

```
In [1]: #Pretrained word2vec embeddings

In [2]: ##from torchtext import data
#from torchtext import datasets
import torch
import spacy
import pandas as pd
import time
import re # for regular expressions
import torch.nn as nn
import torch.nn.functional as F
from sklearn.model_selection import train_test_split
import spacy
nlp = spacy.load('en')

data = pd.read_csv('transferlearning-dl-spring2020/train.csv')

data.head()

questions = data['text']
labels = data['target']

train_data, valid_data, y_train, y_valid = train_test_split(questions, labels,
                                                            random_state=2,
                                                            test_size=0.2)

train_data.head()

Out[1]: 251      @USER All you need to know is he is empty inside
      252      @USER I AM READ: P1 URL
      4776      @USER Go roger i quit watching anyway nrl is o..
      4831      @USER You our dogs should totally fuck
      2835      @USER She is a troll. Not open to facts
      Name: text, dtype: object

In [2]: data.head()

Out[2]:
```

	id	text	target
0	86426	@USER She should ask a few native americans wh...	1
1	16820	Amazon is investigating Chinese employees who...	0
2	62688	@USER Someone should vetaken this piece of sh...	1
3	4305	@USER @USER Obama wanted liberals amp illeg...	0
4	97670	@USER Liberals are all Kokkoo !!	1

```
In [3]: def remove_pattern(input_text, pattern):
    p = re.findall(pattern, input_text)
    for i in range(len(p)):
        input_text = re.sub(i, '', input_text)
    return input_text

data['text'] = np.vectorize(remove_pattern)(data['text'], "%8(w)w")

data['text'] = data['text'].str.replace("%a-zA-58]", " ")
data['text'] = data['text'].str.replace("#", "")

data['text'] = data['text'].apply(lambda x: ' '.join([w for w in x.split() if len(w)>2]))

tokenized_tweet = data['text'].apply(lambda x: x.split()) # tokenizing
tokenized_tweet.head()

len(tokenized_tweet)

#from nltk.stem.porter import *
#stemmer = PorterStemmer()

#tokenized_tweet = tokenized_tweet.apply(lambda x: [stemmer.stem(i) for i in x]) # stemming

from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
tokenized_tweet.apply(lambda x: [lemmatizer.lemmatize(i) for i in x]) # stemming
#lemmatized_output = ' '.join([lemmatizer.lemmatize(token) for token in tokens])

Out[3]: 0      (She, should, ask, few, native, americans, wha...
      1      (Amazon, investigating, chinese, employe...
      2      (Someone, should, vetaken, this, piece, shi...
      3      (Obama, wanted, liberal, amp, illegals, move...
      4      (Was, literally, just, talking, about, this, lo...
      5      (Iwa, literally, just, talking, about, this, lo...
      6      (Buy, more, icecream)
      7      (What, not, fault, you, support, gun, control)
      8      (What, the, difference, between, kavanaugh, an...
      9      (You, are, lying, corrupt, traitor, nobody, wa...
      10     (Like, soda, like, like, boarder, with, lot, tch)
      11     (You, are, also, the, king, taste)
      12     (MAGA, Sing, like, one, listening, Love, like...
      13     (The, time, right, for, this, House, respond, ...
      14     (Beside, jack, mom, and, maybe, Ope, hand, dow...
      15     (gun, control, That, all, these, kid are, ask...
      16     (fuck, going, people, There, the, men, room, a...
      17     (Been, Willie, fan, since, before, most, you, ...
      18     (Tbh, these, day, just, doing, like, people, esp...
      19     (South, Korean, Official, Leaders, will, discus...
      20     (You, can, talk, ill, hooper, too)
      21     (Ireal, nika, better, chasing, the, title)
      22     (She, whom, are, you, referring, Hillary, You, ...
      23     (Glad, see, your, friend, are, supporting, Met...
      24     (Iest, you, are, but, wa, asking, what, about, ...
      25     (Wonder, being, apologetic, and, more, social...
      26     (Any, update, eren, your, blatant, racism, W...
      27     (Blow, hard)
      28     (That, mean, you, are, max, lvl, twitter, user...
      29     (Please, explain, what, controlled, opposition...
      ...
      9316     (Berkeley, Antifa, not, agree, with, you, URL)
      9317     (Advocate, for, gun, control, while, breaking...
      9318     (More, press, talk, about, this, remar...
      9319     (Mohammadhasanahaid, traitor, the, USA, amp...
      9320     (nah, coz, you, guys, fucked)
      9321     (Cats, are, just, special, dog)
      9322     (Best, news, ever, for, GOP, Ready, the, stral...
      9323     (Strength, letting, the, universe, You, are, n...
      9324     (Wow, You, are, good)
      9325     (still, lie, just, like, Obama, talk, right, a...
      9326     (precious)
      9327     (Boonings)
      9328     (Everything, else, wa, ten, year, ago, You, AR...
      9329     (Right, Dang, She, the)
      9330     (McKaven, engaged, publicity, stunt, never, sp...
      9331     (man, You, are, going, trigger, the, fanboys, ...
      9332     (Chelsea, never, end, You, are, always, the, B...
      9333     (I have, the, conservative, accepted, the, title)
      9334     (Alt, Right, amp, Antifa, are, for, coward, fa...
      9335     (Did, you, serve, You, rate, bottom, the, bar...
      9336     (Just, trying, make, good, with, his, libnat...
      9337     (have, the, conservative, accepted, the, title)
      9338     (Can, all, agree, that, Tomlins, seat, heating...
      9339     (Involved, because, wa, there, Now, need, man...
      9340     (How, much, lonely, she, and, how, much, user...
      9341     (BUT, GUN, CONTROL)
      9342     (say, you, are, mad, now, you, will, say, tire...
      9343     (Pretest, complete, amp, followed, all, patriot...
      9344     (And, why, report, this, garbage, don, give, t...
      9345     (Pussy)
      Name: text, dtype: object

In [4]: tokenized_tweet[1:10]

Out[4]: 1      (Amazon, investigating, chinese, employees, wh...
      2      (Someone, should, vetaken, this, piece, shi...
      3      (Obama, wanted, liberals, amp, illegals, move...
      4      (Was, literally, just, talking, about, this, lo...
      5      (Iwa, literally, just, talking, about, this, lo...
      6      (Buy, more, icecream)
      7      (I not, fault, you, support, gun, control)
      8      (What, the, difference, between, kavanaugh, an...
      9      (You, are, lying, corrupt, traitor, nobody, wa...
      Name: text, dtype: object

In [8]: # Using pre-trained word2vec embeddings created by Google

import gensim
model = gensim.models.KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)

#model.train(tokenized_tweet, total_examples= len(combi('description')), epochs=50)

def word_vector_pretrained(tokens, size):
    vec = np.zeros(size).reshape((1, size))
    count = 0
    for word in tokens:
        try:
            vec += model[word].reshape((1, size))
            count += 1
        except KeyError: # handling the case where the token is not in vocabulary
            continue
    if count != 0:
        vec /= count
    return vec

wordvec_arrays = np.zeros((len(tokenized_tweet), 300))

for i in range(len(tokenized_tweet)):
    wordvec_arrays[i,:] = word_vector_pretrained(tokenized_tweet[i], 300)

wordvec_df = pd.DataFrame(wordvec_arrays)
wordvec_df.shape

X_train, X_test, y_train, y_test = train_test_split(wordvec_df, labels, test_size=0.2, random_state=1)

In [9]: np.unique(y_train)

Out[9]: array([0, 1])

In [10]: tokenized_tweet[1:5]

Out[10]: 1      (Amazon, investigating, chinese, employees, wh...
