

```

In [1]: #Pretrained word2vec embeddings

In [11]: from torchtext import data
import torchtext
import torch
import torch.nn as nn
import torch.nn.functional as F
import random
import numpy as np
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

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

data.head()

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

train_data, valid_data, ytrain, yvalid = train_test_split(questions, labels,
                                                         random_state=42,
                                                         test_size=0.2)

train_data.head()

Out[11]: 251      0USER All you need to know is he is empty inside
3720      0USER I. AM. ROAD. 21 URL
8376      0USER Go roger I quit watching anyway nrl is o..
4471      0USER Yoo our dogs should totally fuck
2835      0USER HER He is a troll. Not open to facts
Name: text, dtype: object

In [12]: data.head()

Out[12]:      text      target
0      86426      0USER She should ask a few native Americans wh... 1
1      16820      0Amazon investigating Chinese employees wh... 0
2      62688      0USER Someone should've taken this piece of sh... 1
3      43605      0USER Obama wanted liberals amp illegals move i... 0
4      97670      0USER Liberals are all Kokoo!! 1

In [13]: def remove_pattern(input_txt, pattern):
# re.findall(pattern, input_txt)
for i in range(len(input_txt)):
input_txt = re.sub(i, '', input_txt)

return input_txt

data['text'] = np.vectorize(remove_pattern)(data['text'], "@(\\w|\\W)@")

data['text'] = data['text'].str.replace("@(\\w|\\W)@", " ")

data['text'] = data['text'].str.replace("#", "")

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

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

tokenized_tweet.head()

len(tokenized_tweet)

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

In [14]: tokenized_tweet.apply(lambda x: x.split())

In [15]: # 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(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_df[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 [16]: from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import f1_score

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[16]: 0.6691278341911253

In [17]: data_test = pd.read_csv('transferlearning-dl-spring2020/test.csv')

In [18]: data_test['text'] = np.vectorize(remove_pattern)(data_test['text'], "@(\\w|\\W)@")

data_test['text'] = data_test['text'].str.replace("@(\\w|\\W)@", " ")

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

In [19]: 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

X_train_test, X_test_test, y_train_test, y_test_test = train_test_split(wordvec_df_test, labels, test_size=0.2, random_state=1)

In [20]: from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import f1_score

lreg = LogisticRegression(random_state = 10)
lreg.fit(X_train_test, y_train_test) # training the model

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

prediction_int_test = prediction_test.astype(np.int)

f1_score(y_test_test, prediction_test, average = "macro") # calculating f1 score

Out[20]: 0.5291078614322713

In [21]: # CNN for text classification

In [22]: top_data_df_small = data

In [23]: from gensim.utils import simple_preprocess
# tokenize the text column to get the new column 'tokenized_text'
top_data_df_small['tokenized_text'] = [simple_preprocess(line, deacc=True) for line in top_data_df_small['text']]
print(top_data_df_small['tokenized_text'].head(10))

0 [she, should, ask, few, native, americans, wha...
1 [amazon, investigating, chinese, employees, wh...
2 [someone, should, vetaken, this, piece, shit, ...
3 [obama,wanted, liberals, amp, illegals, move, ...
4 [liberals, are, all, kokoo]
5 [was, literally, just, talking, about, this, l...
6 [buy, more, icecream]
7 [not, fault, you, support, gun, control]
8 [what, the, difference, between, kavanaugh, an...
9 [why, are, lying, corrupt, traitor, nobody, wa...
Name: tokenized_text, dtype: object

In [24]: from gensim.parsing.porter import PorterStemmer
porter_stemmer = PorterStemmer()
# Get the stemmed tokens
top_data_df_small['stemmed_tokens'] = [porter_stemmer.stem(word) for word in tokens] for tokens in top_data_df_small['tokenized_text']
print(top_data_df_small['stemmed_tokens'].head(10))

0 [she, should, ask, few, nativ, american, what...
1 [amazon, investig, chines, employe, who, ar, s...
2 [someon, shoud, vetaken, this, piec, shit, vol...
3 [obama, want, liber, amp, illeg, move, into, r...
