

Natural Language Processing Using TensorFlow

Word Based Embedding using API's.

```
from tensorflow.keras.preprocessing.text import Tokenizer
```

Sentences = ['I love my dog',

'I love my cat',

'You love my dog!']

tokenizer = Tokenizer(num_words=100)

tokenizer.fit_on_texts(Sentences)

word_index = tokenizer.word_index

print(word_index)

specifies max no.
of words - 1 to
keep when gene-
rating sequences.

Output: { 'i': 3, 'my': 2, 'dog': 6, 'love': 1, 'cat': 5,
'dog': 4 }

Represent text in numerical representations which can
be used to train neural network.

Take each word from list of sentences and assign
integer using fit_on_texts(). Get the result by
word_index property.

More frequent words → lower index.

By default → punctuation ignored & converted to
lower. Changed through filters & lower
arguments of Tokenizer.

Text To Sequence

out of vocab

Starter sentences: ['I love my dog',

'I love my cat', 'You love my dog!',]

'Do you think my dog is amazing?']

tokenizer = Tokenizer(num_words=100), 00V-TOKEN
"COOV1")

tokenizer.fit_on_texts(sentences)

word_index = tokenizer.word_index

sequences = tokenizer.texts_to_sequences(sentences)

print(word_index)

print(sequences)

test_data = ['i really love my dog',

'my dog loves my mandate']

test_seq = tokenizer.texts_to_sequences(test_data)

print(test_seq)

To Add Padding:

from tensorflow.keras.preprocessing.sequence import
pad_sequences

padded = pad_sequences(sequences)

For padding after the sentence

padded = pad_sequences(sequences, padding='post')

padded = pad_sequences(sequences, padding='post',

~~truncating='post'~~, maxlen=5)

If losing info from pre, have sentences of
adding post makes it length not more than 5
lose from end of the to be padded.

December 2024							January 2025						
M	T	W	T	F	S	S	M	T	W	T	F	S	S
1	2	3	4	5	6	7	8	9	10	11	12	13	14
15	16	17	18	19	20	21	22	23	24	25	26	27	28
29	30	31	1	2	3	4	5	6	7	8	9	10	11

Tuesday

24

DEC 2024

359-007 • WK 52

NLP Using TensorFlow

M	T	W	T	F	S	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

M	T	W	T	F	S	S
29	30	31	1	2	3	4
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

using datasets.

Monday

28

OCT 2024

→ sarcasm in news headlines.

loading the data:

import json

with open ("sarcasm.json", 'r') as f:

datastore = json.load(f)

sentences = []

labels = []

urls = []

for item in datastore:

sentences.append(item['headline'])

labels.append(item['is-sarcastic'])

urls.append(item['article-link'])

generate word index

tokenizer = Tokenizer.fit_on_texts(sentences)

word_index = tokenizer.word_index

Create sequences & do padding.

NOTES

Word Embeddings29 tfd
Tuesday

	M	T	W	T	F	S	S		M	T	W	T	F	S	S
September 2024							1		1	2	3	4	5	6	
	30	1	2	3	4	5	6	7	7	8	9	10	11	12	13
	2	3	4	5	6	7	8	9	14	15	16	17	18	19	20
	9	10	11	12	13	14	15	16	21	22	23	24	25	26	27
	16	17	18	19	20	21	22	23	28	29	30	31			
	23	24	25	26	27	28	29								
October 2024															

import tensorflow as tf
9 import tensorflow_datasets as tfds
vimb, info = tfds.load("imdb_reviews", with_info=True, as_supervised=True)

10 import numpy as np
11 train_data, test_data = vimb['train'], vimb['test']

12 for training_sentences = []
13 training_labels = []
1 test_sentences = []
2 test_labels = []
3
for s, l in train_data:
4 for training_sentences.append(s.numpy().decode('utf-8'))
5 training_labels.append(l.numpy())
6
for s, l in test_data:
7 test_sentences.append(s.numpy().decode('utf-8'))
8 test_labels.append(l.numpy())
9
training_labels_final = np.array(training_labels)
10 test_labels_final = np.array(test_labels)

hyperparameters
NOTES vocab_size = 10000
embedding_dim = 16
max_length = 120
trunc_type = 'post'
oov_tok = "<OOV>"

word embedding. vector for each word with their associated sentiment.