      2      (Someone, should, vetaken, this, piece, shi...
      3      (Obama, wanted, liberals, amp, illegals, move...
      4      (Was, literally, just, talking, about, this, lo...
      Name: text, dtype: object

In [11]: X_train[1:10]

Out[11]:
```

	0	1	2	3	4	5	6	7	8	9	...	290	291
5693	0.006725	0.010405	-0.052087	0.091202	-0.056473	0.026924	0.077883	-0.019114	0.072357	-0.027234	...	-0.086390	0.052601
4507	0.029175	0.037109	0.032051	0.113220	-0.112213	0.008484	0.028618	0.043793	0.001225	-0.042809	...	0.010437	0.056587
3891	0.017432	0.003791	0.020201	0.063634	-0.040874	0.010851	0.012351	-0.074840	0.095403	0.074498	...	-0.059332	0.024525
5429	0.023418	0.038795	-0.022937	-0.007890	0.017090	-0.051917	0.034424	0.123288	0.058185	0.031116	...	-0.018947	-0.060611
1104	0.023802	0.033569	0.142127	0.134559	-0.098013	-0.016130	0.057371	-0.051827	0.050809	0.048138	...	-0.015437	0.078177
6560	0.061476	0.055678	0.018585	0.084633	-0.116528	-0.019473	-0.007350	-0.011179	0.071787	0.027387	...	-0.034809	0.024687
3332	0.066332	0.052708	0.082484	0.005835	-0.066877	0.020801	0.020983	0.003180	0.003379	0.077595	...	-0.000912	0.021022
6211	0.001829	-0.012245	-0.006155	-0.039429	0.000488	0.053101	-0.025269	0.000488	0.074219	0.051147	...	-0.229482	-0.023682
2384	0.067283	-0.040487	0.053802	0.194775	-0.055338	0.013220	0.075983	-0.038326	0.097266	0.041626	...	-0.051514	-0.028213

9 rows x 300 columns

```
In [12]: from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import f1_score

In [13]: lreg = LogisticRegression(random_state = 10)
lreg.fit(X_train, y_train) # training the model

prediction = lreg.predict(X_test) # predicting on the validation set

prediction_int = prediction.astype(np.int)

f1_score(y_test, prediction, average = 'macro') # calculating f1 score

Out[13]: 0.6691278341911253

In [14]: data_test = pd.read_csv('transferlearning-dl-spring2020/test.csv')

In [15]: data_test['text'] = np.vectorize(remove_pattern)(data_test['text'], "%8(w)w")

data_test['text'] = data_test['text'].str.replace("%a-zA-58]", " ")
data_test['text'] = data_test['text'].str.replace("#", "")

data_test['text'] = data_test['text'].apply(lambda x: ' '.join([w for w in x.split() if len(w)>2]))

tokenized_tweet_test = data_test['text'].apply(lambda x: x.split()) # tokenizing
tokenized_tweet_test.head()

len(tokenized_tweet_test)

#from nltk.stem.porter import *
#stemmer = PorterStemmer()

#tokenized_tweet_test = tokenized_tweet_test.apply(lambda x: [stemmer.stem(i) for i in x]) # stemming

from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
tokenized_tweet_test.apply(lambda x: [lemmatizer.lemmatize(i) for i in x]) # stemming

Out[15]: 0      (home, you, drunk, MAGA, Trump, URL)
      1      (no, Tough, shit)
      2      (Canada, home, need, another, CUCV, already, an...
      3      (should, scare, every, American, She, playing...
      4      (LOL, Throwing, the, BULLSHIT, Flag, such, no...
      5      (You, are, correct)
      6      (Kind, like, when, conservative, wanna, associ...
      7      (The, only, thing, the, Democrats, have, lying...
      8      (You, know, what, going, happen, now, going, h...
      9      (You, are, not, very, smart, are, you, why, fol...
      10     (That, expected, you, placate, the, violent, l...
      11     (Wow, you, liberal, really, don, have, sense, ...
      12     (FUCKING, REDDIT)
      13     (too, Her, wedding, wa, the, best, She, very...
      14     (URL, any, your, announcement, every, come, fr...
      15     (out, British, people, are, basically, full, r...
      16     (GUNCONTROL, advocate, must, STOP, falling, al...
      17     (Fuck, off)
      18     (dumb, and, dumber, all, one, president, two, ...
      19     (shit, Twitter, silence, also, Jones, retail...
      20     (That, because, you, are, old, man)
      21     (Ide, beautiful, person, teach, value, her, ki...
      22     (tear, made, water, and, feeling)
      23     (babytitting, kid, people, how, old, you, thin...
      24     (WTC, this, threat, maga, Qanon, WakeUpAmerica...
      25     (Unfortunately, America, system, like, that, w...
      26     (You, you, are, always, present, bet, why, lol)
      27     (MAGA, YOU, ARE, ALL, FOR, TRUMP, FOLLOW, AND...
      28     (For, the, record, know, Doug, Jones, personal...
      29     (There, are, many, dumb, argument, for, gun, c...
      ...
      3864     (Probably, while, screaming, about, gun, control)
      3865     (Thinking, about, the, Student, Loan, crisis, ...
      3866     (the, lying, increase, the, cost, drug, Canada...
      3867     (damn, straight, Nothing, but, good, vibe, com...
      3868     (THIS, SHOULD, REMIND, ALL, Patriots, maga, WH...
      3869     (sorry, what, B&S, listening, that, amp, n...
      3870     (Americans, make, great, client, buy, many, ap...
      3871     (You, should, know, better, use, common, sense...
      3872     (Resp, lying, bid, one, listening)
      3873     (You, are, great, model, for, inspiration)
      3874     (Chris, Chris, Chris, Aye, you, forgetting, th...
      3875     (Where, will, Antifa, get, their, cloth, now)
      3876     (the, role, model, Adam, you, are, not)
      3877     (The, entire, way, the, dems, have, handed, t...
      3878     (Beto, Rourke, Ted, Cruz, Latest, Polly, democ...
      3879     (Project, much, you, Bvery, campaign, filter, g...
      3880     (Link, and, knuckle, the, you, going, dodge, ...
      3881     (Omg, not, even, interested, his, age, but, da...
      3882     (shut, you, wern, joking, wit...)
      3883     (Bibi, look, like, Stalin, when, Stalin, wa, v...
      3884     (right, say, that, housing, association, shoul...
      3885     (Boise, State, fan, can, tell, you, two, thing...
      3886     (advocating, for, conduct, about, this, lo...
      3887     (when, you, coming, ohio)
      3888     (Liars, like, the, Antifa, twin, you, vigorous...
      3889     (Billy, you, have, short, memory, Obama, already...
      3890     (She, not, the, brightest, light, the, tree)
      3891     (Sometimes, get, strong, vibe, from, people, a...
