive sheet](#)


```
1 pd.crosstab(df.sentiment, columns="Count").sort_values(by='Count', ascending=False)
```



|           | col_0 | Count |   |
|-----------|-------|-------|---|
| sentiment |       |       |  |
| negative  |       | 25000 |   |
| positive  |       | 25000 |   |

## ✓ Data cleansing

```
1 import nltk
2 from nltk.corpus import stopwords
3 import re
4
5 nltk.download('stopwords')
6
7 REPLACE_BY_SPLACE_RE = re.compile(r'[/(){}\\[\]\|@,;]')
8 BAD_SYMBOLS_RE = re.compile(r'^\0-9a-z #+_)')
9 STOPWORDS = set(stopwords.words('english'))
```




```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
```

```
1 def clean_text(text):
2     text = text.lower()
3     text = REPLACE_BY_SPLACE_RE.sub(' ', text)
4     text = BAD_SYMBOLS_RE.sub('', text)
5     text = text.replace('x', '')
6     text = ' '.join(word for word in text.split() if word not in STOPWORDS)
7     return text
```

```
1 df['review'] = df['review'].apply(clean_text)
2 df['review'] = df['review'].str.replace(r'\d+', '')
```

```
1 def print_plot(index):
2     example = df[df.index == index][['review', 'sentiment']].values[0]
3     if len(example) > 0:
4         print(example[0])
5         print('Review: ', example[1])
6
7 print_plot(10)
```



```
phil alien one quirky films humour based around oddness everything rather actual punchlinesbr br first odd pretty funny movie progr
Review: negative
```

## Preprocessing

### Train Tokenizer

```
1 from tensorflow.keras.preprocessing.text import Tokenizer
2
3 MAX_NB_WORDS = 50000
4
5 MAX_SEQUENCE_LENGTH = 250
6
7 EMBEDDING_DIM = 100
8
9 tokenizer = Tokenizer(num_words = MAX_NB_WORDS,
10                       filters = '!"#$%&()*+,-./:;<=>?@[\\]^_`{|}~',
11                       lower=True)
12
13 tokenizer.fit_on_texts(df['review'].values)
14 word_index = tokenizer.word_index
15 print('Found %s unique tokens.'%len(word_index))
```

Found 165306 unique tokens.

### Pad sequence

```
1 from tensorflow.keras.preprocessing.sequence import pad_sequences
2
3 X = tokenizer.texts_to_sequences(df['review'].values)[:2500]
4 X = pad_sequences(X, maxlen = MAX_SEQUENCE_LENGTH)
5 print('Shape of data tensor:', X.shape)
```

Shape of data tensor: (2500, 250)

```
1 print(df['review'].values[0], '\n')
2 X[0]
```

one reviewers mentioned watching 1 oz episode youll hooked right exactly happened mebr br first thing struck oz brutality unflinching

```
array([[ 0,  0,  0,  0,  0,  0,  0,  0,  0,
         0,  0,  0,  0,  0,  0,  0,  0,  0,
         0,  0,  0,  0,  0,  0,  0,  0,  0,
         0,  0,  0,  0,  0,  0,  0,  0,  0,
         0,  0,  0,  0,  0,  0,  0,  0,  0,
         0,  0,  0,  0,  0,  0,  0,  0,  0,
         0,  0,  0,  0,  0,  0,  0,  0,  0,
         0,  0,  0,  0,  0,  0,  0,  0,  0,
         0,  0,  0,  0,  4, 1834, 950, 57, 233,
        3193, 288, 353, 3079, 110, 492, 480, 2107, 1,
         20, 58, 3138, 3193, 5451, 15271, 51, 461, 182,
         110, 560, 53, 1585, 42, 8154, 5657, 11761, 42,
        2394, 5953, 5452, 1337, 264, 461, 3267, 249, 239,
       23365, 1, 364, 3193, 11400, 237, 15905, 6763, 2415,
         947, 2521, 1257, 25213, 425, 4483, 2409, 1081, 6991,
        2860, 12894, 302, 17412, 214, 4942, 3559, 425, 241,
        8294, 40799, 15272, 5061, 7725, 2315, 18295, 224, 9040,
        7356, 13122, 8621, 34707, 35, 128, 5513, 1, 8,
         47, 171, 1191, 42, 557, 95, 163, 159, 439,
        2874, 706, 86, 1150, 4228, 2379, 984, 706, 1295,
         706, 60, 869, 89, 20, 288, 44, 106, 3138,
        1463, 2090, 293, 47, 1437, 178, 1356, 1138, 3193,
         92, 10168, 214, 1973, 1976, 461, 461, 7856, 6992,
        4842, 14080, 2861, 32394, 6934, 14080, 384, 515, 15,
         144, 14, 9815, 639, 703, 6934, 551, 1081, 20459,
         557, 440, 814, 1880, 1081, 448, 57, 3193, 102,
        308, 3653, 3161, 15, 1090, 3906, 394], dtype=int32)
```

```
1 Y = pd.get_dummies(df['sentiment']).values[:2500]
2 print('Shape of label tensor:', Y.shape)
```

