

ML-NLP-sentiment-analysis-of-movie-reviews

March 4, 2023

1 Sentiment Analysis of IMDB Movie Reviews (Part 2: Machine Learning models)

Problem Statement:

In this, we have to predict the number of positive and negative reviews based on sentiments by using different classification models.

Side note for profs - I have done NLP analysis before, so this work is using as a [this template](#) as a baseline. - Feel free to ignore the BERT model code. It works but taking too long to train and often crashes the kernel. Hence, I switch to use HuggingFace training package in Sentiment Analysis of IMDB Movie Reviews (Part 3: BERT model)

1.1 Data Processing

Import necessary libraries

```
[ ]: #Load the libraries
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import nltk
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.preprocessing import LabelBinarizer
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
from wordcloud import WordCloud, STOPWORDS
from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize, sent_tokenize
from bs4 import BeautifulSoup
import spacy
import re, string, unicodedata
from nltk.tokenize.toktok import ToktokTokenizer
from nltk.stem import LancasterStemmer, WordNetLemmatizer
from sklearn.linear_model import LogisticRegression, SGDClassifier
from sklearn.naive_bayes import MultinomialNB
from sklearn.svm import SVC
```

```
# from textblob import TextBlob
# from textblob import Word
from sklearn.metrics import
    classification_report, confusion_matrix, accuracy_score
```

Import the training dataset

```
[ ]: #importing the training data
imdb_data=pd.read_csv('IMDB Dataset.csv')
print(imdb_data.shape)
imdb_data.head(10)
```

(50000, 2)

```
[ ]:                                     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
5  Probably my all-time favorite movie, a story o... positive
6  I sure would like to see a resurrection of a u... positive
7  This show was an amazing, fresh & innovative i... negative
8  Encouraged by the positive comments about this... negative
9  If you like original gut wrenching laughter yo... positive
```

Exploratory data analysis

```
[ ]: #Summary of the dataset
imdb_data.describe()
```

```
[ ]:                                     review sentiment
count                                     50000      50000
unique                                     49582         2
top      Loved today's show!!! It was a variety and not... positive
freq                                     5      25000
```

Sentiment count

```
[ ]: #sentiment count
imdb_data['sentiment'].value_counts()
```

```
[ ]: positive    25000
negative    25000
Name: sentiment, dtype: int64
```

We can see that the dataset is balanced.

Splitting the training dataset

```
[ ]: #split the dataset
#train dataset
train_reviews=imdb_data.review[:40000]
train_sentiments=imdb_data.sentiment[:40000]
#test dataset
test_reviews=imdb_data.review[40000:]
test_sentiments=imdb_data.sentiment[40000:]
print(train_reviews.shape,train_sentiments.shape)
print(test_reviews.shape,test_sentiments.shape)
```

(40000,) (40000,)

(10000,) (10000,)

Removing html strips and noise text

```
[ ]: #Removing the html strips
def strip_html(text):
    soup = BeautifulSoup(text, "html.parser")
    return soup.get_text()

#Removing the square brackets
def remove_between_square_brackets(text):
    return re.sub('\[[^\]]*\]', '', text)

#Removing the noisy text
def denoise_text(text):
    text = strip_html(text)
    text = remove_between_square_brackets(text)
    return text

#Apply function on review column
imdb_data['review']=imdb_data['review'].apply(denoise_text)
```

/Users/swimmingcircle/opt/anaconda3/lib/python3.9/site-packages/bs4/__init__.py:435: MarkupResemblesLocatorWarning: The input looks more like a filename than markup. You may want to open this file and pass the filehandle into BeautifulSoup.

warnings.warn(

Removing special characters

```
[ ]: #Define function for removing special characters
def remove_special_characters(text, remove_digits=True):
    pattern=r'[^a-zA-z0-9\s]'
    text=re.sub(pattern,'',text)
    return text

#Apply function on review column
imdb_data['review']=imdb_data['review'].apply(remove_special_characters)
```