      3892     (Benidorm, Creamfields, Mags, Not, too, shabby...
      3893     (Spanishrevenge, justice, HumanRights, and, fr...
      Name: text, dtype: object

In [16]: wordvec_arrays_test = np.zeros((len(tokenized_tweet_test), 300))

for i in range(len(tokenized_tweet_test)):
    wordvec_arrays_test[i,:] = word_vector_pretrained(tokenized_tweet_test[i], 300)

wordvec_df_test = pd.DataFrame(wordvec_arrays_test)
wordvec_df_test.shape

Out[16]: (3894, 300)

In [17]: X_test1 = np.array(wordvec_df_test)

prediction = lreg.predict(X_test1) # predicting on the validation set

prediction_int1 = prediction.astype(np.int)

In [18]: prediction_int1

Out[18]: array([0, 1, 0, ..., 0, 0, 0])

In [19]: data_test['target'] = prediction_int1
submission = data_test[['id', 'target']]
submission.to_csv('pretrained_word2vec.csv', index=False) # writing data to a CSV file
submission

Out[19]:
```

	id	target
0	90194	0
1	17444	1
2	13384	0
3	54820	1
4	56117	1
5	87757	0
6	12681	0
7	12609	0
8	70380	0
9	12108	0
10	14726	1
11	17477	0
12	48845	0
13	47311	0
14	75689	0
15	84102	0
16	10607	0
17	98992	1
18	53264	1
19	54842	0
20	48995	0
21	72353	0
22	50759	1
23	14574	0
24	93119	0
25	43193	1
26	59537	0
27	32317	0
28	76680	0
29	80561	0
...
3864	54190	0
3865	28037	0
3866	77430	0
3867	75815	0
3868	87290	0
3869	99475	0
3870	43233	0
3871	37686	0
3872	77905	0
3873	83400	0
3874	84081	0
3875	76469	1
3876	20841	0
3877	90599	0
3878	32598	0
3879	43964	0
3880	97745	0
3881	86716	0
3882	21033	1
3883	66832	0
3884	28996	0
3885	64713	0
3886	63482	1
3887	11132	0
3888	84716	0
3889	90041	0
3890	98824	0
3891	95338	1
3892	67210	0
3893	46552	0

3894 rows x 2 columns

```
In [ ]: # Pretrained word2vec embeddings 2

In [20]: # Using pre-trained word2vec embeddings created by Google

import gensim, downloader as api
model = api.load('glove-twitter-200')

#model = gensim.models.KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)

#model.train(tokenized_tweet, total_examples= len(combi('description')), epochs=50)

def word_vector_pretrained(tokens, size):
    vec = np.zeros(size).reshape((1, size))
    count = 0
    for word in tokens:
        try:
            vec += model[word].reshape((1, size))
            count += 1
        except KeyError: # handling the case where the token is not in vocabulary
            continue
    if count != 0:
        vec /= count
    return vec

wordvec_arrays = np.zeros((len(tokenized_tweet), 200))

for i in range(len(tokenized_tweet)):
    wordvec_arrays[i,:] = word_vector_pretrained(tokenized_tweet[i], 200)

wordvec_df = pd.DataFrame(wordvec_arrays)
wordvec_df.shape

X_train, X_test, y_train, y_test = train_test_split(wordvec_df, labels, test_size=0.2, random_state=1)

In [21]: lreg = LogisticRegression(random_state = 10)
lreg.fit(X_train, y_train) # training the model

prediction = lreg.predict(X_test) # predicting on the validation set

prediction_int1 = prediction.astype(np.int)

prediction_int1

f1_score(y_test, prediction, average = 'macro') # calculating f1 score

#anaconda3/lib/python3.7/site-packages/sklearn/linear_model/_logistic.py:940: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG)

Out[21]: 0.6552467092841944

In [23]: data_test = pd.read_csv('transferlearning-dl-spring2020/test.csv')

data_test['text'] = np.vectorize(remove_pattern)(data_test['text'], "%8(w)w")

data_test['text'] = data_test['text'].str.replace("%a-zA-58]", " ")
data_test['text'] = data_test['text'].str.replace("#", "")

data_test['text'] = data_test['text'].apply(lambda x: ' '.join([w for w in x.split() if len(w)>2]))

tokenized_tweet_test = data_test['text'].apply(lambda x: x.split()) # tokenizing
tokenized_tweet_test.head()

len(tokenized_tweet_test)

#from nltk.stem.porter import *
#stemmer = PorterStemmer()

#tokenized_tweet_test = tokenized_tweet_test.apply(lambda x: [stemmer.stem(i) for i in x]) # stemming

from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
tokenized_tweet_test.apply(lambda x: [lemmatizer.lemmatize(i) for i in x]) # stemming

wordvec_arrays_test = np.zeros((len(tokenized_tweet_test), 200))

for i in range(len(tokenized_tweet_test)):
    wordvec_arrays_test[i,:] = word_vector_pretrained(tokenized_tweet_test[i], 200)

wordvec_df_test = pd.DataFrame(wordvec_arrays_test)
wordvec_df_test.shape

X_test1 = np.array(wordvec_df_test)

prediction = lreg.predict(X_test1) # predicting on the validation set

prediction_int1 = prediction.astype(np.int)

prediction_int1

data_test['target'] = prediction_int1
submission = data_test[['id', 'target']]
submission.to_csv('pretrained_word2vec1.csv', index=False) # writing data to a CSV file
submission

Out[23]:
```

	id	target
0	90194	0
1	17444	1
2	13384	0
3	54820	1
4	56117	1
5	87757	0
6	12681	0
7	12609	0
8	70380	0
9	12108	0
10	14726	1
11	17477	0
12	48845	0
13	47311	0
14	75689	0
15	84102	0
16	10607	0
17	98992	1
18	53264	1
19	54842	0
20	48995	0
21	72353	0
22	50759	0
23	14574	0
24	93119	0
25	43193	1
26	59537	0
27	32317	0
28	76680	0
29	80561	0
...
3864	54190	0
3865	28037	0
3866	77430	0
3867	75815	0
3868	87290	0
3869	99475	0
3870	43233	0
3871	37686	0
3872	77905	0
3873	83400	0
3874	84081	0
3875	76469	1
3876	20841	0
3877	90599	0
3878	32598	0
3879	43964	0
3880	97745	1
3881	86716	0
3882	21033	1
3883	66832	0
3884	28996	0
3885	64713	0
3886	63482	1
3887	11132	0
3888	84716	0
3889	90041	0
3890	98824	0
3891	95338	1
3892	67210	0
3893	46552	0

3894 rows x 2 columns

```
In [ ]: # Fine tuned word2vec

In [56]: model_w2v = gensim.models.Word2Vec(tokenized_tweet,
    size=100, # desired no. of features/independent variables
    window=5, # context window size
    min_count=2,
    sg = 1, # 1 for skip-gram model
    hs = 0,
    negative = 10, # for negative sampling
    workers=2, # no. of cores
    seed = 34)

model_w2v.train(tokenized_tweet, total_examples= len(tokenized_tweet), epochs=200)

def word_vector_test(tokens, size):
    vec = np.zeros(size).reshape((1, size))
    count = 0
    for word in tokens:
        try:
            vec += model_w2v[word].reshape((1, size))
            count += 1
        except KeyError: # handling the case where the token is not in vocabulary
            continue
    if count != 0:
        vec /= count
    return vec

wordvec_arrays_test = np.zeros((len(tokenized_tweet), 50))

for i in range(len(tokenized_tweet)):
    wordvec_arrays_test[i,:] = word_vector_test(tokenized_tweet[i], 50)

wordvec_df_test = pd.DataFrame(wordvec_arrays_test)
wordvec_df_test.shape

#anaconda3/lib/python3.7/site-packages/sklearn/linear_model/_logistic.py:940: DeprecationWarning: Call to deprecated 'getitem' (Method will be removed in 4.0.0, use self.wv._getitem__ instead).

