

## Natural Language Processing with Tensorflow

### Introduction:

In this project I will use Tensorflow framework is to all

always has training the same input/output structure. This is extremely useful if you want to use the same model in production, even if you want to access the model using different programming languages.

## Goal:

The goal of this project is to build a NLP API within Tensorflow that can text a string of text and predict if the customer either liked, or disliked the product that they reviewed. For this project I will using product reviews pulled from Amazon.

## Step One: Import required packages

```
import tensorflow as tf
```

```
import
```

```
%load_ext tensorboard

import datetime

from tensorflow import keras

from tensorflow.keras import layers
```

```
from tensorflow.keras.layers import Dense
from tensorflow.keras.models import Sequential
```

```

from tensorflow.keras.preprocessing.text import Tokenizer

from sklearn.metrics import mean_squared_error

from sklearn.feature_extraction.text import CountVectorizer

from sklearn.model_selection import train_test_split

from keras.preprocessing.sequence import pad_sequences

from pprint import pprint

from collections import defaultdict

import pandas as pd

pd.set_option('display.max_rows', 1000)

import numpy as np

from nltk.tokenize import regexp_tokenize

import gensim

from gensim import corpora

from gensim.corpora.dictionary import Dictionary

import nltk

from nltk.corpus import stopwords

from sklearn.feature_extraction.text import TfidfVectorizer

import string

import re

```

## Step Two: Import and prepare the data (Foy, 2020)

```

data = pd.read_csv('amazon_cells_labelled.csv', sep=',')

data.rename(columns={'Liked/Disliked': 'score', 'Review': 'text'}, inplace=True)

ax = data.groupby('score').count().plot.bar(ylim=0)
ax.set_title("Number of reviews by score")
plt.xticks(rotation=0);

```

score	text
0	~500
1	~500

The data contains two columns, first a column of cell phone reviews and then a column of 1's and 0's indicating if the customer liked the phone or not. From the graph above it is clear that the number of positive and negative reviews is equal.

To prepare the data first all NA's will be removed, then all special characters will be removed, then all numbers will be removed, and finally all capital letters will be replaced with lowercase letters.

### Next remove NA's from the data

```
df = data.dropna(subset=['text'])
```

### Next remove any special characters from the text column (Petit 2020)

```

spec_chars = ["!", "@", "#", "$", "%", "&", "'", "(", ")", "*", "+", ",", "-", ".", ":", ";", "<\/pre>
</div>
</html>
```

```
df['text'] = df['text'].str.lower()
```

Next tokenize the text column, for this part I followed the following guide: (Jedamski, 2020)

To tokenize the words in the reviews I will define a function named tokenize, this function will split each sentence into it's individual words. Once we have the words separated, we can filter out stopwords and then apply a vectorization to words replacing them with numerical values. We do this because Tensorflow cannot handle lists of string characters.

```
def tokenize(text):
    tokens = re.split('W+', text)
    return tokens

df['text_tokenized'] = df['text'].apply(lambda x: tokenize(x))
df.head()
```

	text	score	text_tokenized
0	so there is no way for me to plug it in here i	0	[so, there, is, no, way, for, me, to, plug, it-,
1	good case excellent value	1	[good, case, excellent, value, ]

3	tied to charger for conversations lasting more...	0	[tied, to, charger, for, conversations
4	the mic is great	1	[the, mic, is,

## Load stopwords

```
stopwords = nltk.corpus.stopwords.words('english')
```

## Next remove the stopwords

```
def remove_stopwords(text_tokenized):  
    text = [word for word in text_tokenized if word not in stopwords]  
    return text  
  
df['text_nostop'] = df['text_tokenized'].apply(lambda x: remove_stopwords(x))  
df.head()
```

	text	score	text_tokenized
0	so there is no way for me to plug it in here i...	0	[so, there, is, no, way, for, me, to, plug, it, ...]

```
2          great for jawbone      1          [great, for, the, jawbone.]          [great, jawbone.]
3 tied to charger for conversations lasting more... 0 [tied, to, charger, for, conversations, lastin...   [tied, charger, conversations, lasting, minute...
4          the mic is great      1          [the, mic, is, great.]          [mic, great.]

line_num_words = [len(t_line) for t_line in df['text_nostop']]

plt.hist(line_num_words)
plt.xlabel('length of review')
plt.ylabel('# of characters')
plt.xlabel('Frequency');
```

