```
1 !git clone https://github.com/AdvaitDeochakke/SBI_Sentiment_Analysis
```

```
Cloning into 'SBI_Sentiment_Analysis'...
remote: Enumerating objects: 39, done.
remote: Total 39 (delta 0), reused 0 (delta 0), pack-reused 39
Unpacking objects: 100% (39/39), 104.68 MiB | 10.09 MiB/s, done.
```

```python
1 import tensorflow as tf
2
3 gpus = tf.config.list_physical_devices('GPU')
4 if gpus:
5     # GPU/CUDA is available
6     print("GPU is available")
7     # You can also print the details of the available GPUs
8     for gpu in gpus:
9         print("GPU Name:", gpu.name)
10 else:
11     # No GPU/CUDA available
12     print("GPU is not available")
```

```
GPU is available
GPU Name: /physical_device:GPU:0
```

```python
1 import logging
2 logging.basicConfig(format='%(asctime)s : %(levelname)s : %(message)s', level=logging.INFO)
```

```python
1 import pandas as pd
2 df = pd.read_csv('SBI_Sentiment_Analysis/preprocessed_dataset.csv')
3 df.target.value_counts()
```

```
Negative    748869
Positive    734073
Name: target, dtype: int64
```

```python
1 df.head()
```

|   | target | Clean |
|---|--------|-------|
| 0 | Negative | upset cant updat facebook text might cri resul... |
| 1 | Negative | dive mani time ball manag save 50 rest go bound |
| 2 | Negative | whole bodi feel itchi like fire |
| 3 | Negative | behav im mad cant see |
| 4 | Negative | whole crew |

```python
1 print(df.isnull().sum())
2 df.dropna(inplace=True) # clean na
```

```
target    0
Clean     1
dtype: int64
```

```python
1 # simple train test split
2
3 from sklearn.model_selection import train_test_split
4 train_data, test_data = train_test_split(df, test_size=0.3, random_state=42)
5 print(train_data.target.value_counts())
6 print(test_data.target.value_counts())
7
```

```
Negative    524643
Positive    513415
Name: target, dtype: int64
Negative    224225
Positive    220658
Name: target, dtype: int64
```

```python
1 from gensim.models.phrases import Phraser, Phrases
2 unigrams = [_.split() for _ in df.Clean]
3 unigrams[0:5]
```

```
[['upset',
  'cant',
  'updat',
  'facebook',
  'text',
  'might',
  'cri',
  'result',
  'school',
  'today',
  'also',
  'blah'],
 ['dive',
  'mani',
  'time',
  'ball',
  'manag',
  'save',
  '50',
  'rest',
  'go',
  'bound'],
 ['whole', 'bodi', 'feel', 'itchi', 'like', 'fire'],
 ['behav', 'im', 'mad', 'cant', 'see'],
 ['whole', 'crew']]
```

```
1 phrases = Phrases(unigrams, min_count=40, progress_per=10000)
2 bigrams = Phraser(phrases)
3 sentences = bigrams[unigrams]
4 sentences[0:5]
```

```
<gensim.interfaces.TransformedCorpus at 0x7f72d5203df0>
```

```
1 # sanity check time ! im not schizo
2 # get the number of words in the vocabulary
3 from collections import defaultdict
4 wrdfreq = defaultdict(int)
5 for myphrase in sentences:
6     for oneword in myphrase:
7         wrdfreq[oneword] += 1
8 len(wrdfreq)
```

```
454515
```

```
1 # check the most frequent words
2 # we removed stopwords, so not seen 'the' is expected
3 sorted(wrdfreq, key=wrdfreq.get, reverse=True)[:10]