Out[56]: (9346, 50)

In [58]: X_train, X_test, y_train, y_test = train_test_split(wordvec_df_test, labels, test_size=0.2, random_state=1)
lreg = LogisticRegression(random_state = 20)
lreg.fit(X_train, y_train) # training the model

prediction = lreg.predict(X_test) # predicting on the validation set

prediction_int1 = prediction.astype(np.int)

prediction_int1

f1_score(y_test, prediction, average = 'macro') # calculating f1 score

Out[58]: 0.5291078614322713

In [ ]: # CNN for text classification

In [60]: top_data_df_small = data

In [61]: from gensim.utils import simple_preprocess
wordvec_arrays_test = np.zeros((len(tokenized_tweet), 100))

for i in range(len(tokenized_tweet)):
    wordvec_arrays_test[i,:] = word_vector_test(tokenized_tweet[i], 100)

wordvec_df_test = pd.DataFrame(wordvec_arrays_test)
wordvec_df_test.shape

X_train, X_test, y_train, y_test = train_test_split(wordvec_df_test, labels, test_size=0.2, random_state=1)

Value counts for Train sentiments
0      4365
1      2177
Name: target, dtype: int64
<class 'pandas.core.frame.DataFrame'>
index
0      8030      least you are honest lmao
1      3122      [a2v, that, shit, wa, her, url]
2      3222      (cuti, connect, authnet)
3      4232      What interview Ewanne honest and right from th...
4      1678      Warriorcoach fundraising for gun control Brad...
...
index
0      8030      least you are honest lmao
1      3122      [a2v, that, shit, wa, her, url]
2      3222      (cuti, connect, authnet)
3      4232      What interview evann, honest, and, right, f...
4      1678      Warriorcoach, fundatels, for, gun, control, B...
...
index
0      8030      least you are honest lmao
1      3122      [a2v, that, shit, wa, her, url]
2      3222      (cuti, connect, authnet)
3      4232      What interview evann, honest, and, right, f...
4      1678      Warriorcoach, fundatels, for, gun, control, B...
...
index
0      8030      least you are honest lmao
1      3122      [a2v, that, shit, wa, her, url]
2      3222      (cuti, connect, authnet)
3      4232      What interview evann, honest, and, right, f...
4      1678      Warriorcoach, fundatels, for, gun, control, B...
...
index
0      8030      least you are honest lmao
1      3122      [a2v, that, shit, wa, her, url]
2      3222      (cuti, connect, authnet)
3      4232      What interview evann, honest, and, right, f...
4      1678      Warriorcoach, fundatels, for, gun, control, B...
...
index
0      8030      least you are honest lmao
1      3122      [a2v, that, shit, wa, her, url]
2      3222      (cuti, connect, authnet)
3      4232      What interview evann, honest, and, right, f...
4      1678      Warriorcoach, fundatels, for, gun, control, B...
...
index
0      8030      least you are honest lmao
1      3122      [a2v, that, shit, wa, her, url]
2      3222      (cuti, connect, authnet)
3      4232      What interview evann, honest, and, right, f...
4      1678      Warriorcoach, fundatels, for, gun, control, B...
...
index
0      8030      least you are honest lmao
1      3122      [a2v, that, shit, wa, her, url]
2      3222      (cuti, connect, authnet)
3      4232      What interview evann, honest, and, right, f...
4      1678      Warriorcoach, fundatels, for, gun
```



```
[64]: import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim

import torch
# Use cuda if present
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print("Device available for running: ", device)

Device available for running:
cpu

In [67]: from gensim.models import Word2Vec

size = 300
window = 3
min_count = 1
workers = 3
sg = 1

# Function to train word2vec model
def make_word2vec_model(top_data_df_small, padding=True, sg=1, min_count=1, size=500, workers=3, window=3):
    if padding:
        temp_df = (top_data_df_small[['stemmed_tokens']]).values
        temp_df = list(temp_df)
        temp_df.append(['pad'])
        word2vec_file = 'word2vec.' + str(size) + '_PAD.model'
    else:
        temp_df = top_data_df_small[['stemmed_tokens']]
        word2vec_file = 'word2vec.' + str(size) + '.model'
    w2v_model = Word2Vec(temp_df, min_count = min_count, size = size, workers = workers, window = window, sg = sg)

    w2v_model.save(word2vec_file)
    return w2v_model, word2vec_file

# Train Word2Vec model
w2v_model, word2vec_file = make_word2vec_model(top_data_df_small, padding=True, sg=sg, min_count=min_count, size=size, workers=workers, window=window)

9346

In [68]: # Function to get the output tensor
def make_target(label):
    if label == 0:
        return torch.tensor([0], dtype=torch.long, device=device)
    elif label == 0:
        return torch.tensor([1], dtype=torch.long, device=device)
    else:
        return torch.tensor([2], dtype=torch.long, device=device)

In [71]: EMBEDDING_SIZE = 500
NUM_FILTERS = 10
import gensim

class CnnTextClassifier(nn.Module):
    def __init__(self, vocab_size, num_classes, window_size=(1,2,3,5)):
        super(CnnTextClassifier, self).__init__()
        w2vmodel = gensim.models.KeyedVectors.load('word2vec_500_PAD.model')
        weights = w2vmodel.wv
        # With pretrained embeddings
        self.embedding = nn.Embedding.from_pretrained(torch.FloatTensor(weights), padding_idx=w2vmodel.wv.vocab['pad'].index)
        # Without pretrained embeddings
        self.embedding = nn.Embedding(vocab_size, EMBEDDING_SIZE)

        self.conv = nn.ModuleList([
            nn.Conv2d(1, NUM_FILTERS, [window_size, EMBEDDING_SIZE], padding=(window_size - 1, 0))
            for window_size in window_size
        ])

        self.fc = nn.Linear(NUM_FILTERS * len(window_size), num_classes)

    def forward(self, x):
        x = self.embedding(x) # [B, T, E]

        # Apply a convolution + max_pool layer for each window size
        xs = []
        for conv in self.convs:
            x2 = torch.tan(conv(x))
            x2 = torch.squeeze(x2, -1)
            x2 = F.max_pool1d(x2, x2.size(2))
            xs.append(x2)
        x = torch.cat(xs, 2)

        # FC
        x = x.view(x.size(0), -1)
        logits = self.fc(x)

        probs = F.softmax(logits, dim = 1)

