

## ▼ Supervised/unsupervised Sentiment and Topic analysis.

Final project assignment.

Introduction to Natural Language Processing (NLP) (DSCI-D590-31731)

Anitha Ganapathy | [aganapa@iu.edu](mailto:aganapa@iu.edu)

[12/01/2020](https://www.iu.edu/~aganapa/)

```
1 !nvidia-smi
```

```
Sat Dec 12 05:54:33 2020
```

NVIDIA-SMI 455.45.01			Driver Version: 418.67		CUDA Version: 10.1		
GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr.	ECC
Fan	Temp	Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute M.	MIG M.
0	Tesla T4	Off	00000000:00:04.0	Off	0%	Default	0
N/A	55C	P8	10W / 70W	0MiB / 15079MiB			ERR!

Processes:							
GPU	GI	CI	PID	Type	Process name	GPU	Memory
ID	ID					Usage	
No running processes found							

```
1 # !nvcc --version
```

```
1 # !pip install docx2txt
2 # !pip install tensorflow_text
3
```

Organizing imports.

```
1 import nltk
2 import os
3 import sys
4 import pandas as pd
5 import numpy as np
6 import re
7 import nltk
```

```
8 import matplotlib.pyplot as plt
9 import matplotlib.colors as mcolors
10 import docx2txt
11 import tensorflow as tf
12 import tensorflow_hub as hub
13 import tensorflow_text
14 import seaborn as sns
15 import gensim
16 import json
17 import datetime
18
19 from tqdm import tqdm
20 from nltk import word_tokenize
21 from collections import Counter
22 from nltk.corpus import stopwords
23 from nltk import WordPunctTokenizer
24 from nltk.corpus import stopwords
25 from nltk.stem import PorterStemmer, WordNetLemmatizer
26 from nltk import wordnet
27 from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator
28 from tensorflow import keras
29 from keras.preprocessing.text import Tokenizer
30 from keras.preprocessing.sequence import pad_sequences
31 from keras.models import Model, Sequential
32 from keras.layers import Dense, Embedding, LSTM, Dropout
33 from keras.layers import Input, Dense, Embedding, SpatialDropout1D, add, concatenate
34 from keras.layers import Bidirectional, GlobalMaxPooling1D, GlobalAveragePooling1D
35 from keras.callbacks import Callback, ModelCheckpoint
36 from keras.preprocessing import text, sequence
37 from sklearn.preprocessing import OneHotEncoder
38 from keras.utils import to_categorical
39 from sklearn.model_selection import train_test_split
40 from sklearn.metrics import classification_report
41 from sklearn.manifold import TSNE
42 from gensim.models import word2vec
43 from pylab import rcParams
44 from pandas.plotting import register_matplotlib_converters
45
46 # plt.style.use('ggplot')
47 %config InlineBackend.figure_format='retina'
48 %matplotlib inline
49
50 import warnings
51 warnings.filterwarnings('ignore')
52
```

```
1 nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
True
```

```

1 rcParams['figure.figsize'] = 12, 8
2 RANDOM_SEED = 2000737430
3 np.random.seed(RANDOM_SEED)
4 tf.random.set_seed(RANDOM_SEED)
5
6 register_matplotlib_converters()
7 sns.set(style='whitegrid', palette='muted', font_scale=1.2)
8 HAPPY_COLORS_PALETTE = ["#01BEFE", "#FFDD00", "#FF7D00", "#FF006D", "#ADFF02", "#8F00FF"]
9 sns.set_palette(sns.color_palette(HAPPY_COLORS_PALETTE))

1 print(os.listdir("../content"))
2 print(f'Python version      : {sys.version}')
3 print(f'Pandas version     : {pd.__version__}')
4 print(f'Numpy version      : {np.__version__}')
5 print(f'Tensorflow version : {tf.__version__}')
6 print(f'NLTK version       : {nltk.__version__}')
7 print(f'Regex version      : {re.__version__}')
8 print("GPU Device name: ",tf.config.list_physical_devices('GPU'))
9 # print("\nCheck if GPU is available: ", tf.test.is_gpu_available())

['.config', 'Amazon Fine Food Reviews.csv', 'sample_data']
Python version      : 3.6.9 (default, Oct  8 2020, 12:12:24)
[GCC 8.4.0]
Pandas version     : 1.1.5
Numpy version      : 1.18.5
Tensorflow version : 2.3.0
NLTK version       : 3.2.5
Regex version      : 2.2.1
GPU Device name:  [PhysicalDevice(name='/physical_device:GPU:0', device_type='GPU')]

```

**Dataset** The data I am using for the assignment is the Amazon Fine Food Reviews.

<https://www.kaggle.com/snap/amazon-fine-food-reviews>

### Data Reference:

J. McAuley and J. Leskovec. From amateurs to connoisseurs: modeling the evolution of user expertise through online reviews. WWW, 2013.

```
1 df = pd.read_csv("Amazon Fine Food Reviews.csv", parse_dates= True)
```

## ▼ Step 1 : Describe data.

```

1 print("The Amazon Fine Food Reviews dataset brief:")
2 print("\nNumber of reviews:", df.shape[0])
3 print("Number of users: ", len(df.UserId.unique()))
4 print("Number of products: ", len(df.ProductId.unique()))
```

```
5 print("Number of Attributes/Columns in data: ", len(df.columns))
```