	M	T	W	T	F	S	S
1	1	2	3				
2	5	6	7	8	9	10	
3	12	13	14	15	16	17	
4	19	20	21	22	23	24	
5	26	27	28	29	30		

	M	T	W	T	F	S	S
30	31					1	
1	2	3	4	5	6	7	8
2	9	10	11	12	13	14	15
3	16	17	18	19	20	21	22
4	23	24	25	26	27	28	29

[30
Wednesday]

OCT 2024
30-08-2024 • WK 44

- from tensorflow.keras.preprocessing.text import Tokenizer.
- from tensorflow.keras.preprocessing.sequence import pad_sequences.

10 tokenizer = Tokenizer(num_words=vocab_size,
11 oov_tokens=0)

12 tokenizer.fit_on_texts(training_sentences)

13 word_index = tokenizer.word_index

14 sequences = tokenizer.texts_to_sequences(training_sentences)

15 padded = pad_sequences(sequences, maxlen=max_length, truncating='pre')

16 testing_sequences = tokenizer.texts_to_sequences(testing_sentences)

17 testing_padded = pad_sequences(testing_sequences, maxlen=max_length)

18 model = tf.keras.Sequential([

19 tf.keras.layers.Embedding(vocab_size,

20 embedding_dim, input_length=max_length),

21 tf.keras.layers.Flatten(),

22 tf.keras.layers.Dense(6, activation='relu'),

23 tf.keras.layers.Dense(1, activation='sigmoid'))]

24 model.compile(loss='binary_crossentropy', optimizer

25 notes = 'adam', metrics=['accuracy'])

26 model.summary()

27 num_epochs = 10

model.fit(padded, training_labels_final,
9 epochs = num_epochs,
10 validation_data = (testing_padded,
 test_labels_final))

11 e = model.layers[0]
weights = e.get_weights()[0]
12 print(weights.shape) # (10000, 16)

13 reverse_word_index = tokenizer.index_word

14 import io
15 out_v = io.open('vecs.tsv', 'w', encoding='utf-8')
16 out_m = io.open('meta.tsv', 'w', encoding='utf-8')
17 for word_num in range(1, vocab_size):
18 word = reverse_word_index[word_num]
19 embeddings = weights[word_num]
20 out_m.write(word + '\n')
21 out_v.write(' '.join([str(x) for x in embeddings])
22 + "\n")

23 out_v.close()

24 out_m.close()

//du colale

try:

from google.colab import files
except ImportError:

pass

else:

files.download('vecs.tsv')
files.download('meta.tsv')

NOTES

Building a classifier for sarcasm dataset.

Friday

	M	T	W	T	F	S	S
October 2024	1	2	3	4	5	6	
	7	8	9	10	11	12	13
	14	15	16	17	18	19	20
	21	22	23	24	25	26	27
	28	29	30	31			

	M	T	W	T	F	S	S
November 2024	4	5	6	7	8	9	3
	11	12	13	14	15	16	17
	18	19	20	21	22	23	24
	25	26	27	28	29	30	

```

import json
import tensorflow as tf
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences

# hyperparameters
vocab_size = 10000
embedding_dim = 16
max_length = 32
trunc_type = 'post'
padding_type = 'post'
oov_token = '<OOV>'
training_size = 20000

# download the sarcasm dataset.
with open('/tmp/sarcasm.json', 'r') as f:
    datastore = json.load(f)

sentences = []
labels = []

for item in datastore:
    sentences.append(item['headline'])
    labels.append(item['is_sarcastic'])

NOTES
training_sentences = sentences[0:training_size]
training_labels = labels[0:training_size]
testing_sentences = sentences[training_size:]
testing_labels = labels[training_size:]

```

	M	T	W	T	F	S	S
30	31	1	2	3	4	5	
1	2	3	4	5	6	7	8
2	3	4	5	6	7	8	9
9	10	11	12	13	14	15	16
16	17	18	19	20	21	22	23
23	24	25	26	27	28	29	30