```

```
['im', 'go', 'get', 'day', 'work', 'love', 'good', 'like', 'today', 'time']
```

```
1 # for workers in building genshimmodel
2
3 import multiprocessing
4 from gensim.models import Word2Vec
5
6 corecount = multiprocessing.cpu_count()
7 corecount
```

```
2
```

```
1 genshinmodel = Word2Vec(
2     min_count= 30, # ignores words which appear less than 30 times in the corpora
3     window= 7, # context window size
4     vector_size= 300, # size of the vector
5     sample= 5e-5, # random downsampling of high freq words
6     alpha= 0.04, # learning rate
7     min_alpha= 0.005, # minimum rate of learning, where it will stop dropping
8     negative= 10, # negative sampling for drowning
9     workers= corecount-1
10 )
```

```
1 import time
2 t = time.time()
3
4 # monitor time and build the vocabulary
```

```
5 # mesh them up real good
6 genshinmodel.build_vocab(sentences, progress_per=10000)
7 print("\n\ntook time", (time.time()-t), "s")
```

```
    took time 12.653753519058228 s
```

```
 1 t= time.time()
 2 # and now, train it. the important part
 3 # as expected, takes a bit of time.
 4 # more epochs will result in better results as a rule of thumb
 5 # but 30 is fine.
 6 # total examples
 7 genshinmodel.train(
 8     sentences, # the corpus
 9     total_examples=genshinmodel.corpus_count, # number of sentences
10     epochs=30, # how many epochs to train for
11     report_delay=1 # progress report how often, in seconds. can honestly set it to high 10s
12     )
13 print("\n\ntook time", (time.time()-t), "s")
```

```
    took time 1021.8200531005859 s
```

```
 1 from nltk.stem import SnowballStemmer
 2 stemmer = SnowballStemmer('english') # stemmer for similarity
```

```
 1 genshinmodel.wv.most_similar(stemmer.stem('test')) # show the similarity checker
```

```
    [('exam', 0.564516544342041),
     ('math', 0.5120426416397095),
     ('studi', 0.5028963685035706),
     ('fail', 0.4861163794994354),
     ('oral_exam', 0.46020105481147766),
     ('modul', 0.45380353927612305),
     ('math_gcse', 0.452207088470459),
     ('retak', 0.4450313448905945),
     ('class', 0.43802914023399353),
     ('scienc', 0.43654054403305054)]
```

```
 1 # keras modeling.
 2 # i love the simplicity, but not having 12.1 cuda support on my system kinda hurts
 3
 4 from keras.models import Sequential
 5 from keras.layers import Activation, Dense, Dropout, Embedding, Flatten, Conv1D, MaxPooling1D, LSTM
 6 from keras import utils
 7 from keras.callbacks import ReduceLROnPlateau, EarlyStopping
```

```
 1 # tokenize the inputs
 2 # initiliaze ont he dataset
 3
 4 from keras.preprocessing.text import Tokenizer
 5 tknizer = Tokenizer()
 6 tknizer.fit_on_texts(df.Clean)
```

```
 1 # ensure they are of the same length by padding. check pad_sequences docu for various params
 2 # doing simple pre-padding here
 3
 4 from keras.utils import pad_sequences
 5
 6 x_train = pad_sequences(tknizer.texts_to_sequences(train_data.Clean), maxlen=140)
 7 x_test = pad_sequences(tknizer.texts_to_sequences(test_data.Clean), maxlen=140) # twitter absolute max size