The Amazon Fine Food Reviews dataset brief:

```
Number of reviews: 568454  
Number of users: 256059  
Number of products: 74258  
Number of Attributes/Columns in data: 10
```

```
1 # df.head()
```

The column or features in the dataset:

- Id
- ProductId – unique identifier for the product
- UserId – unqiue identifier for the user
- ProfileName
- HelpfulnessNumerator – number of users who found the review helpful
- HelpfulnessDenominator – number of users who indicated whether they found the review helpful or not
- Score – rating between 1 and 5
- Time – timestamp for the review
- Summary – brief summary of the review
- Text – text of the review

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 568454 entries, 0 to 568453  
Data columns (total 10 columns):  
 #   Column           Non-Null Count  Dtype     
---  --  
 0   Id               568454 non-null  int64    
 1   ProductId        568454 non-null  object    
 2   UserId            568454 non-null  object    
 3   ProfileName       568438 non-null  object    
 4   HelpfulnessNumerator  568454 non-null  int64    
 5   HelpfulnessDenominator 568454 non-null  int64    
 6   Score             568454 non-null  int64    
 7   Time              568454 non-null  int64    
 8   Summary            568427 non-null  object    
 9   Text               568454 non-null  object    
dtypes: int64(5), object(5)  
memory usage: 43.4+ MB
```

```
1 df.describe()
```

	<b>Id</b>	<b>HelpfulnessNumerator</b>	<b>HelpfulnessDenominator</b>	<b>Score</b>		
<b>count</b>	568454.000000	568454.000000	568454.000000	568454.000000	568454.000000	5.684
<b>mean</b>	284227.500000	1.743817	2.22881	4.183199	1.296	
<b>std</b>	164098.679298	7.636513	8.28974	1.310436	4.804	
<b>min</b>	1.000000	0.000000	0.000000	1.000000	9.393	
<b>25%</b>	142114.250000	0.000000	0.000000	4.000000	1.271	
<b>50%</b>	284227.500000	0.000000	1.000000	5.000000	1.311	
<b>75%</b>	426340.750000	2.000000	2.000000	5.000000	1.332	
-----	568454.000000	0.000000	0.000000	5.000000	1.251	

Checking for NAN values in the dataset. Nan values in the text column provided difficulty in plotting the wordcloud. So replaced the nan with empty string.

```
1 df.isna().sum()
```

<b>Id</b>	0
<b>ProductId</b>	0
<b>UserId</b>	0
<b>ProfileName</b>	16
<b>HelpfulnessNumerator</b>	0
<b>HelpfulnessDenominator</b>	0
<b>Score</b>	0
<b>Time</b>	0
<b>Summary</b>	27
<b>Text</b>	0
<b>dtype: int64</b>	

## ▼ Step 2: Perform data preprocessing.

```
1 #df = pd.read_csv("Amazon Fine Food Reviews.csv", parse_dates= True)
2
3 df_old = df
4 df.Summary = df.Summary.fillna('')      # replacing the NAN with ''
5 df['review'] = df['Summary'] + df['Text']
6 df.isna().sum()
```

<b>Id</b>	0
<b>ProductId</b>	0
<b>UserId</b>	0
<b>ProfileName</b>	16
<b>HelpfulnessNumerator</b>	0
<b>HelpfulnessDenominator</b>	0
<b>Score</b>	0
<b>Time</b>	0
<b>Summary</b>	0

```
Text          0  
review        0  
dtype: int64
```

```
1 df = df[['Score', 'review']]  
2 df.head()
```

	Score	review
0	5	Good Quality Dog FoodI have bought several of ...
1	1	Not as AdvertisedProduct arrived labeled as Ju...
2	4	"Delight" says it allThis is a confection that...
3	2	Cough Medicinelf you are looking for the secre...
4	5	Great taffyGreat taffy at a great price. Ther...

```
1 df.Score.value_counts()
```

```
5    363122  
4    80655  
1    52268  
3    42640  
2    29769  
Name: Score, dtype: int64
```