Sentence	Length of longest word
1	340
2	320
3	310
4	300
5	290

From the above graph we can see that the longest review is 20 words in the final cleaned dataset

```
tfidf_vec = TfidfVectorizer()
X_tfidf = tfidf_vec.fit_transform(doc['text_final'])
print(X_tfidf.shape)

print(tfidf_vec.get_feature_names())
```

(1000, 1686)

'abbasy', 'ability', 'able', 'abound', 'absolutely', 'absolutely', 'ac', 'accept', 'acceptable', 'access',  
'accessible', 'accessing', 'accessory', 'accessories', 'accidentally', 'accompanied', 'according', 'activities',  
'activated', 'actively', 'actual', 'act', 'action', 'actions', 'add', 'addition', 'additional',  
dress', 'adhesive', 'adorable', 'advertised', 'advise', 'aggravating', 'ago', 'alarm', 'alert', 'allow',  
'allowing', 'almost', 'alone', 'along', 'also', 'also', 'although', 'aluminum', 'always', 'amazed',  
amazon', 'amp', 'among', 'another', 'angle', 'another', 'answer', 'antenna', 'anti', 'and',  
net', 'anything', 'anyway', 'apart', 'apartment', 'apparently', 'appealing', 'appearance', 'appear',  
'appeared', 'areas', 'arabic', 'armband', 'arrival', 'arrive', 'arriving', 'artist', 'artists', 'as'

[illegible]

t', 'defective', 'deffinitely', 'definitely',  
scription', 'design', 'designed', 'designs',  
d', 'development', 'device', 'devices', 'dial',  
'directions', 'directly', 'dirty', 'disapoinme  
ing', 'disappointment', 'discarded', 'discomf

[illegible][illegible][illegible][illegible][illegible][illegible][illegible]

```

The output, X_train, is in sparse matrix format; tensorflow cannot handle this object type.
X_train must be converted to a Pandas dataframe first.

X_train = pd.DataFrame(X_train.toarray())
X_train.head()

```

5 rows × 1686 columns

## Next export the cleaned dataset to csv

```
X_features.to_csv('data_clean.csv')
```

## Building the model in Tensorflow (Load a Pandas.DataFrame | TensorFlow Core, 2021)

```
target = df.pop('score')

X_train, X_test, y_train, y_test = train_test_split(X_features, target, test_size=0.20, random_state=100)
```

For this project the data is split 80/20, 80% of the words will be used for training the model and 20% will be used for testing.

```
print(df.shape) | print(X_train.shape) | print(X_test.shape)
```

```
(1000, 4)
(800, 1686)
(200, 1686)
```

## Next convert the DataFrame to a Tensor object (Ydobon, 2019)

By converting our vectorized list to a tensorflow data object we will not need to pad our sequences. The reason for this is that each word in our dataset is assigned a value by the vectorization function, then each of these values is placed in its own cell within the tensor. This means that when we convert the list to a tensorflow, and because we will not have to worry about the strings having uniform length, I will also not be imposing a stopping criteria on the Tensorflow model, instead I will vary the number the number of epochs and see what effect that on the model accuracy.

```
dataset = tf.data.Dataset.from_tensor_slices((X_train_features.values, target.values))
```

```
train_dataset = dataset.shuffle(len(df)).batch(1)

def get_compiled_model():
    model = tf.keras.Sequential([
        tf.keras.layers.Dense(10, activation='relu'),
        tf.keras.layers.Dense(10, activation='relu'),
        tf.keras.layers.Dense(1)
    ])

    model.compile(optimizer='adam',
                  loss=tf.keras.losses.BinaryCrossentropy(from_logits=True),
                  metrics=['accuracy'])

    return model
```