 8
 9 print("x_train", x_train.shape)
10 print("x_test", x_test.shape)
```

```
    x_train (1038058, 140)
    x_test (444883, 140)
```

```
 1 # encode labels and transform
 2
 3 from sklearn.preprocessing import LabelEncoder
```

```
 4 labels = df.target.unique().tolist()
 5 print(labels)
 6
 7 encoder = LabelEncoder()
 8 encoder.fit(df.target.tolist()) # initialize encoder logic
 9
10 # the usual
11 # x -> features
12 # y -> labels
13 # we encode the labels
14
15 y_train = encoder.transform(train_data.target.tolist()) # encode the target into y_train
16 y_test = encoder.transform(test_data.target.tolist()) # encode to y_test
17
18 y_train = y_train.reshape(-1,1)
19 y_test = y_test.reshape(-1,1)
20
21 print("y_train",y_train.shape)
22 print("y_test",y_test.shape)
```

```
['Negative', 'Positive']
y_train (1038058, 1)
y_test (444883, 1)
```

```
 1 # sanity check
 2 # did we reshape, did we matrixizce the inputs
 3 print(
 4 y_test[0:5],
 5 x_test[0])
```

```
[[1]
 [1]
 [0]
 [1]
 [1]] [   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    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 1279   44   91  278  283  100   42   52]
```

```
 1 import numpy as np
 2 # embed the precious genshim word2vec model
 3
 4 embedding_matrix = np.zeros((len(wrdfreq), 300))
 5 for word, i in tknizer.word_index.items():
 6   if word in genshinmodel.wv:
 7     embedding_matrix[i] = genshinmodel.wv[word]
 8
 9 # the embedding matrix maps index to vectors
10
11 print(embedding_matrix.shape)
```

```
(454515, 300)
```

```
 1 embedding_layer = Embedding(len(wrdfreq), 300, weights=[embedding_matrix], input_length=140, trainable=False)
```

```
 1 # rest of the model layers
 2 model = Sequential()
 3 model.add(embedding_layer) # to make relevant connections
 4 model.add(Dropout(0.2))
 5 model.add(LSTM(200, dropout=0.2, return_sequences=True)) # multiple LSTM layers
 6 # dropouts prevent overfitting, and lower reliance on regularly occuring words
 7 # encourages to actually learn and not take shortcuts
 8 # afterall, its not a real human, we cant have it take shortcuts
 9 model.add(LSTM(100, recurrent_dropout=0.1, return_sequences=True))
10 model.add(Dropout(0.15))
11 # recurrent dropout drops the connection between recurernt cells
12 # i.e; memory and hidden. captures longer term dependancies
13 model.add(LSTM(50))
14 model.add(Dense(1, activation='sigmoid')) # sigmoid because we need binary output
```

```
15
16 model.summary()
```

```
WARNING:tensorflow:Layer lstm_6 will not use cuDNN kernels since it doesn't meet the criteria. It will use a generic GPU kernel as fall
Model: "sequential_2"
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| embedding (Embedding) | (None, 140, 300) | 136354500 |
| dropout_2 (Dropout) | (None, 140, 300) | 0 |
| lstm_5 (LSTM) | (None, 140, 200) | 400800 |
| lstm_6 (LSTM) | (None, 140, 100) | 120400 |
| dropout_3 (Dropout) | (None, 140, 100) | 0 |
| lstm_7 (LSTM) | (None, 50) | 30200 |
| dense_1 (Dense) | (None, 1) | 51 |