```
1 df.isna().sum()
```

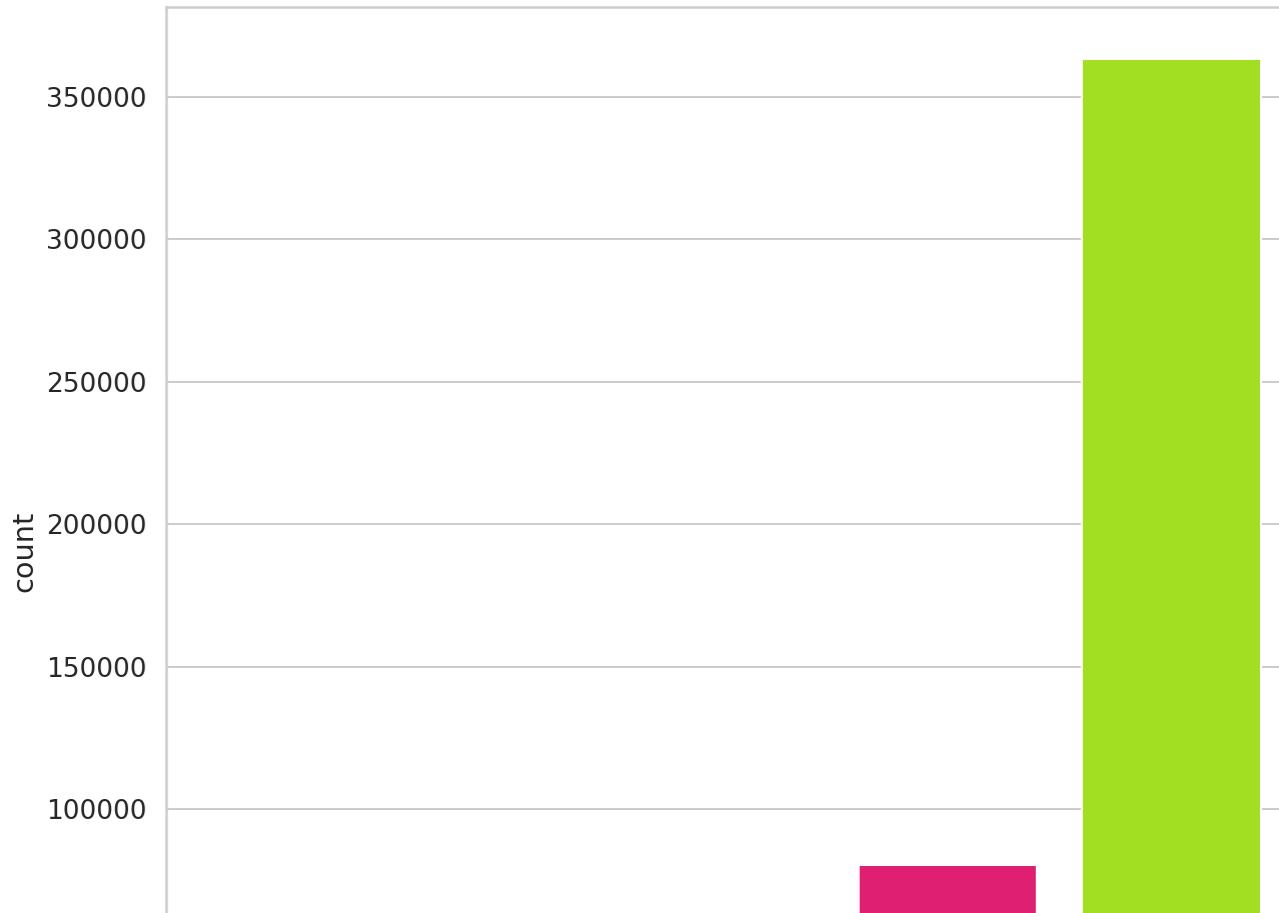
```
Score      0  
review     0  
dtype: int64
```

## ▼ Analysing the data.

```
1 # sns.countplot(df['Score'], order=df.Score.value_counts().index)  
2 sns.countplot(df['Score'])  
3 fig = plt.gcf()  
4 fig.set_size_inches(10,10)  
5 plt.title('Amazon Food and Wine Reviews.')
```

```
Text(0.5, 1.0, 'Amazon Food and Wine Reviews.')
```

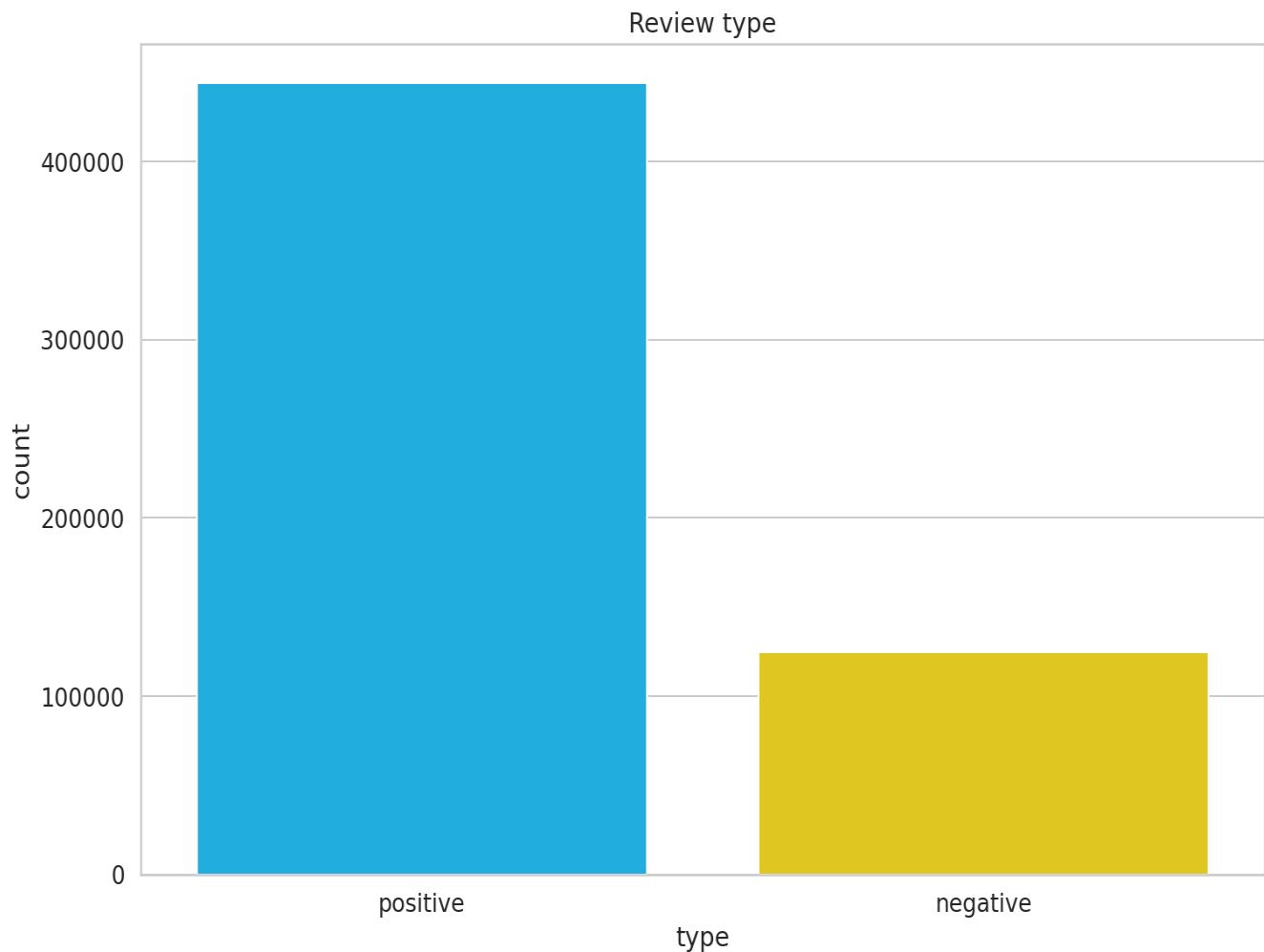
Amazon Food and Wine Reviews.



Just for simplicity and because we need more data for the binary classification of Good or Bad reviews, we shall consider a rating of 3 or less as bad reviews and a rating of 4 and more as good reviews