## Model Features:

This is a sequential model, featuring three layers. The first two layers have ten nodes each, while the final layer features one node. The model metric is accuracy, and the loss function is BinaryCrossentropy. The nodes where chosen by taking the square root of input features, which is  $\sqrt{16384}$  which equals 41. That is then rounded down to the nearest power or ten, which in this case is ten.

```
log_dir = "logs/fit/" + datetime.datetime.now().strftime("%Y%m%d-%H%M%S")
tensorboard_callback = tf.keras.callbacks.TensorBoard(log_dir=log_dir, histogram_freq=1)
```

The model will be fit to the training data, set in the next step.

```
model = get_compiled_model()
history = model.fit(train_dataset, epochs=15, callbacks=[tensorboard_callback])
model.summary()
```

Epoch 1/15  
WARNING:tensorflow:Layer dense\_6 is casting an input tensor from dtype float64 to the layer's dtype of float32 which may cause precision loss. This has been detected when the layer's dtype is different from its inputs' dtype (float64 vs tf.float32). The layer has dtype float32 because its dtype defaults to float32.

If you intended to run this layer in float32, you can safely ignore this warning. If in doubt, this warning is likely only an issue if you are porting a TensorFlow 1.x model to TensorFlow 2.x.

To change all layers to have dtype float64 by default, call 'tf.keras.backend.set\_floatx('float64')'. To change just this layer, call 'dtype='float64' to the layer constructor. If you are the author of this layer, you can disable autocasting by passing autocast=False to the base Layer constructor.

```

27/1000 [.....] 100% ETA: 41m 41s 0.6917 accuracy: 0.5000 Attention:tensorflow
1 epoch ended on train batch end 100% ETA: 41m 41s 0.6915 loss on train
batch_end_time: 0.5638s. Check your callbacks.
Epoch 1/10: 100% 1 ms/step - loss: 6.5755 - accuracy: 0.5620
Repoch 2/15
1000/1000 [.....] 100%
Epoch 3/15: 100% 1 m16us/step - loss: 0.3446 - accuracy: 0.8990
1000/1000 [.....] 100%
Epoch 4/15: 100% 1 m59us/step - loss: 0.1292 - accuracy: 0.9570
1000/1000 [.....] 100%
Epoch 5/15: 100% 1 m12us/step - loss: 0.0565 - accuracy: 0.9840
1000/1000 [.....] 100%
Epoch 6/15: 100% 1 m19us/step - loss: 0.3032 - accuracy: 0.9910
1000/1000 [.....] 100%
Epoch 7/15: 100% 1 m61us/step - loss: 0.0184 - accuracy: 0.9950
1000/1000 [.....] 100%
Epoch 8/15: 100% 1 m60us/step - loss: 0.0130 - accuracy: 0.9980
1000/1000 [.....] 100%
Epoch 9/15: 100% 1 m58us/step - loss: 0.0099 - accuracy: 0.9990

```

```
Epoch 9/15 [=====] - lrs: 6170s/step - loss: 0.0077 - accuracy: 0.9970
1000/1000 [=====] - lrs: 6170s/step - loss: 0.0096 - accuracy: 0.9960
Epoch 10/15 [=====] - lrs: 6170s/step - loss: 0.0096 - accuracy: 0.9960
1000/1000 [=====] - lrs: 6170s/step - loss: 0.0096 - accuracy: 0.9960
Epoch 11/15 [=====] - lrs: 624us/step - loss: 0.0068 - accuracy: 0.9960
1000/1000 [=====] - lrs: 624us/step - loss: 0.0068 - accuracy: 0.9960
Epoch 12/15 [=====] - lrs: 654us/step - loss: 0.0073 - accuracy: 0.9980
1000/1000 [=====] - lrs: 654us/step - loss: 0.0067 - accuracy: 0.9970
Epoch 13/15 [=====] - lrs: 651us/step - loss: 0.0067 - accuracy: 0.9970
1000/1000 [=====] - lrs: 651us/step - loss: 0.0067 - accuracy: 0.9970
Epoch 14/15 [=====] - lrs: 704us/step - loss: 0.0057 - accuracy: 0.9970
1000/1000 [=====] - lrs: 704us/step - loss: 0.0057 - accuracy: 0.9970
Epoch 15/15 [=====] - lrs: 599us/step - loss: 0.0064 - accuracy: 0.9960
1000/1000 [=====] - lrs: 599us/step - loss: 0.0064 - accuracy: 0.9960
Model: "sequential_2"
```