```
=================================================================
Total params: 136,905,951
Trainable params: 551,451
Non-trainable params: 136,354,500
```

```
1 # compile the model
2
3 model.compile(loss='binary_crossentropy',
4               optimizer="adam",
5               metrics=['accuracy'])
```

```
1 # callback so we dont train forever
2 # model checkpoint added retrospectively
3 # now im stuck waiting for this to be finished, which will take forever. Sadge
4 from keras.callbacks import ModelCheckpoint
5
6 callbacks = [ ReduceLROnPlateau(monitor='val_loss', patience=3, cooldown=0),
7               EarlyStopping(monitor='val_accuracy', min_delta=2e-3, patience=3),
8               ModelCheckpoint(filepath='model_checkpoint.h5', save_best_only=True, save_weights_only=False)
9              ]
```

```
 1 t = time.time()
 2 # train. my system takes ~ 3 hrs or so
 3 history = model.fit(x_train, y_train,
 4                     batch_size=1024,
 5                     epochs=10,
 6                     validation_split=0.12,
 7                     verbose=1,
 8                     callbacks=callbacks)
 9
10 print("time taken", time.time() - t, "s")
```

```
Epoch 1/10
893/893 [==============================] - ETA: 0s - loss: 0.4947 - accuracy: 0.7579WARNING:tensorflow:Early stopping conditioned on me
893/893 [==============================] - 648s 702ms/step - loss: 0.4947 - accuracy: 0.7579 - val_loss: 0.4595 - val_accuracy: 0.7821
Epoch 2/10
893/893 [==============================] - ETA: 0s - loss: 0.4608 - accuracy: 0.7807WARNING:tensorflow:Early stopping conditioned on me
893/893 [==============================] - 649s 726ms/step - loss: 0.4608 - accuracy: 0.7807 - val_loss: 0.4469 - val_accuracy: 0.7885
Epoch 3/10
893/893 [==============================] - ETA: 0s - loss: 0.4487 - accuracy: 0.7880WARNING:tensorflow:Early stopping conditioned on me
893/893 [==============================] - 643s 720ms/step - loss: 0.4487 - accuracy: 0.7880 - val_loss: 0.4408 - val_accuracy: 0.7931
Epoch 4/10
893/893 [==============================] - ETA: 0s - loss: 0.4405 - accuracy: 0.7929WARNING:tensorflow:Early stopping conditioned on me
893/893 [==============================] - 646s 724ms/step - loss: 0.4405 - accuracy: 0.7929 - val_loss: 0.4377 - val_accuracy: 0.7952
Epoch 5/10
893/893 [==============================] - ETA: 0s - loss: 0.4338 - accuracy: 0.7968WARNING:tensorflow:Early stopping conditioned on me
893/893 [==============================] - 616s 690ms/step - loss: 0.4338 - accuracy: 0.7968 - val_loss: 0.4379 - val_accuracy: 0.7962
Epoch 6/10
893/893 [==============================] - ETA: 0s - loss: 0.4280 - accuracy: 0.7999WARNING:tensorflow:Early stopping conditioned on me
893/893 [==============================] - 633s 709ms/step - loss: 0.4280 - accuracy: 0.7999 - val_loss: 0.4359 - val_accuracy: 0.7973
Epoch 7/10
893/893 [==============================] - ETA: 0s - loss: 0.4240 - accuracy: 0.8022WARNING:tensorflow:Early stopping conditioned on me
893/893 [==============================] - 611s 684ms/step - loss: 0.4240 - accuracy: 0.8022 - val_loss: 0.4364 - val_accuracy: 0.7960
Epoch 8/10
893/893 [==============================] - ETA: 0s - loss: 0.4201 - accuracy: 0.8045WARNING:tensorflow:Early stopping conditioned on me
```

```
893/893 [==============================] - 631s 707ms/step - loss: 0.4201 - accuracy: 0.8045 - val_loss: 0.4355 - val_accuracy: 0.7978
Epoch 9/10
893/893 [==============================] - ETA: 0s - loss: 0.4162 - accuracy: 0.8067WARNING:tensorflow:Early stopping conditioned on me
893/893 [==============================] - 630s 705ms/step - loss: 0.4162 - accuracy: 0.8067 - val_loss: 0.4353 - val_accuracy: 0.7974
Epoch 10/10
893/893 [==============================] - ETA: 0s - loss: 0.4129 - accuracy: 0.8083WARNING:tensorflow:Early stopping conditioned on me
893/893 [==============================] - 622s 697ms/step - loss: 0.4129 - accuracy: 0.8083 - val_loss: 0.4341 - val_accuracy: 0.7986
time taken 6342.21458530426 s
```

```python
1  # simple evaluation
2  score = model.evaluate(x_test, y_test, batch_size=512)
3  print("accuracy:",score[1])
4  print("loss:",score[0])
```
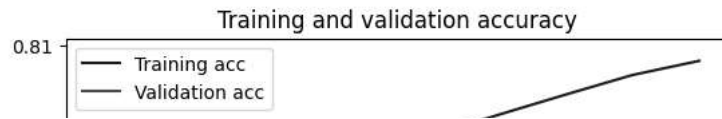
```
869/869 [==============================] - 57s 65ms/step - loss: 0.4352 - accuracy: 0.7981
accuracy: 0.7980952262878418
loss: 0.43522322177886963
```

```python
1  import pickle
2
3  # Serialize the objects (x_test and y_test)
4  serialized_x_test = pickle.dumps(x_test)
5  serialized_y_test = pickle.dumps(y_test)
6
7  # Create a zip file
8  with zipfile.ZipFile('serialized_data.zip', 'w') as zip_file:
9      # Add the serialized objects to the zip file
10     zip_file.writestr('x_test.pickle', serialized_x_test)
11     zip_file.writestr('y_test.pickle', serialized_y_test)
```

```python
1  import zipfile
2  # lets save the model because we dont want to train it for 2 hours every time we open the notebook
3  model.save('snetiment_analzyer.h5')
4  # lets zip it too, so we can save space
5  with zipfile.ZipFile('snetiment_analzyer.zip', 'w', zipfile.ZIP_DEFLATED) as zipf:
6      zipf.write('snetiment_analzyer.h5')
```