```
1 # Give reviews with Score > 3 a positive rating('good'),  
2 # and reviews with a score<=3 a negative review_type('bad').  
3  
4 df["review_type"] = df["Score"].apply(lambda x: "negative" if x <= 3 else "positive")  
  
1 df.review_type.value_counts()  
  
positive    443777  
negative    124677  
Name: review_type, dtype: int64  
  
1 sns.countplot(  
2     x='review_type',  
3     data=df,  
4     order=df.review_type.value_counts().index  
5 )  
6 plt.xlabel("type")  
7 plt.title("Review type")
```

```
Text(0.5, 1.0, 'Review type')
```



```
1 positive_reviews = df[df.review_type == "positive"]
2 negative_reviews = df[df.review_type == "negative"]
3
4 print("Good reviews shape: ", positive_reviews.shape)
5 print("Bad reviews shape: ", negative_reviews.shape)
```

```
Good reviews shape: (443777, 3)
Bad reviews shape: (124677, 3)
```

We need the good and bad reviews to have same review counts, so we will take equal amount of good reviews as bad.

```
1 positive_reviews_text = " ".join(positive_reviews.review.to_numpy().tolist())
2 negative_reviews_text = " ".join(negative_reviews.review.to_numpy().tolist())
3
4 from nltk.corpus import stopwords
5 wpt = nltk.WordPunctTokenizer()
6 stop_words = nltk.corpus.stopwords.words('english')
```

```
4 def normalize_document(doc):
5     # lowercase and remove special characters\whitespace
6     doc = re.sub(r'[^a-zA-Z\s]', '', doc, flags =re.I)
7     doc = doc.lower()
8     doc = doc.strip()
9     # tokenize document
10    tokens = wpt.tokenize(doc)
11    # filter stopwords out of document
12    filtered_tokens = [token for token in tokens if token not in stop_words]
13    # re-create document from filtered tokens
14    doc = ' '.join(filtered_tokens)
15    return doc
16
17
```

```
1 positive_normalize_text = normalize_document(positive_reviews_text)
2 negative_normalize_text = normalize_document(negative_reviews_text)
3
```

```
1 datetime.datetime.now()
2 positive_reviews_cloud = WordCloud(stopwords=STOPWORDS,background_color="black").generate(po
3 negative_reviews_cloud = WordCloud(stopwords=STOPWORDS, background_color="black").generate(n
4 datetime.datetime.now()
```

```
1 print(datetime.datetime.now())
2 def show_word_cloud(cloud, title):
3     plt.figure(figsize = (20, 20))
4     plt.imshow(cloud, interpolation='bilinear')
5     plt.title(title)
6     plt.axis("off")
7     plt.show()
8 print(datetime.datetime.now())
```

```
2020-12-12 06:02:10.231269
```

```
2020-12-12 06:02:10.231701
```

## ▼ Step 3: Topic Analysis

```
1 show_word_cloud(positive_reviews_cloud, "Positive Reviews cloud display.")
```

Positive Reviews cloud display.



```
1 show_word_cloud(negative_reviews.cloud, "Negative Reviews cloud display.")
```

Negative Reviews cloud display.



We'll deal with the review type count imbalance by sampling the number of good ones to that of the bad ones. We need same amount of positive and negative reviews.