Shape of label tensor: (2500, 2)

```
1 Y
```

```
array([[False,  True],
       [False,  True],
```

```
[False, True],
...,
[False, True],
[ True, False],
[False, True]])
```

## ✓ Training Model

```
1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.3, random_state=42)
4
5 print(X_train.shape, Y_train.shape)
6 print(X_test.shape, Y_test.shape)
```

```
→ (1750, 250) (1750, 2)
   (750, 250) (750, 2)
```

```
1 from keras import Input
2 from keras.models import Sequential
3 from keras.layers import Embedding, SpatialDropout1D, LSTM, Dense
4 from keras.callbacks import EarlyStopping
5
6 model = Sequential()
7 model.add(Input(shape=(X.shape[1], )))
8 model.add(Embedding(MAX_NB_WORDS, EMBEDDING_DIM))
9 model.add(SpatialDropout1D(0.2))
10 model.add(LSTM(100, dropout=0.2, recurrent_dropout=0.2))
11 model.add(Dense(2, activation='softmax'))
12 model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
13
14 epochs=3
15 batch_size=64
16
17 history = model.fit(X_train, Y_train,
18                     epochs=epochs,
19                     batch_size=batch_size,
20                     validation_split=0.1,
21                     callbacks=[EarlyStopping(monitor='val_loss', patience=3, min_delta=0.0001)])
```

```
→ Epoch 1/3
25/25 ————— 16s 504ms/step - accuracy: 0.5439 - loss: 0.6913 - val_accuracy: 0.5943 - val_loss: 0.6696
Epoch 2/3
25/25 ————— 25s 682ms/step - accuracy: 0.7983 - loss: 0.5829 - val_accuracy: 0.7429 - val_loss: 0.5268
Epoch 3/3
25/25 ————— 16s 480ms/step - accuracy: 0.9142 - loss: 0.2584 - val_accuracy: 0.8057 - val_loss: 0.4707
```

```
1 model.summary()
```

```
→ Model: "sequential_1"
```

| Layer (type)                           | Output Shape     | Param #   |
|--|------------------|-----------|
| embedding_1 (Embedding)                | (None, 250, 100) | 5,000,000 |
| spatial_dropout1d_1 (SpatialDropout1D) | (None, 250, 100) | 0         |
| lstm_1 (LSTM)                          | (None, 100)      | 80,400    |
| dense_1 (Dense)                        | (None, 2)        | 202       |

```
Total params: 15,241,808 (58.14 MB)
Trainable params: 5,080,602 (19.38 MB)
Non-trainable params: 0 (0.00 B)
Optimizer params: 10,161,206 (38.76 MB)
```