Layer (type)	Output Shape	Param #
--------------	--------------	---------

```

dense_6 (Dense)                (None, 10)                10070
dense_7 (Dense)                (None, 10)                110
dense_8 (Dense)                (None, 1)                 11
-----
Total params: 16,991
Trainable params: 16,991
Non-trainable params: 0

```

---

Next we need to make sure that the model is not overfitting the data. According to (Tensorflow.org, 2020) "If the validation metric is going in the wrong direction, the model is clearly

accuracy drops going from epoch 13 to epoch 14. So yes there is some overfitting going on. This is another reason to use a smaller number of epochs.


Next we will visualize the training for different epochs versus the Mean Absolute Error (Chris & Sara, 2021)

```
plt.plot(history['loss'], label='MSE (training data)')
plt.plot(history['val_loss'], label='MSE (validation data)')
plt.xlabel('Epochs')
```

```
plt.legend(loc="upper left")
plt.show()
```

The plot shows the Mean Absolute Error (MAE) for training data. The y-axis is labeled 'MAE value' and ranges from 0.2 to 0.6. The x-axis is unlabeled but represents an index or iteration. A single blue line, identified in the legend as 'MAE (training data)', starts at a value of approximately 0.6 at the first point and decreases rapidly to about 0.25 at the second point.

Index	MAE (training data)
0	~0.60
1	~0.25



From the above graph the MAE reaches a minimum at around 8 epochs

Next create a tensorboard object analyze the model (Get Started with TensorBoard |, 2020)

```
%tensorboard --logdir logs/fit;
```

Jupyter notebook gives an error when opening the TensorBoard object but it opens successfully on <http://localhost:6006/>

From looking at the graph of Epoch Accuracy it is clear that 8 epochs is the value that maxes out the accuracy

Next create a new model with 8 epochs

```
model = get_compiled_model()
model.fit(train_dataset, epochs=8, callbacks=[tensorboard_callback])
model.summary()
```

Epoch 1/8  
WARNING:tensorflow:Layer dense\_3 is casting an input tensor from dtype float64 to the layer's dtype of float32

```

X
If you intended to run this layer in float32, you can safely ignore this warning. If in doubt, this warning
is likely only an issue if you are porting a TensorFlow 1.x model to TensorFlow 2.

To change all layers to have dtype float64 by default, call tf.keras.backend.set_floatx('float64'). To
change a single layer, pass dtype='float64' to the layer constructor. If you are the author of this layer,
it can disable autocasting by passing autocast=False to the base Layer constructor.

2/10000 [.....] ETA: 516s - loss: 0.6876 - accuracy: 0.5000 WARNING:tensorflow:
Callables method 'on_train_batch_end' is slow compared to the batch time (batch time: 0.0010s vs 'on_train_
batch_end' time: 0.6324s). Check your callbacks.

10000/10000 [.....] - loss: 0.6531 - accuracy: 0.5380
Epoch 2/8
10000/10000 [.....] - loss: 0.6705/step - loss: 0.3430 - accuracy: 0.8710
Epoch 3/8
10000/10000 [.....] - loss: 0.6398/step - loss: 0.1295 - accuracy: 0.9560

```

```

1000/1000 (=====) - 1s 606us/step - loss: 0.0619 - accuracy: 0.9830
Epoch 5/8
1000/1000 (=====) - 1s 618us/step - loss: 0.0339 - accuracy: 0.9920
Epoch 6/8
1000/1000 (=====) - 1s 712us/step - loss: 0.0222 - accuracy: 0.9950
Epoch 7/8
1000/1000 (=====) - 1s 672us/step - loss: 0.0151 - accuracy: 0.9970
Epoch 8/8
1000/1000 (=====) - 1s 620us/step - loss: 0.0122 - accuracy: 0.9950
Model: "sequential_3"

Layer (type)                Output Shape                Param #
-----
dense_9 (Dense)              (None, 10)                  16870
dense_10 (Dense)              (None, 10)                  110