4 [liberal, ar, all, kokoo]
5 [wa, liter, just, talk, about, thi, lol, all...
6 [buy, more, icecream]
7 [not, fault, you, support, gun, control]
8 [what, the, differ, between, kavanaugh, and, o...
9 [you, ar, ly, corrupt, traitor, nobodi, want, ...
Name: stemmed_tokens, dtype: object

In [25]: from sklearn.model_selection import train_test_split
# Train test split function
def split_train_test(top_data_df_small, test_size=0.3, shuffle_state=True):
X_train, X_test, y_train, y_test = train_test_split(top_data_df_small[['text', 'stemmed_tokens']]
                                                    , 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_test = Y_test.to_frame()
print(X_train.head())
return X_train, X_test, y_train, y_test

# Call the train test split
X_train, X_test, Y_train, Y_test = split_train_test(top_data_df_small)

Value counts for Train sentiments
0 4365
1 2177
Name: target, dtype: int64
Value counts for Test sentiments
0 1855
1 949
Name: target, dtype: int64
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
Index
0 8360      0      least you are honest, lmao
1 3722      0      All that shit was hers URL
2 3322      0      Cutie connection authentic
3 4232      0      What interview Evanne honest and right from th...
4 1678      0      Warriorscoach fundraising for gun control Brad...

return X_train, X_test, Y_train, Y_test

X_train, X_test, Y_train, Y_test = split_train_test(top_data_df_small)

Value counts for Train sentiments
0 4365
1 2177
Name: target, dtype: int64
Value counts for Test sentiments
0 1855
1 949
Name: target, dtype: int64
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
Index
0 8360      0      least you are honest, lmao
1 3722      0      All that shit was hers URL
2 3322      0      Cutie connection authentic
3 4232      0      What interview Evanne honest and right from th...
4 1678      0      Warriorscoach fundraising for gun control Brad...

return X_train, X_test, Y_train, Y_test

X_train, X_test, Y_train, Y_test = split_train_test(top_data_df_small)

Value counts for Train sentiments
0 4365
1 2177
Name: target, dtype: int64
Value counts for Test sentiments
0 1855
1 949
Name: target, dtype: int64
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
Index
0 8360      0      least you are honest, lmao
1 3722      0      All that shit was hers URL
2 3322      0      Cutie connection authentic
3 4232      0      What interview Evanne honest and right from th...
4 1678      0      Warriorscoach fundraising for gun control Brad...

return X_train, X_test, Y_train, Y_test

X_train, X_test, Y_train, Y_test = split_train_test(top_data_df_small)

Value counts for Train sentiments
0 4365
1 2177
Name: target, dtype: int64
Value counts for Test sentiments
0 1855
1 949
Name: target, dtype: int64
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
Index
0 8360      0      least you are honest, lmao
1 3722      0      All that shit was hers URL
2 3322      0      Cutie connection authentic
3 4232      0      What interview Evanne honest and right from th...
4 1678      0      Warriorscoach fundraising for gun control Brad...

return X_train, X_test, Y_train, Y_test

X_train, X_test, Y_train, Y_test = split_train_test(top_data_df_small)

Value counts for Train sentiments
0 4365
1 2177
Name: target, dtype: int64
Value counts for Test sentiments
0 1855
1 949
Name: target, dtype: int64
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
Index
0 8360      0      least you are honest, lmao
1 3722      0      All that shit was hers URL
2 3322      0      Cutie connection authentic
3 4232      0      What interview Evanne honest and right from th...
4 1678      0      Warriorscoach fundraising for gun control Brad...

return X_train, X_test, Y_train, Y_test

X_train, X_test, Y_train, Y_test = split_train_test(top_data_df_small)

Value counts for Train sentiments
0 4365
1 2177
Name: target, dtype: int64
Value counts for Test sentiments
0 1855
1 949
Name: target, dtype: int64
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
Index
0 8360      0      least you are honest, lmao
1 3722      0      All that shit was hers URL
2 3322      0      Cutie connection authentic
3 4232      0      What interview Evanne honest and right from th...
4 1678      0      Warriorscoach fundraising for gun control Brad...