## ✓ Evaluation

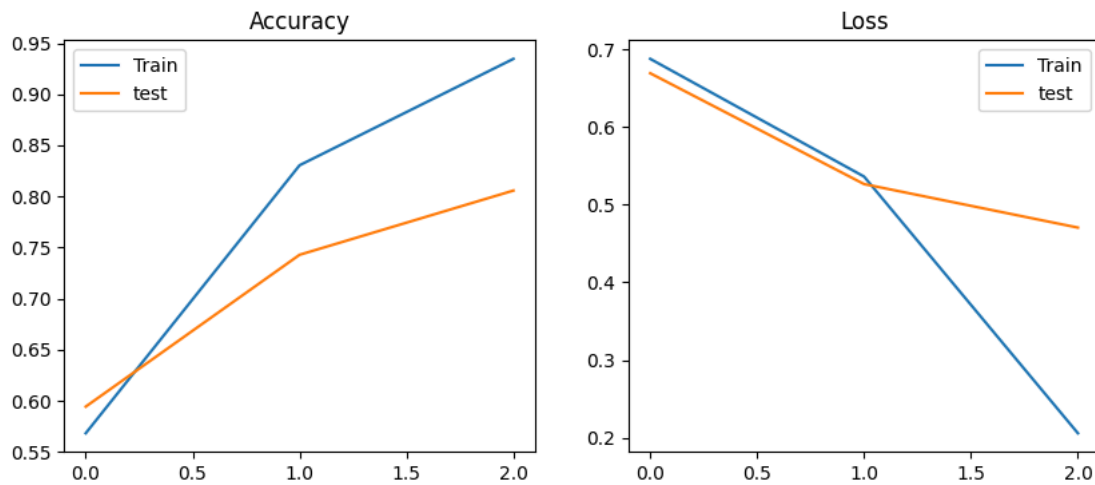
```
1 acc = model.evaluate(X_test, Y_test)
2 print('Test set \n\tLoss: {:.3f}\n\tAccuracy: {:.3f}'.format(acc[0], acc[1]))
```

```
→ 24/24 ————— 1s 55ms/step - accuracy: 0.8079 - loss: 0.4661
Test set
Loss: 0.476
Accuracy: 0.801
```

```

1 def plot_history(history):
2     plt.figure(figsize=(10, 4))
3
4     plt.subplot(1, 2, 1)
5     plt.title('Accuracy')
6     plt.plot(history.history['accuracy'], label='Train')
7     plt.plot(history.history['val_accuracy'], label='test')
8     plt.legend()
9
10    plt.subplot(1, 2, 2)
11    plt.title('Loss')
12    plt.plot(history.history['loss'], label='Train')
13    plt.plot(history.history['val_loss'], label='test')
14    plt.legend()
15
16
17 plot_history(history)

```



### Confusion matrix

```

1 labels = pd.get_dummies(df['sentiment']).columns
2
3 from sklearn.metrics import confusion_matrix
4 y_pred = model.predict(X_test)
5
6 pd.DataFrame(confusion_matrix(Y_test.argmax(axis=1),
7                               y_pred.argmax(axis=1)),
8              index=labels, columns=labels)

```



24/24 ————— 6s 224ms/step

|          | negative | positive |
|----------|----------|----------|
| negative | 295      | 82       |
| positive | 67       | 306      |

### Classification Report

```

1 from sklearn.metrics import classification_report
2
3 y_pred = model.predict(X_test)
4 print(classification_report(y_true=Y_test.argmax(axis=1),
5                             y_pred=y_pred.argmax(axis=1)))

```



24/24 ————— 2s 87ms/step

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.81      | 0.78   | 0.80     | 377     |
| 1            | 0.79      | 0.82   | 0.80     | 373     |
| accuracy     |           |        | 0.80     | 750     |
| macro avg    | 0.80      | 0.80   | 0.80     | 750     |
| weighted avg | 0.80      | 0.80   | 0.80     | 750     |

## ✓ Test using new review

```
1 new_review = ['''Visually beautiful but unfortunately the dialogues are poorly written and superficial,  
2 Galedriel who has a lot of screen time and massive importance unfortunately is portrayed very badly makes me questio  
3 how could they allow such performance, not to insult the actress at all but the performance was horrible almost  
4 everyone in the show was better performing than her. Overall fun to watch but very superficial and Shallow.  
5 The plot is yet to reveal itself but so far nothing interesting. Good to mention the good portrayal of elrond.  
6 elendil. Music is also good and appropriate. The dark theme is a bit overlooked.''']
```

```
1 seq = tokenizer.texts_to_sequences(new_review)  
2 padded = pad_sequences(seq, maxlen=MAX_SEQUENCE_LENGTH)  
3 pred = model.predict(padded)  
4 print(pred, labels[np.argmax(pred)])
```

→ 1/1 ————— 0s 151ms/step  
[[0.6472341 0.3527659]] negative

1 Start coding or [generate](#) with AI.