```
1 positive_df = positive_reviews.sample(n=len(negative_reviews), random_state=RANDOM_SEED)
2 negative_df = negative_reviews

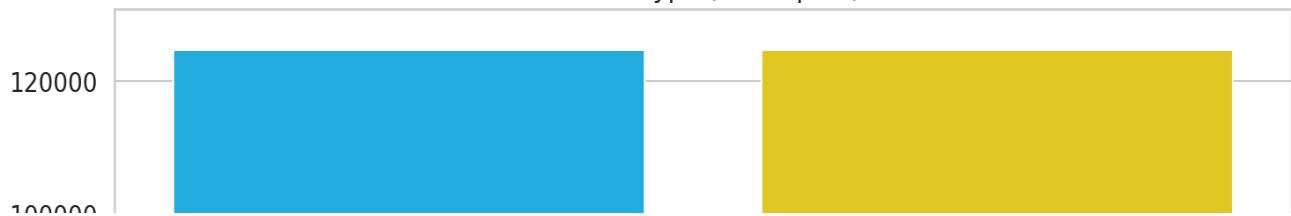
1 review_df = positive_df.append(negative_df).reset_index(drop=True)
2 print("Final review data df shape: ", review_df.shape)

Final review data df shape: (249354, 3)

1 sns.countplot(
2     x='review_type',
3     data=review_df,
4     order=review_df.review_type.value_counts().index
5 )
6
7 plt.xlabel("type")
8 plt.title("Review type (resampled)")
```

```
Text(0.5, 1.0, 'Review type (resampled)')
```

Review type (resampled)



Double-click (or enter) to edit



google/universal-sentence-encoder/4 - a much larger model yielding 512 dimensional embeddings trained with a deep averaging network (DAN) encoder.

The Universal Sentence Encoder encodes text into high dimensional vectors that can be used for text classification, semantic similarity, clustering, and other natural language tasks.



```
1 import tensorflow_hub as hub
2
3 # it took 14 mins to load this module
4
5 print(datetime.datetime.now())
6 use = hub.load("https://tfhub.dev/google/universal-sentence-encoder-multilingual-large/3")
7 print(datetime.datetime.now())
```

2020-12-12 06:02:59.901877

2020-12-12 06:16:37.210808

```
1 sent_1 = ["the location is great"]
2 sent_2 = ["amazing location"]
3
4 emb_1 = use(sent_1)
5 emb_2 = use(sent_2)
6 emb_1.shape
```

TensorShape([1, 512])

```
1 np.inner(emb_1, emb_2).flatten()[0]
```

0.79254675

## ▼ Step 4. Splitting data into training and testing.

Split in 80-20 training or UseCross-validation

```
1 type_one_hot = OneHotEncoder(sparse=False).fit_transform(
2     review_df.review_type.to_numpy().reshape(-1, 1)
```

3 )

```
1 train_reviews, test_reviews, y_train, y_test =\  
2     train_test_split(  
3         review_df.review,  
4         type_one_hot,  
5         test_size=.2,  
6         random_state=RANDOM_SEED  
7     )
```

## ▼ Step 5: Perform data vectorization.

```
1 print(datetime.datetime.now())  
2 X_train = []  
3 for r in tqdm(train_reviews):  
4     embedded = use(r)  
5     review_emb = tf.reshape(embedded, [-1]).numpy()  
6     X_train.append(review_emb)  
7  
8 X_train = np.array(X_train)  
9 print(datetime.datetime.now())
```

```
0%|          | 5/199483 [00:00<1:14:57, 44.35it/s]2020-12-12 06:20:09.520324  
100%|██████████| 199483/199483 [51:58<00:00, 63.97it/s]  
2020-12-12 07:12:08.403913
```

```
1 print(datetime.datetime.now())  
2 X_test = []  
3 for r in tqdm(test_reviews):  
4     emb = use(r)  
5     review_emb = tf.reshape(emb, [-1]).numpy()  
6     X_test.append(review_emb)  
7  
8 X_test = np.array(X_test)  
9 print(datetime.datetime.now())
```

```
0%|          | 5/49871 [00:00<17:36, 47.19it/s]2020-12-12 07:12:38.939502  
100%|██████████| 49871/49871 [13:01<00:00, 63.82it/s]2020-12-12 07:25:40.443657
```

```
1 print(X_train.shape, X_test.shape)  
  
(199483, 512) (49871, 512)
```

```
1 print(y_train.shape, y_test.shape)  
  
(199483, 2) (49871, 2)
```

## ▼ Step 6: Sentiment analysis

Sentiment Analysis is a binary classification problem. Let's use Keras to build a model.