        return probs

In [73]: max_sen_len = top_data_df_small.stemmed_tokens.map(len).max()
padding_idx = w2vmodel.wv.vocab['pad'].index
def make_word2vec_vector_cnn(sentence):
    padded_x = [padding_idx for i in range(max_sen_len)]
    i = 0
    for word in sentence:
        if word not in w2vmodel.wv.vocab:
            padded_x[i] = 0
            print(word)
        else:
            padded_x[i] = w2vmodel.wv.vocab[word].index
            i += 1
    return torch.tensor(padded_x, dtype=torch.long, device=device).view(1, -1)

In [74]: NUM_CLASSES = 3
VOCAB_SIZE = len(w2vmodel.wv.vocab)

cnn_model = CnnTextClassifier(vocab_size=VOCAB_SIZE, num_classes=NUM_CLASSES)
loss_function = nn.CrossEntropyLoss()
optimizer = optim.Adam(cnn_model.parameters(), lr=0.001)
num_epochs = 30

# Open the file for writing loss
loss_file_name = 'cnn_class_big_loss_with_padding.csv'
f = open(loss_file_name, 'w')
f.write('Iter, loss')
f.write('\n')
losses = []
cnn_model.train()

for epoch in range(num_epochs):
    print("Epoch" + str(epoch + 1))
    train_loss = 0
    for index, row in X_train.iterrows():
        # Clearing the accumulated gradients
        cnn_model.zero_grad()

        # Make the bag of words vector for stemmed tokens
        bow_vec = make_word2vec_vector_cnn(row['stemmed_tokens'])

        # Forward pass to get output
        probs = cnn_model(bow_vec)

        # Get the target label
        target = make_target(Y_train['target'][index])

        # Calculate Loss: softmax --> cross entropy loss
        loss = loss_function(probs, target)
        train_loss += loss.item()

        # Getting gradients w.r.t. parameters
        loss.backward()

        # Updating parameters
        optimizer.step()

    # if index == 0:
    #     continue
    print("Epoch ran "+" str(epoch+1))
    f.write(str(epoch+1)) + ", " + str(train_loss / len(X_train))
    f.write('\n')
    train_loss = 0

torch.save(cnn_model, 'cnn_big_model_500_with_padding.pth')

f.close()
print("Input vector")
print(bow_vec.cpu().numpy())
print("Probs")
print(probs)
print(torch.argmax(probs, dim=1).cpu().numpy()[0])

Epoch1
Epoch ran :1
Epoch2
Epoch ran :2
Epoch3
Epoch ran :3
Epoch4
Epoch ran :4
Epoch5
Epoch ran :5
Epoch6
Epoch ran :6
Epoch7
Epoch ran :7
Epoch8
Epoch ran :8
Epoch9
Epoch ran :9
Epoch10
Epoch ran :10
Epoch11
Epoch ran :11
Epoch12
Epoch ran :12
Epoch13
Epoch ran :13
Epoch14
Epoch ran :14
Epoch15
Epoch ran :15
Epoch16
Epoch ran :16
Epoch17
Epoch ran :17
Epoch18
Epoch ran :18
Epoch19
Epoch ran :19
Epoch20
Epoch ran :20
Epoch21
Epoch ran :21
Epoch22
Epoch ran :22
Epoch23
Epoch ran :23
Epoch24
Epoch ran :24
Epoch25
Epoch ran :25
Epoch26
Epoch ran :26
Epoch27
Epoch ran :27
Epoch28
Epoch ran :28
Epoch29
Epoch ran :29
Epoch30
Epoch ran :30
Input vector
[[ 235  13  16  0  928  172  358  244  7  7 10887 10887
 10887 10887 10887 10887 10887 10887 10887 10887 10887 10887 10887
 10887 10887 10887 10887 10887 10887 10887 10887 10887 10887
 10887]
Probs
tensor([[1.1411e-11, 9.5633e-01, 4.3666e-02]], grad_fn=<SoftmaxBackward>)
1

In [79]: from sklearn.metrics import classification_report
bow_cnn_predictions = []
original_labels_cnn_bow = []
cnn_model.eval()
loss_df = pd.read_csv('cnn_class_big_loss_with_padding.csv')
print(loss_df.columns)
# loss_df.plot('loss')
for torch.no_grad():
    for index, row in X_test.iterrows():
        bow_vec = make_word2vec_vector_cnn(row['stemmed_tokens'])
        probs = cnn_model(bow_vec)
        _, predicted = torch.max(probs.data, 1)
        bow_cnn_predictions.append(predicted.cpu().numpy()[0])
        original_labels_cnn_bow.append(make_target(Y_test['target'][index]).cpu().numpy()[0])

print(classification_report(original_labels_cnn_bow, bow_cnn_predictions))
loss_file_name = 'cnn_class_big_loss_with_padding.csv'
loss_df = pd.read_csv(loss_file_name)
print(loss_df.columns)
plt.500_padding_30_epochs = loss_df[' loss'].plot()
fig = plt.500_padding_30_epochs.get_figure()
fig.savefig('loss_plt_500_padding_30_epochs.pdf')

Index(['Iter', ' loss'], dtype='object')

/anaconda3/lib/python3.7/site-packages/sklearn/metrics/classification.py:1272: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use zero_division parameter to control this behavior.
  warn_prf(average, modifier, msg_start, len(result))

              precision    recall  f1-score   support

             1         0.66         1.00         0.80         1855
             2         0.00         0.00         0.00          949

 accuracy         0.33         0.50         0.40         2804
 macro avg         0.33         0.50         0.40         2804
weighted avg         0.44         0.66         0.53         2804

Index(['Iter', ' loss'], dtype='object')

In [88]: data_test = pd.read_csv('transferlearning-dl-spring2020/test.csv')

In [89]: from gensim.utils import simple_preprocess

# Tokenize the text column to get the new column 'tokenized_text'
data_test['tokenized_text'] = [simple_preprocess(line, deacc=True) for line in data_test['text']]
print(data_test['tokenized_text'].head(10))

Out [90]: 0 [user, user, go, home, you, re, drunk, user, m...
1 [user, user, oh, noe, tough, shit]
2 [user, canada, doesn, need, another, cuck, we...
3 [user, user, user, it, should, scare, every, a...
4 [user, user, user, user, lol, throwing, the, b...
5 [user, user, you, are, correct]
6 [user, user, kind, of, like, when, conservativ...
7 [the, only, thing, the, democrat, have, is, ly...
8 [user, user, user, user, user, user, user, use...
9 [user, user, user, user, user, user, user, use...
Name: tokenized_text, dtype: object

In [90]: from gensim.parsing.porter import PorterStemmer

porter_stemmer = PorterStemmer()
# Get the stemmed tokens
data_test['stemmed_tokens'] = [[porter_stemmer.stem(word) for word in tokens] for tokens in data_test['tokenized_text']]
print(data_test['stemmed_tokens'].head(10))

Out [90]: 0 [user, user, go, home, you, re, drunk, user, m...
1 [user, user, oh, noe, tough, shit]
2 [user, canada, doesn, need, another, cuck, we...
3 [user, user, user, it, should, scare, every, a...
4 [user, user, user, user, lol, throw, the, bull...