```

```
=====
Total params: 16,991
Trainable params: 16,991
Non-trainable params: 0
=====
```

From the above model 8 epochs are enough to achieve an accuracy of 99.7%

## Network Architecture

The model uses the following architecture:

Activation: Relu, this is a rectified linear activation function. It outputs the input value for

number of nodes per layer: For this take the square root of the number of features and the try values near to this and see how it effects the accuracy. Ten nodes for each of the first two layers, and one node in the final layer produce an accuracy of over 99%

Loss function = BinaryCrossentropy. This problem is a binary classification problem, and so BinaryCrossentropy is perfect for this type of problem.

Optimizer = Adam. "computationally efficient, has little memory requirement, invariant to diagonal rescaling of gradients, and is well suited for problems that are large in terms of data/parameters". (Kingma & Ba, 2017)

## Network Functionality:

This is a very functional neural network. It produces very high levels of accuracy over a large range of words. Taking this level of accuracy to a different set of data would be very easy.

## Course of action:

The goal of this project was to show that it is possible to predict if a review would be positive or negative based on analyzing key words. The results of the model applied to the

action would be for a company to use this information to train bots to create fake positive reviews. This would allow the company to increase profits by having people that check reviews online before making purchases see these created reviews and trust that they are real reviews.

## References

N. (2020, June 2). 4 Ways to Use Pandas to Select Columns in a Dataframe. Datagy. <https://datagy.io/pandas-select-columns/>

Removing URL from a column in Pandas Dataframe. (2018, August 23). Stack Overflow. <https://stackoverflow.com/questions/51994254/removing-url-from-a-column-in-pandas->

Category: nltk. (2020, June 26). Pythonspot. <https://pythonspot.com/category/nltk/>

Team, K. (2020). Keras documentation: The Sequential model. K. Team. [https://keras.io/guides/sequential\\_model/](https://keras.io/guides/sequential_model/)

Petit, U. (2020, November 13). Simplify your Dataset Cleaning with Pandas - Towards Data Science. Medium. <https://towardsdatascience.com/simplify-your-dataset-cleaning-with-pandas-75951b23568e>

Load a pandas.DataFrame | TensorFlow Core. (2021). TensorFlow. [https://www.tensorflow.org/tutorials/load/data/pandas\\_dataframe](https://www.tensorflow.org/tutorials/load/data/pandas_dataframe)

Yobdon, A. (2019, September 26). [Tensorflow 2.0] Load Pandas dataframe to Tensorflow. Medium, <https://financial-engineering.medium.com/load-pandas-2-0-load-pandas-dataframe-to-tensorflow-202c5d48966b>

Jedamski, D. (2020, August 10). Compare all methods using key performance metrics - Advanced NLP with Python for Machine Learning. LinkedIn, <https://www.linkedin.com/learning/advanced-nlp-with-python-for-machine-learning/compare-all-methods-using-key-performance-metrics?u=2045532>

Get started with TensorBoard |. (2020). TensorFlow. [https://www.tensorflow.org/tensorboard/get\\_started](https://www.tensorflow.org/tensorboard/get_started)

Brownlee, J. (2020, August 20). A Gentle Introduction to the Rectified Linear Unit (ReLU).  
 Machine Learning Mastery. <https://machinelearningmastery.com/rectified-linear-activation-function-for-deep-learning-neural-networks/>

[j/jax/output/CommonHTML/fonts/TeX/fontdata.js](#)



Kingma, D., & Ba, J. (2017, January 30). Adam: A method for stochastic optimization. Retrieved March 11, 2021, from <https://arxiv.org/abs/1412.6980>

Chris, & Sara. (2021, January 26). Visualizing training performance with TensorFlow 2 and Keras. Retrieved March 11, 2021, from <https://www.machinecurve.com/index.php/2019/10/08/how-to-visualize-the-training-process-in-keras/>