```
1 model = keras.Sequential()
2
3 model.add(keras.layers.Dense(units=256, input_shape=(X_train.shape[1], ), activation='relu'))
4 model.add(keras.layers.Dropout(rate=0.2))
5 model.add(keras.layers.Dense(units=128, activation='relu'))
6 model.add(keras.layers.Dropout(rate=0.2))
7 model.add(keras.layers.Dense(2, activation='softmax'))
8 model.compile(loss='categorical_crossentropy', optimizer=keras.optimizers.Adam(0.001), metrics=[f1])
9 model.summary()
```

Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
dense (Dense)	(None, 256)	131328
dropout (Dropout)	(None, 256)	0
dense_1 (Dense)	(None, 128)	32896
dropout_1 (Dropout)	(None, 128)	0
dense_2 (Dense)	(None, 2)	258
=====		
Total params: 164,482		
Trainable params: 164,482		
Non-trainable params: 0		

The model is composed of 2 fully-connected hidden layers. Dropout is used for regularization.

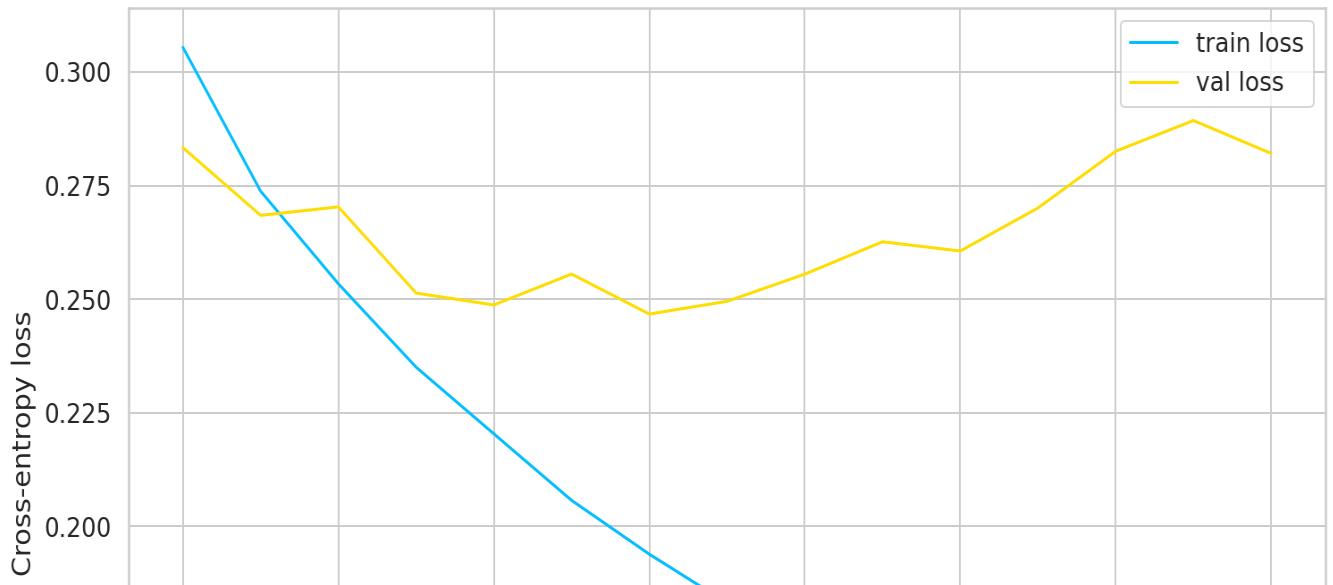
We'll train for 15 epochs and use 10% of the data for validation:

```
1 print(datetime.datetime.now())
2
3 history = model.fit(
4     X_train, y_train,
5     epochs=15,
6     batch_size=16,
7     validation_split=0.1,
8     verbose=1,
9     shuffle=True
10 )
11
```

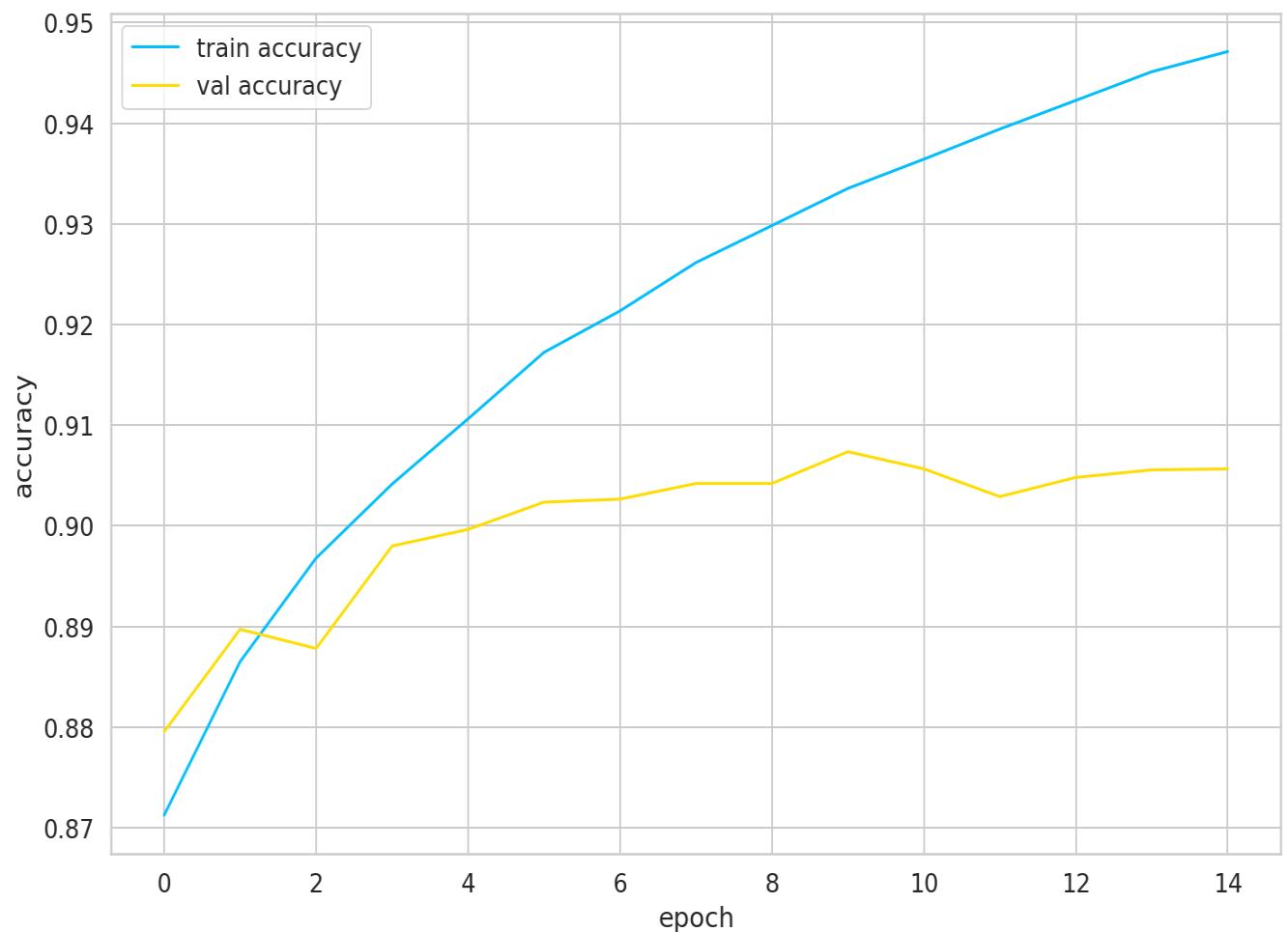
```
12 print(datetime.datetime.now())