5 [user, user, you, ar, correct]
6 [user, user, kind, of, like, when, conserv, we...
7 [the, onli, thing, the, democrat, have, is, ly...
8 [user, user, user, user, user, user, user, use...
9 [user, user, user, user, user, user, user, use...
Name: stemmed_tokens, dtype: object

In [92]: from sklearn.model_selection import train_test_split

# Train Test Split Function
def split_train_test(top_data_df_small, test_size=1, shuffle_state=True):
    X_train, X_test, Y_train, Y_test = train_test_split(top_data_df_small[['text', 'stemmed_tokens']],
                                                        top_data_df_small[['target']],
                                                        shuffle=shuffle_state,
                                                        test_size=test_size,
                                                        random_state=15)

    print("Value counts for Train sentiments")
    print(Y_train.value_counts())
    print("Value counts for Test sentiments")
    print(Y_test.value_counts())
    print(type(X_train))
    print(type(Y_train))
    X_train = X_train.reset_index()
    X_test = X_test.reset_index()
    Y_train = Y_train.to_frame()
    Y_train = Y_train.reset_index()
    Y_test = Y_test.to_frame()
    Y_test = Y_test.reset_index()
    print(X_train.head())
    return X_train, X_test, Y_train, Y_test

# Call the train test split
X_test = data_test[['text', 'stemmed_tokens']]

In [93]: import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim

import torch
# Use cuda if present
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print("Device available for running: ")
print(device)

Device available for running:
cpu

In [94]: from gensim.models import Word2Vec

size = 300
window = 3
min_count = 1
workers = 3
sg = 1

# Function to train word2vec model
def make_word2vec_model(top_data_df_small, padding=True, sg=1, min_count=1, size=500, workers=3, window=3):
    if padding:
        temp_df = (top_data_df_small[['stemmed_tokens']]).values
        temp_df = list(temp_df)
        temp_df.append(['pad'])
        word2vec_file = 'word2vec.' + str(size) + '_PAD.model'
    else:
        temp_df = top_data_df_small[['stemmed_tokens']]
        word2vec_file = 'word2vec.' + str(size) + '.model'
    w2v_model = Word2Vec(temp_df, min_count = min_count, size = size, workers = workers, window = window, sg = sg)

    w2v_model.save(word2vec_file)
    return w2v_model, word2vec_file

# Train Word2Vec model
w2_model, word2vec_file = make_word2vec_model(X_test, padding=True, sg=sg, min_count=min_count, size=size, workers=workers, window=window)

3894

In [97]: max_sen_len = X_test.stemmed_tokens.map(len).max()
padding_idx = w2_model.wv.vocab['pad'].index
def make_word2vec_vector_cnn(sentence):
    padded_x = [padding_idx for i in range(max_sen_len)]
    i = 0
    for word in sentence:
        if word not in w2_model.wv.vocab:
            padded_x[i] = 0
            print(word)
        else:
            padded_x[i] = w2_model.wv.vocab[word].index
            i += 1
    return torch.tensor(padded_x, dtype=torch.long, device=device).view(1, -1)

from sklearn.metrics import classification_report
bow_cnn_predictions = []
original_labels_cnn_bow = []
cnn_model.eval()
# loss_df = pd.read_csv('cnn_class_big_loss_with_padding.csv')
#print(loss_df.columns)
# loss_df.plot('loss')
for torch.no_grad():
    for index, row in X_test.iterrows():
        bow_vec = make_word2vec_vector_cnn(row['stemmed_tokens'])
        probs = cnn_model(bow_vec)
        _, predicted = torch.max(probs.data, 1)
        bow_cnn_predictions.append(predicted.cpu().numpy()[0])
        original_labels_cnn_bow.append(make_target(Y_test['target'][index]).cpu().numpy()[0])

#print(classification_report(original_labels_cnn_bow, bow_cnn_predictions))
#loss_file_name = 'cnn_class_big_loss_with_padding.csv'
#loss_df = pd.read_csv(loss_file_name)
#print(loss_df.columns)
# plt.500_padding_30_epochs = loss_df[' loss'].plot()
# fig = plt.500_padding_30_epochs.get_figure()
# fig.savefig('loss_plt_500_padding_30_epochs.pdf')
```



```
In [127]: # LSTM with Dropout for sequence classification in the IMDB dataset
import numpy
from keras.datasets import imdb
from keras.models import Sequential
from keras.layers import Dense
from keras.layers import LSTM
from keras.layers import Dropout
from keras.layers import Embedding
from keras.preprocessing import sequence
# fix random seed for reproducibility
numpy.random.seed(7)
# load the dataset but only keep the top n words, zero the rest
top_words = 5000

X = data.text
Y = data.target
le = LabelEncoder()
Y = le.fit_transform(Y)
Y = Y.reshape(-1,1)

X_train,X_test,Y_train,Y_test = train_test_split(X,Y,test_size=0.2)

max_words = 1000
max_len = 200
tok = Tokenizer(num_words=max_words)
tok.fit_on_texts(X_train)
sequences = tok.texts_to_sequences(X_train)
X_train = sequence.pad_sequences(sequences,maxlen=max_len)
print(X_train.shape)
print(Y_train.shape)

# create the model
embedding_vector_length = 32
model = Sequential()
model.add(Embedding(top_words, embedding_vector_length, input_length=max_len))
model.add(Dropout(0.2))
model.add(LSTM(100))
model.add(Dropout(0.2))
model.add(Dense(1, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
print(model.summary())
model.fit(X_train, y_train, epochs=20, batch_size=64)
# Final evaluation of the model
scores = model.evaluate(X_test, y_test, verbose=0)
print("Accuracy: %.2f%%" % (scores[1]*100))

(7476, 200)
(7476,)
Model: "sequential_34"

Layer (type) Output Shape Param #
-----
embedding_25 (Embedding) (None, 200, 32) 160000
dropout_6 (Dropout) (None, 200, 32) 0
lstm_32 (LSTM) (None, 100) 53200
dropout_7 (Dropout) (None, 1) 0
dense_28 (Dense) (None, 1) 101
-----
Total params: 213,301
Trainable params: 213,301
Non-trainable params: 0

None
Epoch 1/20
7476/7476 [=====] - 113s 15ms/step - loss: 0.6294 - accuracy: 0.6664
Epoch 2/20
7476/7476 [=====] - 109s 15ms/step - loss: 0.5412 - accuracy: 0.7330
Epoch 3/20
7476/7476 [=====] - 108s 15ms/step - loss: 0.4781 - accuracy: 0.7824
Epoch 4/20
7476/7476 [=====] - 109s 15ms/step - loss: 0.4582 - accuracy: 0.7956
Epoch 5/20
7476/7476 [=====] - 114s 15ms/step - loss: 0.4473 - accuracy: 0.8028
Epoch 6/20
7476/7476 [=====] - 109s 14ms/step - loss: 0.4334 - accuracy: 0.8089
Epoch 7/20
7476/7476 [=====] - 123s 15ms/step - loss: 0.4216 - accuracy: 0.8131
Epoch 8/20
7476/7476 [=====] - 122s 16ms/step - loss: 0.4085 - accuracy: 0.8178
Epoch 9/20
7476/7476 [=====] - 118s 16ms/step - loss: 0.4059 - accuracy: 0.8202
Epoch 10/20
7476/7476 [=====] - 127s 16ms/step - loss: 0.3894 - accuracy: 0.8287
Epoch 11/20
7476/7476 [=====] - 127s 15ms/step - loss: 0.3802 - accuracy: 0.8300
Epoch 12/20
7476/7476 [=====] - 129s 17ms/step - loss: 0.3721 - accuracy: 0.8356
Epoch 13/20
7476/7476 [=====] - 122s 16ms/step - loss: 0.3612 - accuracy: 0.8426