return X_train, X_test, Y_train, Y_test

X_train, X_test, Y_train, Y_test = split_train_test(top_data_df_small)

Value counts for Train sentiments
0 4365
1 2177
Name: target, dtype: int64
Value counts for Test sentiments
0 1855
1 949
Name: target, dtype: int64
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
Index
0 8360      0      least you are honest, lmao
1 3722      0      All that shit was hers URL
2 3322      0      Cutie connection authentic
3 4232      0      What interview Evanne honest and right from th...
4 1678      0      Warriorscoach fundraising for gun control Brad...

return X_train, X_test, Y_train, Y_test

X_train, X_test, Y_train, Y_test = split_train_test(top_data_df_small)

Value counts for Train sentiments
0 4365
1 2177
Name: target, dtype: int64
Value counts for Test sentiments
0 1855
1 949
Name: target, dtype: int64
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
Index
0 8360      0      least you are honest, lmao
1 3722      0      All that shit was hers URL
2 3322      0      Cutie connection authentic
3 4232      0      What interview Evanne honest and right from th...
4 1678      0      Warriorscoach fundraising for gun control Brad...

return X_train, X_test, Y_train, Y_test

X_train, X_test, Y_train, Y_test = split_train_test(top_data_df_small)

Value counts for Train sentiments
0 4365
1 2177
Name: target, dtype: int64
Value counts for Test sentiments
0 1855
1 949
Name: target, dtype: int64
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
Index
0 8360      0      least you are honest, lmao
1 3722     
```



```
[43]: import torch

#Reproducing same results
SEED = 2019

#torch
torch.manual_seed(SEED)

#Cuda algorithms
torch.backends.cudnn.deterministic = True

In [46]: #loading custom dataset
training_data=data

#print preprocessed text
print(vars(training_data[0:1]))

{'is_copy': <weakref at 0x1a30693c98; to 'dataFrame' at 0x1a3055a20>, '_data': BlockManager
Items: Index(['id', 'text', 'target'], dtype=object)
Axis 1: RangeIndex(start=0, stop=1, step=1)
Axis 0: slice(0, 4, 2), 2 x 1, dtype: int64
IntBlock: slice(1, 2, 1), 1 x 1, dtype: object, '_item_cache': {}

In [49]: import random

train_data,X_test, valid_data,Y_test = train_test_split(training_data['text'],training_data['target'],t
est_size=0.2)

In [50]: # LSTM for text classification

In [118]: # LSTM 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.embeddings 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)
#X_train = sequence.pad_sequences(X_train, maxlen=max_review_length)
#X_test = sequence.pad_sequences(X_test, maxlen=max_review_length)
# create the model
embedding_vecor_length = 32
model = Sequential()
model.add(Embedding(top_words, embedding_vecor_length, input_length=max_len))
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)
(25009,)
Model: "sequential_31"

Layer (type) Output Shape Param #
-----
embedding_22 (Embedding) (None, 200, 32) 160000
lstm_29 (LSTM) (None, 100) 53200
dense_25 (Dense) (None, 1) 101
-----
Total params: 213,301
Trainable params: 213,301
Non-trainable params: 0

None
Epoch 1/10
7476/7476 [=====] - 120s 16ms/step - loss: 0.6291 - accuracy: 0.6651
Epoch 2/10
7476/7476 [=====] - 112s 15ms/step - loss: 0.5440 - accuracy: 0.7289
Epoch 3/10
7476/7476 [=====] - 113s 15ms/step - loss: 0.4807 - accuracy: 0.7746
Epoch 4/10
7476/7476 [=====] - 125s 17ms/step - loss: 0.4601 - accuracy: 0.7925
Epoch 5/10
7476/7476 [=====] - 111s 15ms/step - loss: 0.4495 - accuracy: 0.7968
Epoch 6/10
7476/7476 [=====] - 105s 14ms/step - loss: 0.4337 - accuracy: 0.8091
Epoch 7/10
7476/7476 [=====] - 106s 14ms/step - loss: 0.4179 - accuracy: 0.8146
Epoch 8/10
7476/7476 [=====] - 105s 14ms/step - loss: 0.4056 - accuracy: 0.8234
Epoch 9/10
7476/7476 [=====] - 106s 14ms/step - loss: 0.3928 - accuracy: 0.8261
Epoch 10/10
7476/7476 [=====] - 112s 15ms/step - loss: 0.3918 - accuracy: 0.8277

-----
ValueError Traceback (most recent call last)
<ipython-input-118-32a60363e4e4> in <module>
    146 model.fit(X_train, Y_train, epochs=10, batch_size=64)
    147 # Final evaluation of the model
--> 48 scores = model.evaluate(X_test, Y_test, verbose=0)
    49 print("Accuracy: %.2f%%" % (scores[1]*100))

/anaconda3/lib/python3.7/site-packages/keras/engine/training.py in evaluate(self, x, y, batch_size, v
erbos, sample_weight, steps, callbacks, max_queue_size, workers, use_multiprocessing)
    1347         sample_weight=sample_weight,
    1348         batch_size=batch_size)
    1349 # Prepare inputs, delegate logic to 'test_loop'.