Text stemming

```
[ ]: #Stemming the text
def simple_stemmer(text):
    ps=nlTK.porter.PorterStemmer()
    text= ' '.join([ps.stem(word) for word in text.split()])
    return text
#Apply function on review column
imdb_data['review']=imdb_data['review'].apply(simple_stemmer)
```

Removing stopwords

```
[ ]: #Tokenization of text
tokenizer=ToktokTokenizer()
#Setting English stopwords
stopword_list=nlTK.corpus.stopwords.words('english')
```

```
[ ]: #set stopwords to english
stop=set(stopwords.words('english'))
print(stop)

#removing the stopwords
def remove_stopwords(text, is_lower_case=False):
    tokens = tokenizer.tokenize(text)
    tokens = [token.strip() for token in tokens]
    if is_lower_case:
        filtered_tokens = [token for token in tokens if token not in
↪stopword_list]
    else:
        filtered_tokens = [token for token in tokens if token.lower() not in
↪stopword_list]
    filtered_text = ' '.join(filtered_tokens)
    return filtered_text
#Apply function on review column
imdb_data['review']=imdb_data['review'].apply(remove_stopwords)
```

```
{'until', 'all', 'didn', 'when', 'before', "mightn't", "you're", 'herself',
'that', "you'd", 'doing', 't', 'while', 'needn', 'ma', "isn't", 'their', 'has',
'him', 'mustn', 'an', 'd', 'in', 'same', 'then', 'being', 'both', 'itself',
'very', "shouldn't", 'hasn', "wouldn't", 'having', 'some', "weren't",
'ourselves', 'ours', 'against', 's', 'through', "haven't", "wasn't", "doesn't",
'ain', 'been', "she's", 'a', 'why', 'because', 'down', 'further', 'hadn',
'wasn', 'theirs', 'its', 'themselves', 'were', 'between', "you've", 'from',
'doesn', 've', 'should', 'so', "mustn't", 'himself', 'nor', 'more', 'couldn',
'had', "didn't", 'can', 'them', 'shouldn', 'for', 'o', 'don', 'which', 'to',
'shan', 'again', "it's", 'isn', 'of', 'yourselves', "shan't", 'no', 'yours',
'under', 'you', 'up', 'are', 'i', 'on', 'aren', 'and', 'about', 'y', 'll',
"needn't", 'once', 'will', 'at', 'her', 'by', 'wouldn', 're', "won't", 'other',
'there', 'what', 'she', 'or', 'haven', "hasn't", 'the', 'me', 'if', 'mightn',
'they', 'but', 'now', "hadn't", 'who', 'won', 'off', 'each', 'too', 'he', 'his',
```

```
"don't", 'it', 'how', 'does', 'is', 'few', 'be', 'below', 'not', 'into', 'am',
'as', 'over', "should've", 'most', 'hers', 'm', 'weren', 'after', 'do', 'our',
'only', 'these', 'those', 'own', 'have', 'did', 'any', "couldn't", "that'll",
'just', 'than', 'myself', 'yourself', 'whom', 'your', 'my', 'with', "you'll",
'we', 'this', 'where', 'such', "aren't", 'above', 'was', 'here', 'out',
'during'}
```

Normalized train reviews

```
[ ]: #normalized train reviews
norm_train_reviews=imdb_data.review[:40000]
norm_train_reviews[0]
```

```
[ ]: 'one review ha mention watch 1 oz episod youll hook right thi exactli happen
meth first thing struck oz wa brutal unflinch scene violenc set right word go
trust thi show faint heart timid thi show pull punch regard drug sex violenc
hardcor classic use wordit call oz nicknam given oswald maximum secur state
penitentari focus mainli emerald citi experiment section prison cell glass front
face inward privaci high agenda em citi home manyaryan muslim gangsta latino
christian italian irish moreso scuffl death stare dodgi deal shadi agreement
never far awayi would say main appeal show due fact goe show wouldnt dare forget
pretti pictur paint mainstream audienc forget charm forget romanceoz doesnt mess
around first episod ever saw struck nasti wa surreal couldnt say wa readi watch
develop tast oz got accustom high level graphic violenc violenc injustic crook
guard wholl sold nickel inmat wholl kill order get away well manner middl class
inmat turn prison bitch due lack street skill prison experi watch oz may becom
comfort uncomfort viewingthat get touch darker side'
```