```python
1  # since we have a history object from earlier, we can plot how our model evolved
2  # really helpful to see some cool outputs
3  import matplotlib.pyplot as plt
4  %matplotlib inline
5
6  acc = history.history['accuracy']
7  val_acc = history.history['val_accuracy']
8  loss = history.history['loss']
9  val_loss = history.history['val_loss']
10
11 epochs = range(len(acc)) # 3 for us, so not real 'graph', just a few points
12
13 plt.plot(epochs, acc, 'b', label='Training acc')
14 plt.plot(epochs, val_acc, 'r', label='Validation acc')
15 plt.title('Training and validation accuracy')
16 plt.legend()
17
18 plt.show()
```

Training and validation accuracy

0.81

— Training acc
— Validation acc

```
1 def score_to_sentiment(score, include_neutral=True):
2   if include_neutral:
3       label = 'Neutral'
4       if score <= 0.42:
5           label = 'Negative'
6       elif score >= 0.58:
7           label = 'Positive'
8       return label
9   else:
10      return 'Negative' if score < 0.5 else 'Positive'
11
12 def predict(text, include_neutral=True):
13     # Tokenize text
14     x_test = pad_sequences(tknizer.texts_to_sequences([text]), maxlen=140)
15     # Predict
16     score = model.predict([x_test])[0]
17     # Decode sentiment
18     label = score_to_sentiment(score, include_neutral)
19
20     return {"text": text, "label": label, "score": float(score)}
```

```
1 predict("I love the music")
```

```
1/1 [==============================] - 1s 803ms/step
{'text': 'I love the music', 'label': 'Positive', 'score': 0.9515533447265625}
```

```
1 predict("It's impressive how you always manage to find new ways to avoid taking responsibility.")
```

```
1/1 [==============================] - 0s 193ms/step
{'text': "It's impressive how you always manage to find new ways to avoid taking responsibility.",
 'label': 'Neutral',
 'score': 0.5647751092910767}
```

```
1 predict("Every link in Wikipedia is around a six chain away from Hitler")
```

```
1/1 [==============================] - 0s 100ms/step
{'text': 'Every link in Wikipedia is around a six chain away from Hitler',
 'label': 'Positive',
 'score': 0.6825700998306274}
```

```
1 predict("My dad beats me, my mom beats me, my friends beat me. But I feel safe with Vitality because they can't beat anyone omeglaul")
```

```
1/1 [==============================] - 0s 111ms/step
{'text': "My dad beats me, my mom beats me, my friends beat me. But I feel safe with Vitality because they can't beat anyone
omeglaul",
 'label': 'Positive',
 'score': 0.7683548927307129}
```

```
1 predict("doing drugs sucks...")
```

```
1/1 [==============================] - 0s 124ms/step
{'text': 'doing drugs sucks...',
 'label': 'Positive',
 'score': 0.5968612432479858}
```

```
1 predict("i hate my life")
```

```
1/1 [==============================] - 0s 133ms/step
{'text': 'i hate my life', 'label': 'Negative', 'score': 0.03444806858897209}
```

```
1 # try to add confusion matrix later
2
3 # y_pred_1d = []
4 # y_test_1d = list(test_data.target)
5 # scores = model.predict(x_test)
6 # y_pred_1d = [score_to_sentiment(score, include_neutral=False) for score in scores]
```

```
---------------------------------------------------------------------------
InternalError                             Traceback (most recent call last)
<ipython-input-54-27fc5bc5a1fe> in <cell line: 3>()
      1 y_pred_1d = []
      2 y_test_1d = list(df.target)
----> 3 scores = model.predict(x_test)
      4 y_pred_1d = [score_to_sentiment(score, include_neutral=False) for score in
scores]
```

‹› 1 frames

```
/usr/local/lib/python3.10/dist-packages/tensorflow/python/framework/ops.py in
_numpy(self)
   1126         return self._numpy_internal()
   1127     except core._NotOkStatusException as e:  # pylint: disable=protected-
access
-> 1128         raise core._status_to_exception(e) from None  # pylint:
disable=protected-access
   1129
```

1

⚠ 1s    completed at 6:56 PM    ● ✕