2020-12-12 07:30:47.334659
Epoch 1/15
11221/11221 [=====] - 29s 3ms/step - loss: 0.3054 - accuracy: 0.1939
Epoch 2/15
11221/11221 [=====] - 31s 3ms/step - loss: 0.2737 - accuracy: 0.2057
Epoch 3/15
11221/11221 [=====] - 31s 3ms/step - loss: 0.2534 - accuracy: 0.2057
Epoch 4/15
11221/11221 [=====] - 31s 3ms/step - loss: 0.2350 - accuracy: 0.1939
Epoch 5/15
11221/11221 [=====] - 31s 3ms/step - loss: 0.2204 - accuracy: 0.1939
Epoch 6/15
11221/11221 [=====] - 29s 3ms/step - loss: 0.2057 - accuracy: 0.1831
Epoch 7/15
11221/11221 [=====] - 30s 3ms/step - loss: 0.1939 - accuracy: 0.1831
Epoch 8/15
11221/11221 [=====] - 30s 3ms/step - loss: 0.1831 - accuracy: 0.1723
Epoch 9/15
11221/11221 [=====] - 30s 3ms/step - loss: 0.1723 - accuracy: 0.1649
Epoch 10/15
11221/11221 [=====] - 30s 3ms/step - loss: 0.1649 - accuracy: 0.1565
Epoch 11/15
11221/11221 [=====] - 30s 3ms/step - loss: 0.1565 - accuracy: 0.1493
Epoch 12/15
11221/11221 [=====] - 30s 3ms/step - loss: 0.1493 - accuracy: 0.1426
Epoch 13/15
11221/11221 [=====] - 30s 3ms/step - loss: 0.1426 - accuracy: 0.1368
Epoch 14/15
11221/11221 [=====] - 30s 3ms/step - loss: 0.1368 - accuracy: 0.1314
Epoch 15/15
11221/11221 [=====] - 30s 3ms/step - loss: 0.1314 - accuracy: 0.1314
2020-12-12 07:38:21.637799
```

```
1 rcParams['figure.figsize'] = 12, 8
2 plt.plot(history.history['loss'], label='train loss')
3 plt.plot(history.history['val_loss'], label='val loss')
4 plt.xlabel("epoch")
5 plt.ylabel("Cross-entropy loss")
6 plt.legend();
```



```
1 rcParams['figure.figsize'] = 12, 8
2 plt.plot(history.history['accuracy'], label='train accuracy')
3 plt.plot(history.history['val_accuracy'], label='val accuracy')
4 plt.xlabel("epoch")
5 plt.ylabel("accuracy")
6 plt.legend();
```



```
1 loss, accuracy = model.evaluate(X_train, y_train, verbose=False)
```

```
2 loss, accuracy = model.evaluate(X_test, y_test, verbose=False)
3 print("Training Accuracy: {:.4f}".format(accuracy))
4 print("Testing Accuracy: {:.4f}".format(accuracy))

Training Accuracy: 0.9047
Testing Accuracy: 0.9047
```

1

```
Testing Accuracy: 0.9047
```

## ▼ Step 7: Predicting Sentiment

Example : 1

```
1 print(test_reviews.iloc[0])
2 print("Sentiment Analysis is :")
3 print("Bad" if y_test[0][0] == 1 else "Good")

ODNESSWhere has this butter been all my life!? I LOVE IT. It's like coconut flavored wh
ent Analysis is :
```

```
1 y_pred = model.predict(X_test[:1])
2 print(y_pred)
3 "Bad" if np.argmax(y_pred) == 0 else "Good"

[[1.3080695e-05 9.9998689e-01]]
'Good'
```

Example 2:

```
1 print(test_reviews.iloc[1])
2 print("Bad" if y_test[1][0] == 1 else "Good")
3

these are actually awfuli gave it 2 stars because i was able to swallow 2 spoonfuls. so
Bad
```

```
1 y_pred = model.predict(X_test[1:2])
2 print(y_pred)
3 "Bad" if np.argmax(y_pred) == 0 else "Good"
```

```
[[9.999907e-01 9.355102e-06]]  
'Bad'
```

Example 3:

```
1 print(test_reviews.iloc[5])  
2 print("Bad" if y_test[5][0] == 1 else "Good")
```

```
Change in blend?I am a long time lover of Caffe Verona whole beans. However, I bought  
Bad
```