Normalized test reviews

```
[ ]: #Normalized test reviews
norm_test_reviews=imdb_data.review[40000:]
norm_test_reviews[45005]
```

```
[ ]: 'read review watch thi piec cinemat garbag took least 2 page find somebodi els
didnt think thi appallingli unfunni montag wasnt acm humour 70 inde ani era thi
isnt least funni set sketch comedi ive ever seen itll till come along half skit
alreadi done infinit better act monti python woodi allen wa say nice piec anim
last 90 second highlight thi film would still get close sum mindless
drivelridden thi wast 75 minut semin comedi onli world semin realli doe mean
semen scatolog humour onli world scat actual fece precursor joke onli mean thi
handbook comedi tit bum odd beaver niceif pubesc boy least one hand free havent
found playboy exist give break becaus wa earli 70 way sketch comedi go back
least ten year prior onli way could even forgiv thi film even made wa gunpoint
retro hardli sketch clown subtli pervert children may cut edg circl could actual
funni come realli quit sad kept go throughout entir 75 minut sheer belief may
save genuin funni skit end gave film 1 becaus wa lower scoreand onli recommend
insomniac coma patientsor perhap peopl suffer lockjawtheir jaw would final drop
open disbelief'
```

Bags of words model

It is used to convert text documents to numerical vectors or bag of words.

```
[ ]: #Count vectorizer for bag of words
cv=CountVectorizer(min_df=0,max_df=1,binary=False,ngram_range=(1,3))
#transformed train reviews
cv_train_reviews=cv.fit_transform(norm_train_reviews)
#transformed test reviews
cv_test_reviews=cv.transform(norm_test_reviews)

print('BOW_cv_train:',cv_train_reviews.shape)
print('BOW_cv_test:',cv_test_reviews.shape)
#vocab=cv.get_feature_names()-toget feature names
```

BOW_cv_train: (40000, 6209089)

BOW_cv_test: (10000, 6209089)

Term Frequency-Inverse Document Frequency model (TFIDF)

It is used to convert text documents to matrix of tfidf features.

```
[ ]: #Tfidf vectorizer
tv=TfidfVectorizer(min_df=0,max_df=1,use_idf=True,ngram_range=(1,3))
#transformed train reviews
tv_train_reviews=tv.fit_transform(norm_train_reviews)
#transformed test reviews
tv_test_reviews=tv.transform(norm_test_reviews)
print('Tfidf_train:',tv_train_reviews.shape)
print('Tfidf_test:',tv_test_reviews.shape)
```

Tfidf_train: (40000, 6209089)

Tfidf_test: (10000, 6209089)

1.2 Logistic Regression

```
[ ]: #training the model
lr=LogisticRegression(penalty='l2',max_iter=500,C=1,random_state=42)
#Fitting the model for Bag of words
lr_bow=lr.fit(cv_train_reviews,train_sentiments)
print(lr_bow)
#Fitting the model for tfidf features
lr_tfidf=lr.fit(tv_train_reviews,train_sentiments)
print(lr_tfidf)
```

LogisticRegression(C=1, max_iter=500, random_state=42)

LogisticRegression(C=1, max_iter=500, random_state=42)

```
[ ]: #Predicting the model for bag of words
lr_bow_predict=lr.predict(cv_test_reviews)
```

```
print(lr_bow_predict)
##Predicting the model for tfidf features
lr_tfidf_predict=lr.predict(tv_test_reviews)
print(lr_tfidf_predict)
```

```
['negative' 'negative' 'negative' ... 'negative' 'positive' 'positive']
['negative' 'negative' 'negative' ... 'negative' 'positive' 'positive']
```

```
[ ]: #Accuracy score for bag of words
lr_bow_score=accuracy_score(test_sentiments,lr_bow_predict)
print("lr_bow_score :",lr_bow_score)
#Accuracy score for tfidf features
lr_tfidf_score=accuracy_score(test_sentiments,lr_tfidf_predict)
print("lr_tfidf_score :",lr_tfidf_score)
```

```
lr_bow_score : 0.7512
lr_tfidf_score : 0.75
```

```
[ ]: #Classification report for bag of words
lr_bow_report=classification_report(test_sentiments,lr_bow_predict,target_names=['Positive','Negative'])
print(lr_bow_report)

#Classification report for tfidf features
lr_tfidf_report=classification_report(test_sentiments,lr_tfidf_predict,target_names=['Positive','Negative'])
print(lr_tfidf_report)
```