January 2025

	M	T	W	T	F	S	S
1	2	3	4	5	6	7	8
6	7	8	9	10	11	12	13
13	14	15	16	17	18	19	20
20	21	22	23	24	25	26	27
27	28	29	30	31			

NOV 2024
307-059 • WK 44

NOV 2024

307-059 • WK 44

02

Saturday

- # tokenize training-sentences
- # sequence training-sentences & testing-sentences
- # pad both.
- # convert the labels list to np-array.
- # create usual neural network - sequential
- # use GlobalAveragePooling1D.
- # compile with binary_crossentropy
- num_epochs = 30
- model.fit (training_padded, training_labels, epochs = num_epochs, validation_data = (testing_padded, testing_labels), verbose = 2)

Plotting the results.

```
import matplotlib.pyplot as plt
def plot_graphs(history, string):
    plt.plot (history.history [string])
    " " ( " " [ " " val_ + string ] )
```

Sunday 03

```
plt.xlabel ('Epochs')
```

```
plt.ylabel (String)
```

```
plt.legend ([string, 'val_ ' + string])
```

```
plt.show()
```

NOTES

```
plot_graphs (history, 'acc')
plot_graphs (history, 'loss')
```

Validation loss found to be increasing why?

~~Progress~~04 Solution: tweak the Monday hyperparameters

vocab_size = 1000 embedding_dim = 32
 max_length = 16

in all subwords 8K using single layer LSTM

RNN - takes ordering of inputs into account.

LSTM - computes the state of current timestep & passes it on to the next timestep where the state is also updated. The process repeats until final timestep where the output computation is affected by all previous timestep states.

import tensorflow as tfds.

dataset, info = tfds.load('imdb_reviews/subwords_8K', with_info=True, as_supervised=True)

tokenizer = info.features['text'].encoder

BUFFER_SIZE = 10000

BATCH_SIZE = 256

train_data, test_data = dataset['train'], dataset['test']

train_dataset = train_data.shuffle(BUFFER_SIZE)

train_dataset = train_dataset.padded_batch(BATCH_SIZE)

test_dataset = test_data.padded_batch(BATCH_SIZE)

NOTES

M	T	W	T	F	S	S	M	T	W	F	S	S
1	2	3	4	5	6	7	1	2	3	4	5	6
8	9	10	11	12	13	14	11	12	13	14	15	16
15	16	17	18	19	20	1	18	19	20	21	22	23
22	23	24	25	26	27	29	26	27	28	29	30	31

M	T	W	T	F	S	S	M	T	W	F	S	S
4	5	6	7	8	9	3	4	5	6	7	8	9
11	12	13	14	15	16	17	18	19	20	21	22	23
18	19	20	21	22	23	24	25	26	27	28	29	30

Model Building : Swap Flatten() &

Global Average Pooling 1D ()
with LSTM Layer.

NOV 2024
310-558 - WK 45

05

nest it inside Bidirectional layer

so passing the sequence info goes both forward & backward.

embedding_dim = 64

lstm_dim = 64

dense_dim = 64

model = tf.keras.Sequential [

tf.keras.layers.Embedding(tokenizer.vocab_size,
embedding_dim),

tf.keras.layers.Bidirectional(tf.keras.layers.

LSTM(lstm_dim))

tf.keras.layers.Dense(dense_dim, activation='relu')

tf.keras.layers.Dense(1, activation='sigmoid')])

To build Multilayer LSTM

simply append another LSTM layer in your sequential model & enable return_sequences
as True.

Why?

Because LSTM layer expects a sequence input,

so if previous input is also an LSTM, then it should output a sequence as well.

Multilayer LSTM model

embedding_dim = 64

NOTES lstm1_dim = 64

lstm2_dim = 32

dense_dim = 64

model = tf.keras.Sequential [

tf.keras.layers.Embedding(tokenizer.vocab_size,
embedding_dim),

06

Wednesday

	M	T	W	T	F	S	S	M	T	W	T	F	S	S
October 2024	1	2	3	4	5	6		4	5	6	7	8	9	3
	7	8	9	10	11	12	13	11	12	13	14	15	16	10
	14	15	16	17	18	19	20	18	19	20	21	22	23	17
	21	22	23	24	25	26	27	25	26	27	28	29	30	24
	28	29	30	31										

- 9 tf. keras. layers. Bidirectional (tf. keras. layers. LSTM
 (lstm1. dim, return_sequences = True)),
 10 tf. keras. layers. Bidirectional (tf. keras. layers.
 LSTM (lstm2. dim)),
 11 tf. keras. layers. Dense (dense_dim, activation = 'tanh')
 12 tf. keras. layers. Dense(1, activation = 'sigmoid')
 13

Generating text with Neural Networks

⇒ Building the vocab

1 data = "Song lyrics separated by \n")

2 corpus = data. split(lower()). split('\n')

3 tokenizer = Tokenizer()

4 tokenizer. fit_on_texts(corpus)

5 total_words = len(tokenizer.word_index) + 1

⇒ Preprocessing the data.