Epoch 14/20
7476/7476 [=====] - 123s 16ms/step - loss: 0.3537 - accuracy: 0.8467
Epoch 15/20
7476/7476 [=====] - 117s 16ms/step - loss: 0.3493 - accuracy: 0.8487
Epoch 16/20
7476/7476 [=====] - 121s 16ms/step - loss: 0.3281 - accuracy: 0.8559
Epoch 17/20
7476/7476 [=====] - 116s 15ms/step - loss: 0.3257 - accuracy: 0.8598
Epoch 18/20
7476/7476 [=====] - 124s 17ms/step - loss: 0.3156 - accuracy: 0.8666
Epoch 19/20
7476/7476 [=====] - 123s 16ms/step - loss: 0.3006 - accuracy: 0.8713
Epoch 20/20
7476/7476 [=====] - 113s 15ms/step - loss: 0.3016 - accuracy: 0.8732

-----
ValueError: Traceback (most recent call last)
<ipython-input-127-b855e3d580bc> in <module>
    43 model.fit(X_train, y_train, epochs=20, batch_size=64)
    44 # Final evaluation of the model
--> 45 scores = model.evaluate(X_test, y_test, verbose=0)
    46 print("Accuracy: %.2f%%" % (scores[1]*100))

/anaconda3/lib/python3.7/site-packages/keras/engine/training.py in evaluate(self, x, y, batch_size,
sample_weight, steps, callbacks, max_queue_size, workers, use_multiprocessing)
    1347 x, y,
    1348 sample_weight=sample_weight,
-> 1349 batch_size=batch_size)
    1350 # Prepare inputs, delegate logic to 'test_loop'.
    1351 if self._uses_dynamic_learning_phase():

/anaconda3/lib/python3.7/site-packages/keras/engine/training.py in _standardize_user_data(self, x, y,
sample_weight, class_weight, check_array_lengths, batch_size)
    577 feed_input_shapes,
    578 check_batch_axis=False, # Don't enforce the batch size.
--> 579 exception_prefix='input')
    580
    581 if y is not None:

/anaconda3/lib/python3.7/site-packages/keras/engine/training_utils.py in standardize_input_data(data,
names, shapes, check_batch_axis, exception_prefix)
    143 '': expected ' + names[i] + ' to have shape ' +
    144 str(shape) + ' but got array with shape ' +
--> 145 str(data.shape))
    146 return data
    147

ValueError: Error when checking input: expected embedding_25_input to have shape (200,) but got array
with shape (1,)
```

```
In [128]: #tok.fit_on_texts(X)
sequences = tok.texts_to_sequences(X_test)
X_test = sequence.pad_sequences(sequences,maxlen=max_len)

scores = model.evaluate(X_test, Y_test, verbose=0)
print("Accuracy: %.2f%%" % (scores[1]*100))

test_sequences = tok.texts_to_sequences(data_test['text'])
test_sequences_matrix = sequence.pad_sequences(test_sequences,maxlen=max_len)

test_sequences_matrix.shape
pred = model.predict(test_sequences_matrix)

predicted = []
for i in range(len(pred)):
    if pred[i] >= 0.5:
        predicted += [1]
    else:
        predicted += [0]

data_test['Target'] = predicted
submission = data_test[['Id','Target']]
submission.to_csv('lstm4.csv', index=False) # writing data to a CSV file
submission

Accuracy: 71.39%
```

```
Out[128]:
```

	id	Target
0	90194	0
1	77444	1
2	13384	0
3	54920	0
4	56117	1
5	67757	0
6	12681	0
7	12609	1
8	70380	1
9	12108	1
10	14726	0
11	74477	0
12	49845	1
13	47311	0
14	75689	0
15	84102	0
16	10697	0
17	98992	1
18	53264	1
19	54842	0
20	48995	0
21	72353	1
22	50759	0
23	14574	0
24	89119	0
25	43133	1
26	59537	0
27	32317	0
28	76680	0
29	80561	1
--	--	--
3864	54190	0
3865	28037	0
3866	77430	0
3867	75815	1
3868	87290	0
3869	89475	0
3870	43323	0
3871	37666	0
3872	77905	0
3873	83400	0
3874	84081	1
3875	71649	0
3876	20841	0
3877	90959	0
3878	32598	0
3879	43964	0
3880	97745	0
3881	86716	1
3882	21033	1
3883	66632	0
3884	28996	0
3885	64713	0
3886	63482	1
3887	11132	0
3888	87416	0
3889	90041	0
3890	98824	0
3891	95338	1
3892	67210	0
3893	46552	0
3894	rows x 2 columns	

```
In [ ]: # LSTM and Convolutional Neural Network For Sequence Classification
```

```
In [130]: # LSTM and CNN for sequence classification in the IMDB dataset
import numpy
from keras.datasets import imdb
from keras.models import Sequential
from keras.layers import Dense
from keras.layers import Conv1D
from keras.layers import LSTM
from keras.layers import Embedding
from keras.preprocessing import sequence
# fix random seed for reproducibility
numpy.random.seed(7)
# load the dataset but only keep the top n words, zero the rest
top_words = 5000

X = data.text
Y = data.target
le = LabelEncoder()
Y = le.fit_transform(Y)
Y = Y.reshape(-1,1)

X_train,X_test,Y_train,Y_test = train_test_split(X,Y,test_size=0.2)

max_words = 1000
max_len = 200
tok = Tokenizer(num_words=max_words)
tok.fit_on_texts(X_train)
sequences = tok.texts_to_sequences(X_train)
X_train = sequence.pad_sequences(sequences,maxlen=max_len)
print(X_train.shape)
print(Y_train.shape)

# create the model
embedding_vector_length = 32
model = Sequential()
model.add(Embedding(top_words, embedding_vector_length, input_length=max_len))
model.add(Conv1D(filters=32, kernel_size=3, padding='same', activation='relu'))
model.add(MaxPooling1D(pool_size=2))
model.add(LSTM(100))
model.add(Dense(1, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
print(model.summary())
model.fit(X_train, y_train, epochs=10, batch_size=64)
# Final evaluation of the model
scores = model.evaluate(X_test, y_test, verbose=0)
print("Accuracy: %.2f%%" % (scores[1]*100))

(7476, 200)
(7476,)
Model: "sequential_36"

Layer (type) Output Shape Param #
-----
embedding_27 (Embedding) (None, 200, 32) 160000
conv1d_2 (Conv1D) (None, 200, 32) 3104
max_pooling1d_2 (MaxPooling1D) (None, 100, 32) 0
lstm_34 (LSTM) (None, 100) 53200
dense_30 (Dense) (None, 1) 101
-----
Total params: 216,405
Trainable params: 216,405
Non-trainable params: 0

None
Epoch 1/10
7476/7476 [=====] - 76s 10ms/step - loss: 0.6315 - accuracy: 0.6663
Epoch 2/10
7476/7476 [=====] - 66s 9ms/step - loss: 0.5387 - accuracy: 0.7334
Epoch 3/10
7476/7476 [=====] - 66s 9ms/step - loss: 0.4650 - accuracy: 0.7875
Epoch 4/10
7476/7476 [=====] - 59s 8ms/step - loss: 0.4393 - accuracy: 0.8039
Epoch 5/10
7476/7476 [=====] - 61s 8ms/step - loss: 0.3923 - accuracy: 0.8121
Epoch 6/10
7476/7476 [=====] - 61s 8ms/step - loss: 0.3923 - accuracy: 0.8222
Epoch 7/10
7476/7476 [=====] - 58s 8ms/step - loss: 0.3610 - accuracy: 0.8404
Epoch 8/10
7476/7476 [=====] - 59s 8ms/step - loss: 0.3101 - accuracy: 0.8670
Epoch 9/10
7476/7476 [=====] - 59s 8ms/step - loss: 0.2623 - accuracy: 0.8953
Epoch 10/10
7476/7476 [=====] - 68s 9ms/step - loss: 0.2010 - accuracy: 0.9240

-----
ValueError: Traceback (most recent call last)
<ipython-input-130-2c043b74a2dd> in <module>
    43 model.fit(X_train, y_train, epochs=10, batch_size=64)
    44 # Final evaluation of the model
--> 45 scores = model.evaluate(X_test, y_test, verbose=0)
    46 print("Accuracy: %.2f%%" % (scores[1]*100))

/anaconda3/lib/python3.7/site-packages/keras/engine/training.py in evaluate(self, x, y, batch_size,
sample_weight, steps, callbacks, max_queue_size, workers, use_multiprocessing)
    1347 x, y,
    1348 sample_weight=sample_weight,
-> 1349 batch_size=batch_size)
    1350 # Prepare inputs, delegate logic to 'test_loop'.