	precision	recall	f1-score	support
Positive	0.75	0.75	0.75	4993
Negative	0.75	0.75	0.75	5007
accuracy			0.75	10000
macro avg	0.75	0.75	0.75	10000
weighted avg	0.75	0.75	0.75	10000

	precision	recall	f1-score	support
Positive	0.74	0.77	0.75	4993
Negative	0.76	0.73	0.75	5007
accuracy			0.75	10000
macro avg	0.75	0.75	0.75	10000
weighted avg	0.75	0.75	0.75	10000

Logistic analysis shows around 75% accuracy for both bags of words and tf-idf data which is quite impressive!

1.3 SDGC: Stochastic gradient descent or Linear support vector machines for bag of words and tfidf features

SGDClassifier is a linear classifier in scikit-learn that uses stochastic gradient descent (SGD) as the optimization algorithm. It is a type of online learning algorithm that can handle large-scale datasets efficiently, by processing one instance at a time and updating the model parameters incrementally.

SGDClassifier can be used for binary classification, multi-class classification, and regression tasks. It supports a variety of loss functions, such as hinge loss (for linear SVM), log loss (for logistic regression), and squared loss (for linear regression). It also supports various regularization methods, such as L1 and L2 regularization, to prevent overfitting.

SGDClassifier can be a good choice for large datasets or streaming data, where batch learning algorithms may not be suitable due to memory constraints or processing time. However, it may require more hyperparameter tuning and preprocessing than other classifiers, as it is more sensitive to the scaling and distribution of the features.

```
[ ]: #training the linear svm
svm=SGDClassifier(loss='hinge',max_iter=500,random_state=42)
#fitting the svm for bag of words
svm_bow=svm.fit(cv_train_reviews,train_sentiments)
print(svm_bow)
#fitting the svm for tfidf features
svm_tfidf=svm.fit(tv_train_reviews,train_sentiments)
print(svm_tfidf)
```

```
SGDClassifier(max_iter=500, random_state=42)
SGDClassifier(max_iter=500, random_state=42)
```

```
[ ]: #Predicting the model for bag of words
svm_bow_predict=svm.predict(cv_test_reviews)
print(svm_bow_predict)
#Predicting the model for tfidf features
svm_tfidf_predict=svm.predict(tv_test_reviews)
print(svm_tfidf_predict)
```

```
['positive' 'positive' 'negative' ... 'positive' 'positive' 'positive']
['positive' 'positive' 'positive' ... 'positive' 'positive' 'positive']
```

```
[ ]: #Accuracy score for bag of words
svm_bow_score=accuracy_score(test_sentiments,svm_bow_predict)
print("svm_bow_score :",svm_bow_score)
#Accuracy score for tfidf features
svm_tfidf_score=accuracy_score(test_sentiments,svm_tfidf_predict)
print("svm_tfidf_score :",svm_tfidf_score)
```

```
svm_bow_score : 0.5829
svm_tfidf_score : 0.5112
```



```
[ ]: #Classification report for bag of words
svm_bow_report=classification_report(test_sentiments,svm_bow_predict,target_names=['Positive',
print(svm_bow_report)
#Classification report for tfidf features
svm_tfidf_report=classification_report(test_sentiments,svm_tfidf_predict,target_names=['Positi
print(svm_tfidf_report)
```

	precision	recall	f1-score	support
Positive	0.94	0.18	0.30	4993
Negative	0.55	0.99	0.70	5007
accuracy			0.58	10000
macro avg	0.74	0.58	0.50	10000
weighted avg	0.74	0.58	0.50	10000

	precision	recall	f1-score	support
Positive	1.00	0.02	0.04	4993
Negative	0.51	1.00	0.67	5007
accuracy			0.51	10000
macro avg	0.75	0.51	0.36	10000
weighted avg	0.75	0.51	0.36	10000

Since random guess will have 50% accuracy, we can conclude that SDGC (bags of word accuracy: 58%, tf-idf accuracy: 51%) isn't a good model to implement.

