# LSTM(Long Short Term Memory) on Amazon Fine Food Reviews

The Amazon Fine Food Reviews dataset consists of reviews of fine foods from Amazon.

Number of reviews: 568,454 Number of users: 256,059 Number of products: 74,258 Timespan: Oct 1999 - Oct 2012 Number of Attributes/Columns in data: 10

## Attribute Information:

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

## Objective:

Given a review, determine whether the review is positive (Rating of 4 or 5) or negative (rating of 1 or 2).

## ▾ Loading,Cleaning & Preprocessing the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

In order to load the data, We have used the SQLITE dataset as it easier to query the data and visualise the data efficiently.

Here as we only want to get the global sentiment of the recommendations (positive or negative), we will purposefully ignore all Scores equal to 3. If the score id above 3, then the recommendation wil be set to "positive". Otherwise, it will be set to "negative".*italicized text*

```
1 # This Python 3 environment comes with many helpful analytics libraries installed
2 # It is defined by the kaggle/python Docker image: https://github.com/kaggle/docker-python
3 # For example, here's several helpful packages to load
4
5 import numpy as np # linear algebra
6 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
7
8 # Input data files are available in the read-only "../input/" directory
9 # For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input directory
10
11 import os
12 for dirname, _, filenames in os.walk('/kaggle/input'):
13     for filename in filenames:
14         print(os.path.join(dirname, filename))
15
16 # You can write up to 5GB to the current directory (/kaggle/working/) that gets preserved as output when you creat
17 # You can also write temporary files to /kaggle/temp/, but they won't be saved outside of the current session
```

```
/kaggle/input/amazon-fine-food-reviews/database.sqlite
/kaggle/input/amazon-fine-food-reviews/hashes.txt
/kaggle/input/amazon-fine-food-reviews/Reviews.csv
```

## ▾ Import All Required Libraries

```
1 %matplotlib inline
```

```
2 import warnings
3
4 warnings.filterwarnings("ignore")
```

```
1 import sqlite3
2 import numpy as np
3 from keras.models import Sequential
4 from keras.layers import Dense
5 from keras.layers import LSTM,Dropout
6 from keras.layers.embeddings import Embedding
7 from keras.preprocessing import sequence
8 import string
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11 import pickle
12 from tqdm import tqdm
13 import os
14 from tqdm import tqdm
15 import os
16 import pdb
17 import pickle
18 import re
19 import nltk
20 from collections import Counter
21 from itertools import islice
22 from sklearn.model_selection import train_test_split
23 from keras.models import Sequential
24 from keras.preprocessing import sequence
25 from keras.initializers import he_normal
26 from keras.layers import BatchNormalization, Dense, Dropout, Flatten, LSTM
27 from keras.layers.embeddings import Embedding
28 from keras.regularizers import L1L2
```

## ▼ Connect to sqlite and fetch the data using SQL Query

```
1 con=sqlite3.connect('../input/amazon-fine-food-reviews/database.sqlite')
2
3 filtered_data=pd.read_sql_query("""SELECT * FROM Reviews WHERE Score !=3""",con)
4 filtered_data.head(3)
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 | 5 | 13038 |

```
1 filtered_data.shape
```

   (525814, 10)

```
1 def partition(x):
2   if x < 3 :
3     return 'negative'
4   return 'positive'
5
6 actualScore=filtered_data['Score']
```

```
7 positive_negative=actualScore.map(partition)
8 filtered_data['Score']=positive_negative
9 print("Number of datapoints",filtered_data.shape)
10 filtered_data.head(3)
```

Number of datapoints (525814, 10)

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 | positive | 130 |

```
1 display = pd.read_sql_query("""
2 SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
3 FROM Reviews
4 GROUP BY UserId
5 HAVING COUNT(*)>1
6 """, con)
```

```
1 print(display.shape)
2 display.head(3)
```

(80668, 7)

| | UserId | ProductId | ProfileName | Time | Score | Text | COUNT(*) |
|---|---|---|---|---|---|---|---|
| **0** | #oc-R115TNMSPFT9I7 | B005ZBZLT4 | Breyton | 1331510400 | 2 | Overall its just OK when considering the price... | 2 |
| **1** | #oc-R11D9D7SHXIJB9 | B005HG9ESG | Louis E. Emory "hoppy" | 1342396800 | 5 | My wife has recurring extreme muscle spasms, u... | 3 |

```
1 display[display["UserId"]=='AZY10LLTJ71NX']
```

| | UserId | ProductId | ProfileName | Time | Score | Text | COUNT(*) |
|---|---|---|---|---|---|---|---|
| **80638** | AZY10LLTJ71NX | B001ATMQK2 | undertheshrine "undertheshrine" | 1296691200 | 5 | I bought this 6 pack because for the price tha... | 5 |