    1351 if self._uses_dynamic_learning_phase():

/anaconda3/lib/python3.7/site-packages/keras/engine/training.py in _standardize_user_data(self, x, y,
sample_weight, class_weight, check_array_lengths, batch_size)
    577 feed_input_shapes,
    578 check_batch_axis=False, # Don't enforce the batch size.
--> 579 exception_prefix='input')
    580
    581 if y is not None:

/anaconda3/lib/python3.7/site-packages/keras/engine/training_utils.py in standardize_input_data(data,
names, shapes, check_batch_axis, exception_prefix)
    143 '': expected ' + names[i] + ' to have shape ' +
    144 str(shape) + ' but got array with shape ' +
--> 145 str(data.shape))
    146 return data
    147

ValueError: Error when checking input: expected embedding_27_input to have shape (200,) but got array
with shape (1,)
```

```
In [131]: #tok.fit_on_texts(X)
sequences = tok.texts_to_sequences(X_test)
X_test = sequence.pad_sequences(sequences,maxlen=max_len)

scores = model.evaluate(X_test, Y_test, verbose=0)
print("Accuracy: %.2f%%" % (scores[1]*100))

test_sequences = tok.texts_to_sequences(data_test['text'])
test_sequences_matrix = sequence.pad_sequences(test_sequences,maxlen=max_len)

test_sequences_matrix.shape
pred = model.predict(test_sequences_matrix)

predicted = []
for i in range(len(pred)):
    if pred[i] >= 0.5:
        predicted += [1]
    else:
        predicted += [0]

data_test['Target'] = predicted
submission = data_test[['Id','Target']]
submission.to_csv('lstm5.csv', index=False) # writing data to a CSV file
submission

Accuracy: 67.70%
```

	id	Target
0	90194	1
1	77444	1
2	13384	0
3	54920	1
4	56117	1
5	67757	0
6	12681	0
7	12609	0
8	70380	0
9	12108	1
10	14726	0
11	74477	0
12	49845	1
13	47311	1
14	75689	0
15	84102	0
16	10607	0
17	98992	1
18	53264	1
19	54842	0
20	48995	1
21	72353	1
22	50759	0
23	14574	0
24	89119	0
25	43133	1
26	59537	0
27	32317	0
28	76680	1
29	80561	1
--	--	--
3864	54190	0
3865	28037	0
3866	77430	0
3867	75815	1
3868	87290	0
3869	89475	0
3870	43323	1
3871	37666	1
3872	77905	1
3873	83400	0
3875	71649	0
3876	20841	0
3877	90959	0
3878	32598	0
3879	43964	0
3880	97745	1
3881	86716	1
3882	21033	1
3883	66632	0
3884	28996	0
3885	64713	1
3886	63482	0
3887	11132	0
3888	87416	0
3889	90041	1
3890	98824	0
3891	95338	1
3892	67210	0
3893	46552	0
3894	rows x 2 columns	

```
In [ ]: # LSTM with Dropout and Random for sequence classification in the IMDB dataset
import numpy
from numpy import array
from numpy import cumsum
from keras.models import Sequential
from keras.layers import Dense
from keras.layers import LSTM
from keras.layers import Embedding
from keras.layers import Bidirectional

X = data.text
Y = data.target
le = LabelEncoder()
Y = le.fit_transform(Y)
Y = Y.reshape(-1,1)

X_train,X_test,Y_train,Y_test = train_test_split(X,Y,test_size=0.2)

max_words = 1000
max_len = 200
tok = Tokenizer(num_words=max_words)
tok.fit_on_texts(X_train)
sequences = tok.texts_to_sequences(X_train)
X_train = sequence.pad_sequences(sequences,maxlen=max_len)
print(X_train.shape)
print(Y_train.shape)

# define LSTM
model = Sequential()
model.add(Bidirectional(LSTM(20, return_sequences=True), input_shape=(max_len, 1)))
model.add(LSTM(20, return_sequences=True, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
# train LSTM
model.fit(X, y, epochs=1, batch_size=1, verbose=2)
# evaluate LSTM
X,y = get_sequence(n,timesteps)
yhat = model.predict_classes(X, verbose=0)
for i in range(n_timesteps):
    print("Expected:", y[i0, 1], "Predicted:", yhat[i0, i])
```

```
In [ ]: #tok.fit_on_texts(X)
sequences = tok.texts_to_sequences(X_test)
X_test = sequence.pad_sequences(sequences,maxlen=max_len)

scores = model.evaluate(X_test, Y_test, verbose=0)
print("Accuracy: %.2f%%" % (scores[1]*100))

test_sequences = tok.texts_to_sequences(data_test['text'])
test_sequences_matrix = sequence.pad_sequences(test_sequences,maxlen=max_len)

test_sequences_matrix.shape
pred = model.predict(test_sequences_matrix)

predicted = []
for i in range(len(pred)):
    if pred[i] >= 0.5:
        predicted += [1]
    else:
        predicted += [0]

data_test['Target'] = predicted
submission = data_test[['Id','Target']]
submission.to_csv('lstm2.csv', index=False) # writing data to a CSV file
submission
```