1.4 XGboost

```
[ ]: train_sentiments_label = train_sentiments.apply(lambda x: 1 if x == 'positive'␣
↪else 0)
test_sentiments_label = test_sentiments.apply(lambda x: 1 if x == 'positive'␣
↪else 0)
train_sentiments_label.value_counts(), test_sentiments_label.value_counts()
```

```
[ ]: (0    20007
      1    19993
      Name: sentiment, dtype: int64,
      1     5007
      0     4993
      Name: sentiment, dtype: int64)
```

```
[ ]: import xgboost as xgb
```

```
[ ]: cv_classifier = xgb.XGBClassifier(max_depth = 7, eta = 0.9, objective= 'binary:
↪hinge', n_estimators = 200,
```

```

                                use_label_encoder=False, eval_metric = 'auc')
tv_classifier = xgb.XGBClassifier(max_depth = 10, eta = 0.2, objective= 'binary:
↳hinge', n_estimators = 200,
                                use_label_encoder=False, eval_metric = 'auc')

cv_bow = cv_classifier.fit(cv_train_reviews, train_sentiments_label)
cv_tfidf = tv_classifier.fit(tv_train_reviews, train_sentiments_label)

```

```

[ ]: xgb_bow_pred = cv_classifier.predict(cv_test_reviews)
xgb_tfidf_pred = tv_classifier.predict(tv_test_reviews)

# evaluate predictions
cv_score = classification_report(xgb_bow_pred, test_sentiments_label,
↳target_names=['Positive', 'Negative'])
print(cv_score)

tv_score = classification_report(xgb_tfidf_pred, test_sentiments_label,
↳target_names=['Positive', 'Negative'])
print(tv_score)

```

	precision	recall	f1-score	support
Positive	1.00	0.50	0.67	9972
Negative	0.00	0.61	0.01	28
accuracy			0.50	10000
macro avg	0.50	0.55	0.34	10000
weighted avg	1.00	0.50	0.66	10000

	precision	recall	f1-score	support
Positive	0.00	0.00	0.00	0
Negative	1.00	0.50	0.67	10000
accuracy			0.50	10000
macro avg	0.50	0.25	0.33	10000
weighted avg	1.00	0.50	0.67	10000

/Users/swimmingcircle/opt/anaconda3/lib/python3.9/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use `zero_division` parameter to control this behavior.

```

_warn_prf(average, modifier, msg_start, len(result))
/Users/swimmingcircle/opt/anaconda3/lib/python3.9/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use `zero_division` parameter to control this behavior.

```

```

_warn_prf(average, modifier, msg_start, len(result))
/Users/swimmingcircle/opt/anaconda3/lib/python3.9/site-
packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Recall
and F-score are ill-defined and being set to 0.0 in labels with no true samples.
Use `zero_division` parameter to control this behavior.
_warn_prf(average, modifier, msg_start, len(result))

```

XGboost doesn't seem to be good at handling language transformed texts. It is very sensitive to parameter changes. It often predict the all the data into one class or another when changing the parameters. We conclude that it isn't a good model for our sentiment classification.

1.5 BERT model

1.5.1 What's special about BERT?

- Context-free models: generate a single word embedding representation for each word in the vocabulary, such as word2vec or GloVe. For example, the word “bank” would have the same representation in “bank deposit” and in “riverbank”
- Contextual models instead generate a representation of each word that is based on the other words in the sentence, such as BERT.