6 input_sequences = []

for line in corpus:

7 token_list = tokenizer.fit_texts_to_sequences
 ([line])[0]

8 for i in range(1, len(token_list)):

n_gram_sequences = token_list[:i+1]

input_sequences.append(n_gram_sequences)

max_seq_length = max([len(x) for x in input_sequences])

NOTES

M	T	W	T	F	S	S
30	31	1	2	3	4	5
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

January 2025

M	T	W	T	F	S	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

Thursday

07

NOV 2024
J-2024-WK15

input-sequences = np.array (pad-sequences (input-
 sequences, maxlen = max-seq-length,
 padding = 'pre'))

xs, labels = input-sequences [:, :-1],
 input-sequences [:, -1]

ys = tf.keras.utils.to_categorical (labels,
 num_classes = total_words)

sentence = corpus [0].split()

token-list = []

for word in sentence:

token-list.append (toknizer.word_index [word])

elem-number = 6

print (len (token-list))

print (toknizer.sequences_to_texts ([xs [elem-number]]))

print (ys [elem-number])

print (np.argmax (ys [elem-number]))

⇒ Build sequential model with embeddings
 (totalwords, 64, input_length = max-seq-
 length - 1), LSTM (10) bidirectional
 and dense (total words, softmax)

⇒ compile with categorical_crossentropy,
 -adam & accuracy metrics.

	M	T	W	T	F	S	S	
October 2024	1	2	3	4	5	6	7	1 2 3
	7	8	9	10	11	12	13	8 9 10
	14	15	16	17	18	19	20	14 15 16 17
	21	22	23	24	25	26	27	21 22 23 24
	28	29	30	31				29 30

	M	T	W	T	F	S	S	
November 2024	4	5	6	7	8	9	10	1 2 3
	11	12	13	14	15	16	17	11 12 13 14 15 16 17
	18	19	20	21	22	23	24	18 19 20 21 22 23 24
	25	26	27	28	29	30		25 26 27 28 29 30

Generating the first Algorithm

- 9 → Feed the seed text to initiate the process
- 10 → Model predicts the index of the most probable next word.
- 11 → Look up the index in reverse word index dictionary
- 12 → Append the next word to the seed text.
- 13 → Feed the result to the model again.

seed_text = "Pratik went to Dublin"

next_word = 100

for - in range(next_words):

token_list = tokenizer.texts_to_sequences([seed_text])[0]

token_list = pad_sequences([token_list], maxlen=

new_seq_length - 1, padding: 'pre')

probabilities = model.predict(token_list)

{ predicted = np.argmax(probabilities, axis=-1)}

feed the model to get proper probabilities

for each token and get the index with

highest probability

if predicted != 0:

output_word = tokenizer.index_word[predicted]

seed_text += " " + output_word

NOTES

	M	T	W	T	F	S	S
December 2024	30	31		1			
	2	3	4	5	6	7	8
	9	10	11	12	13	14	15
	16	17	18	19	20	21	22
	23	24	25	26	27	28	29

	M	T	W	T	F	S	S
January 2025	6	7	8	9	10	11	12
	13	14	15	16	17	18	19
	20	21	22	23	24	25	26
	27	28	29	30	31		

Saturday

09

NOV 2024
314-052 • WK 45

Text Generation with RNN

- 9 → Import Tensorflow & other libraries
- 10 → Download Shakespeare dataset
- 11 → Read the data
- 12 → Process the text
 - Vectorise the text
 - Prediction task
 - Create training example & targets
 - 1 → Create training batches
- 13 → Build the model.
- 2 → Train the ~~model~~ model
 - Attach an optimizer, & loss func.
 - 3 → Config checkpoints
 - Execute training
- 4 → Generate text
- Export generator.
- ~~to~~ customize training.

Sunday 10

NOTES