    1350 if self._uses_dynamic_learning_phase():
    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, expected_shape, names_to_have_shape)
    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_22_input to have shape (200,) but got array
with shape (1,)
```

```
In [120]: #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))

Accuracy: 71.98%
```

```
In [ ]: 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
```

```
In [891]: # LSTM For Text Classification With Dropout
```

```
In [ ]: # 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.embeddings 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_vecor_length = 32
model = Sequential()
model.add(Embedding(top_words, embedding_vecor_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, 200)
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, 100) 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: 0.7956
Epoch 4/20
7476/7476 [=====] - 110s 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 [=====] - 108s 14ms/step - loss: 0.4334 - accuracy: 0.8089
Epoch 7/20
7476/7476 [=====] - 109s 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
6336/7476 [=====,.....] - ETA: 17s - loss: 0.4055 - accuracy: 0.8212

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
```

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

```
In [ ]: # 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 LSTM
from keras.layers.convolutional import Conv1D
from keras.layers.convolutional import MaxPooling1D
from keras.layers.embeddings 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_train,X_test,Y_train,Y_test = train_test_split(data['text'],data['target'],test_size=0.2)
# truncate and pad input sequences
max_review_length = 500
X_train = sequence.pad_sequences(X_train, maxlen=max_review_length)
X_test = sequence.pad_sequences(X_test, maxlen=max_review_length)
# create the model
embedding_vecor_length = 32
model = Sequential()
model.add(Embedding(top_words, embedding_vecor_length, input_length=max_review_length))
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=3, 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, 200)
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, 100) 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: 0.7956
Epoch 4/20
7476/7476 [=====] - 110s 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 [=====] - 108s 14ms/step - loss: 0.4334 - accuracy: 0.8089
Epoch 7/20
7476/7476 [=====] - 109s 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
6336/7476 [=====,.....] - ETA: 17s - loss: 0.4055 - accuracy: 0.8212
```

```
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
```

```
In [ ]: from random import random
from numpy import array
from numpy import cumsum
from keras.models import Sequential
from keras.layers import LSTM
from keras.layers import Dense
from keras.layers import TimeDistributed
from keras.layers import Bidirectional

# create a sequence classification instance
def get_sequence(n_timesteps):
    # create a sequence of random numbers in [0,1)
    X = array([random() for _ in range(n_timesteps)])
    # calculate cut-off value to change class values
    limit = n_timesteps/4.0
    # determine the class outcome for each item in cumulative sequence
    y = array([0 if x < limit else 1 for x in cumsum(X)])
    # reshape input and output data to be suitable for LSTMs
    X = X.reshape(1, n_timesteps, 1)
    y = y.reshape(1, n_timesteps, 1)
    return X, y

# define problem properties
n_timesteps = 10
model = Sequential()
model.add(Bidirectional(LSTM(20, return_sequences=True), input_shape=(n_timesteps, 1)))
model.add(TimeDistributed(Dense(1, activation='sigmoid')))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])

for epoch in range(1000):
    # generate new random sequence
    X,y = get_sequence(n_timesteps)
    # fit model for one epoch on this sequence
    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, i], "Predicted:", yhat[0, 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
```