1.5.2 Understand how BERT works

1. Token embeddings: A [CLS] token is added to the input word tokens at the beginning of the first sentence and a [SEP] token is inserted at the end of each sentence.
 2. Segment embeddings: A marker indicating Sentence A or Sentence B is added to each token. **This allows the encoder to distinguish between sentences.**
 3. Positional embeddings: A positional embedding is added to each token to indicate its position in the sentence.
1. Masked LM (MLM) The idea here is “simple”: Randomly mask out 15% of the words in the input — replacing them with a [MASK] token. Loss function considers only the prediction of the masked tokens and ignores the prediction of the non-masked ones.
 2. Next Sentence Prediction (NSP) In order to understand relationship between two sentences, BERT training process also uses next sentence prediction, BERT separates sentences with a special [SEP] token. During training the model is fed with two input sentences at a time such that:
 - 50% of the time the second sentence comes after the first one.
 - 50% of the time it is a a random sentence from the full corpus.

Example: predict if the next sentence is random or not

Important note: BERT does not try to predict the next word in the sentence!!

1.5.3 Tokenizer for BERT

BERT uses what is called a WordPiece tokenizer. It works by splitting words either into the full forms (e.g., one word becomes one token) or into word pieces — where one word can be broken into multiple tokens.

Word	Token(s)
surf	['surf']
surfing	['surf', '##ing']
surfboarding	['surf', '##board', '##ing']
surfboard	['surf', '##board']
snowboard	['snow', '##board']
snowboarding	['snow', '##board', '##ing']
snow	['snow']
snowing	['snow', '##ing']

By splitting words into word pieces, we have already identified that the words “surfboard” and “snowboard” share meaning through the wordpiece “##board” We have done this without even encoding our tokens or processing them in any way through BERT.

1.5.4 BERT model choice

BERT model we choose **DistilBERT** vs BERT - DistilBERT is a small, fast, cheap and light Transformer model trained by distilling BERT base. It has 40% less parameters than bert-base-uncased, runs 60% faster while preserving over 95% of BERT’s performances as measured on the GLUE language understanding benchmark.

BERT-base vs BERT-large: BERT-based - BERT-Base: 12-layer, 768-hidden-nodes, 12-attention-heads, 110M parameters - BERT-Large: 24-layer, 1024-hidden-nodes, 16-attention-heads, 340M parameters

BERT-based-case vs **BERT-base-uncased**: - We don’t differentiate between cased and uncased data (english vs English)

1.5.5 BERT input

Input IDs – The input ids are often the only required parameters to be passed to the model as input. Token indices, numerical representations of tokens building the sequences that will be used as input by the model.

Attention mask – Attention Mask is used to avoid performing attention on padding token indices. Mask value can be either 0 or 1, 1 for tokens that are NOT MASKED, 0 for MASKED tokens.

Token type ids – It is used in use cases like sequence classification or question answering. As these require two different sequences to be encoded in the same input IDs. Special tokens, such as the classifier[CLS] and separator[SEP] tokens are used to separate the sequences.

Note: Padding is a special form of masking where the masked steps are at the start or the end of a sequence. Padding comes from the need to encode sequence data into contiguous batches: in order to make all sequences in a batch fit a given standard length, it is necessary to pad or truncate some sequences

1.5.6 BERT tokens

CLS: The [CLS] token, short for “classification,” is a special token used in BERT to represent the entire input sequence for classification tasks.

When training a classification model using BERT, the [CLS] token is added to the beginning of the input sequence, and the final hidden state corresponding to this token is used as the input to a classifier. This allows the model to make a prediction for the entire input sequence.

SEP: The [SEP] token, short for “separator,” is used to separate two different segments of a sentence or document.

In BERT, the [SEP] token is used to separate the two segments when performing tasks like question answering or natural language inference, where the model needs to understand the relationship between two different segments of text.

MASK: [MASK] is used during pre-training to randomly mask some of the input tokens, forcing the model to learn to predict the masked tokens based on the surrounding context.

1.5.7 Understanding the parameters

`max_length` is a parameter used to define the maximum length of an input sequence.

`pad_to_max_length` is a Boolean parameter used to indicate whether sequences shorter than the `max_length` should be padded with a special token, usually [PAD], to make them the same length as the longest sequence in the batch.