```
1 y_pred = model.predict(X_test[5:6])  
2 print(y_pred)  
3 "Bad" if np.argmax(y_pred) == 0 else "Good"
```

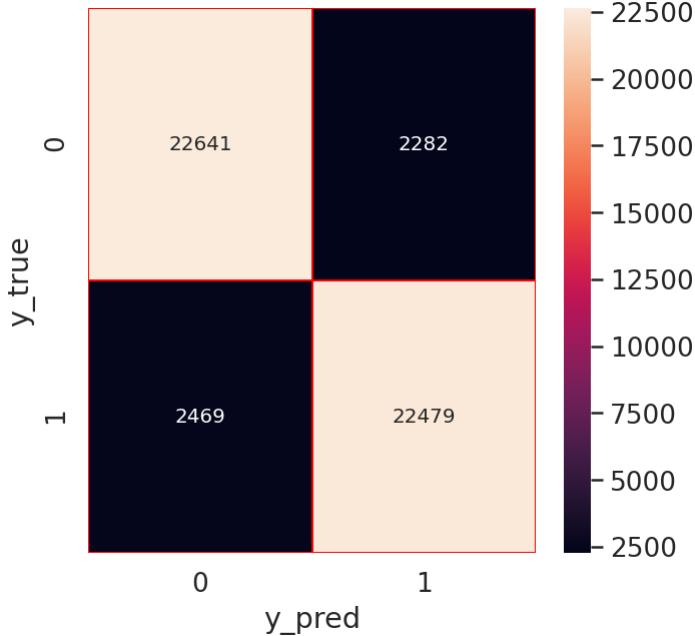
```
[[0.99064845 0.00935157]]  
'Bad'
```

## ▼ Step 8: Confusion matrix

```
1 # confusion matrix in sklearn  
2 from sklearn.metrics import confusion_matrix  
3 from sklearn.metrics import classification_report  
4 from keras import backend as K  
5 from keras.layers.convolutional import Convolution2D, MaxPooling2D  
6 from keras.preprocessing.image import ImageDataGenerator  
7  
8 labels=['positive','negative']  
9  
10 #Confusion Matrix and Classification Report  
11 #confusion matrix  
12 y_pred=model.predict(X_test)  
13 #y_pred = (y_pred > 0.5)  
14  
15 #cm =confusion_matrix(y_test.argmax(axis=1), y_pred.argmax(axis=1))  
16 #print("Confusion Matrix : \n", cm)  
  
1 confusion_df = pd.DataFrame(confusion_matrix(y_test.argmax(axis=1), y_pred.argmax(axis=1))  
2                 columns=["Predicted Class " + str(class_name) for class_name in [0,1]],  
3                 index = ["Class " + str(class_name) for class_name in [0,1]])  
4  
5 print(confusion_df)  
  
          Predicted Class 0  Predicted Class 1  
Class 0           22641           2282
```

```
Class 1          2469          22479
```

```
1 f, ax=plt.subplots(figsize=(5,5))
2 sns.heatmap(cm, annot=True, linewidths=0.5, linecolor="red", fmt=".0f", ax=ax)
3 plt.xlabel("y_pred")
4 plt.ylabel("y_true")
5 plt.show()
```



## ▼ Step 9: Classification Report.

Precision Score

TP – True Positives

FP – False Positives

Precision – Accuracy of positive predictions.

$$\text{Precision} = \text{TP}/(\text{TP} + \text{FP})$$

```
1 from sklearn.metrics import precision_score
2
3 print("Precision score: {}".format(precision_score(y_test.argmax(axis=1), y_pred.argmax(axis=1))))
```

Precision score: 0.9078389402689714

Recall Score

FN – False Negatives

Recall (aka sensitivity or true positive rate):

Fraction of positives That were correctly identified.

Recall = TP/(TP+FN)

```
1 from sklearn.metrics import recall_score  
2  
3 print("Recall score: {}".format(recall_score(y_test.argmax(axis=1),y_pred.argmax(axis=1)))  
  
Recall score: 0.9010341510341511
```

## F1 Score

F1 Score (aka F-Score or F-Measure) – A helpful metric for comparing two classifiers. F1 Score takes into account precision and the recall. It is created by finding the the harmonic mean of precision and recall.

$F1 = 2 \times (\text{precision} \times \text{recall}) / (\text{precision} + \text{recall})$

```
1 from sklearn.metrics import f1_score  
2  
3 print("F1 Score: {}".format(f1_score(y_test.argmax(axis=1),y_pred.argmax(axis=1)))  
  
F1 Score: 0.9044237462029008
```

## ▼ Classification Report

Report which includes Precision, Recall and F1-Score.

```
1 print("CLASSIFICATION REPORT: \n")  
2 cr = classification_report(y_test.argmax(axis=1), y_pred.argmax(axis=1))  
3 print(cr)
```

CLASSIFICATION REPORT:

	precision	recall	f1-score	support
0	0.90	0.91	0.91	24923
1	0.91	0.90	0.90	24948
accuracy			0.90	49871
macro avg	0.90	0.90	0.90	49871
weighted avg	0.90	0.90	0.90	49871

## Step 10: Conclusion

1. We started with analyzing and describing the data and providing a few metrics.

2. We then proceeded with data preprocessing and creating the word clouds for the positive and negative reviews.
3. We used the USE for word embeddings and tested our data with the Keras model.
4. We achieved an accuracy of 91% and 90% for the negative and positive reviews respectively.
5. The full confusion matrix and the classification reports are detailed as above.
6. We also checked few test data as well.

## Step 11: Reference:

<https://curiously.com/posts/sentiment-analysis-with-tensorflow-2-and-keras-using-python/>

<https://realpython.com/python-keras-text-classification/>

<https://joshlawman.com/metrics-classification-report-breakdown-precision-recall-f1/>

## Step 12: The END