`return_tensors` parameter specifies that we want the encoded data to be returned as TensorFlow tensor

`attention_mask`: 1 indicates a value that should be attended to, while 0 indicates a padded value.

Example:

```
sequence_a = "This is a short sequence."          sequence_b = "This is a rather long
sequence. It is at least longer than the sequence A."  len(encoded_sequence_a),
len(encoded_sequence_b)
```

(8, 19)

```
padded_sequences = tokenizer([sequence_a, sequence_b], padding=True)
padded_sequences["input_ids"]
```

```
[[101, 1188, 1110, 170, 1603, 4954, 119, 102, 0, 0, 0, 0, 0, 0, 0, 0], [101, 1188, 1110, 170,
1897, 1263, 4954, 119, 1135, 1110, 1120, 1655, 2039, 1190, 1103, 4954, 138, 119, 102]]
```

```
padded_sequences["attention_mask"]
```

```
[[1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0], [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]]
```

```
[ ]: #test with a smaller dataset
imdb_data = imdb_data[:300]
```

```
[ ]: import transformers
from tokenizers import BertWordPieceTokenizer
# First load the real tokenizer
tokenizer = transformers.BertTokenizerFast.from_pretrained('distilbert-base-uncased', lower = True)
# Save the loaded tokenizer locally
```

```
tokenizer.save_pretrained('.')
# Reload it with the huggingface tokenizers library
```

```
[ ]: Tokenizer(vocabulary_size=30522, model=BertWordPiece, unk_token=[UNK],
sep_token=[SEP], cls_token=[CLS], pad_token=[PAD], mask_token=[MASK],
clean_text=True, handle_chinese_chars=True, strip_accents=None, lowercase=True,
wordpieces_prefix=##)
```

Resources - [BERT Explained: A Complete Guide with Theory and Tutorial](#)

```
[ ]: # Encode the training data
encoded_train_data = tokenizer(train_reviews.values.tolist(), padding=True,
↪truncation=True, return_tensors='pt')
```

```
[ ]: from torch.utils.data import TensorDataset, DataLoader
import torch
from tqdm import tqdm
from transformers import BertModel

# Load the pre-trained BERT model
bert_model = BertModel.from_pretrained('bert-base-uncased')

# Convert data to PyTorch tensors
input_ids = torch.tensor(encoded_train_data['input_ids'])
attention_masks = torch.tensor(encoded_train_data['attention_mask'])
train_sentiments = train_sentiments.apply(lambda x: 1 if x == 'positive' else 0)
labels = torch.tensor(train_sentiments.values.tolist())

# Create a TensorDataset
dataset = TensorDataset(input_ids, attention_masks, labels)

# Define batch size
batch_size = 8

# Create a DataLoader
dataloader = DataLoader(
    dataset,
    batch_size=batch_size,
    shuffle=True
)

# Iterate over batches
for batch in tqdm(dataloader):
    batch_input_ids = batch[0]
    batch_attention_masks = batch[1]
    batch_labels = batch[2]
    outputs = bert_model(batch_input_ids, attention_mask=batch_attention_masks)
```

```
/var/folders/0h/xyv81g2n7sj6zr0c9cw30gkc0000gn/T/ipykernel_8024/2238400034.py:7:
UserWarning: To copy construct from a tensor, it is recommended to use
sourceTensor.clone().detach() or
sourceTensor.clone().detach().requires_grad_(True), rather than
torch.tensor(sourceTensor).
    input_ids = torch.tensor(encoded_train_data['input_ids'])
/var/folders/0h/xyv81g2n7sj6zr0c9cw30gkc0000gn/T/ipykernel_8024/2238400034.py:8:
UserWarning: To copy construct from a tensor, it is recommended to use
sourceTensor.clone().detach() or
sourceTensor.clone().detach().requires_grad_(True), rather than
torch.tensor(sourceTensor).
    attention_masks = torch.tensor(encoded_train_data['attention_mask'])
100%|      | 38/38 [05:53<00:00,  9.30s/it]
```

```
[ ]:
```

```
[ ]:
```