```
1 display['COUNT(*)'].sum()
```

393063

```
1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND UserId="AR5J8UI46CURR"
5 ORDER BY ProductID
6 """, con)
7 display.head()
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | |
|---|---|---|---|---|---|---|---|---|
| **0** | 78445 | B000HDL1RQ | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199 |
| **1** | 138317 | B000HDOPYC | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199 |

```
1 sorted_data=filtered_data.sort_values('ProductId',axis=0,ascending=True,inplace=False,kind='quicksort',na_position
2
3 final_data=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"},keep='first',inplace=False)
4 final_data.shape
```

(364173, 10)

```
1 (final_data['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

69.25890143662969

```
1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND Id=44737 OR Id=64422
5 ORDER BY ProductID
6 """, con)
7
8 display.head()
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | |
|---|---|---|---|---|---|---|---|---|
| **0** | 64422 | B000MIDROQ | A161DK06JJMCYF | J. E. Stephens "Jeanne" | 3 | 1 | 5 | 1224 |

```
1 final_data=final_data[final_data.HelpfulnessNumerator<=final_data.HelpfulnessDenominator]
```

## ▾ Count the Positive and Negative Review Counts

```
1 print(final_data.shape)
2
3 #How many positive and negative reviews are present in our dataset?
4 final_data['Score'].value_counts()
```

(364171, 10)
positive     307061
negative      57110
Name: Score, dtype: int64

## ▾ Import nltk library

```
1
2 from nltk.corpus import stopwords
3 stopping_words=set(stopwords.words('english'))
4 print(stopping_words)
```

{'yourselves', "won't", 'ma', 'herself', 'them', 'didn', 'should', 'your', 'an', "haven't", 'wasn', 'his', 'who

```python
1 def clean_html(text):
2     clean_r = re.compile('<,*?>')
3     clean_text = re.sub(clean_r,'',text)
4     return clean_text
5
6 def Clean_punc(text):
7     clean_sentence = re.sub(r'[?|!|\'|"|#]',r' ',text)
8     clean_data = re.sub(r'[.|,|)|(|\|/)]',r' ',clean_sentence)
9     return clean_data
```

```python
1 stem_no = nltk.stem.SnowballStemmer('english')
2
3 if not os.path.isfile('final_data.sqlite'):
4     final_string=[]
5     all_positive_words=[]
6     all_negative_words=[]
7     for i,sentence in enumerate(tqdm(final_data['Text'].values)):
8         filtered_sentence=[]
9         sent_without_html_tags=clean_html(sentence)
10         #pdb.set_trace()
11         for w in sent_without_html_tags.split():
12             for cleaned_words in Clean_punc(w).split():
13                 if ((cleaned_words.isalpha()) & (len(cleaned_words) > 2)):
14                     if(cleaned_words.lower() not in stopping_words) :
15                         stemming=(stem_no.stem(cleaned_words.lower())).encode('utf8')
16                         filtered_sentence.append(stemming)
17                         if(final_data['Score'].values)[i]=='positive':
18                             all_positive_words.append(stemming)
19                         if(final_data['Score'].values)[i]=='negative':
20                             all_negative_words.append(stemming)
21         str1 = b" ".join(filtered_sentence)
22         final_string.append(str1)
23
24     final_data['Cleaned_text']=final_string
25     final_data['Cleaned_text']=final_data['Cleaned_text'].str.decode("utf-8")
26
27     conn = sqlite3.connect('final_data.sqlite')
28     cursor=conn.cursor
29     conn.text_factory = str
30     final_data.to_sql('Reviews',conn,schema=None,if_exists='replace',index=True,index_label=None,chunksize=None,dt
31     conn.close()
32
33
34     with open('positive_words.pkl','wb') as f :
35         pickle.dump(all_positive_words,f)
36     with open('negative_words.pkl','wb') as f :
37         pickle.dump(all_negative_words,f)
```

100%|████████████| 364171/364171 [08:20<00:00, 727.25it/s]

```python
1 final_data['total_words'] = [len(x.split()) for x in final_data['Cleaned_text'].tolist()]
```

```python
1 final_data.head(3)
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score |
|---|---|---|---|---|---|---|---|
| **138706** | 150524 | 0006641040 | ACITT7DI6IDDL | shari zychinski | 0 | 0 | positive |
| **138688** | 150506 | 0006641040 | A2IW4PEEKO2R0U | Tracy | 1 | 1 | positive |
| **138689** | 150507 | 0006641040 | A1S4A3IQ2MU7V4 | sally sue "sally sue" | 1 | 1 | positive |

```
1 final_data.shape
```

(364171, 12)

## ▾ Pick the top 100K Data Points from the Dataset

```
1 final_data_100K=final_data[0:100000]
2 amazon_polarity_labels=final_data_100K['Score'].values
3 final_data_100K.head(2)
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score |
|---|---|---|---|---|---|---|---|
| **138706** | 150524 | 0006641040 | ACITT7DI6IDDL | shari zychinski | 0 | 0 | positive |
| **138688** | 150506 | 0006641040 | A2IW4PEEKO2R0U | Tracy | 1 | 1 | positive |

```
1 X_Data_Text=final_data_100K['Cleaned_text']
```

```
1 Y_Data_Score=final_data_100K['Score']
```

```
1 X_Data_Text.head(3)
```

```
    138706    witti littl book make son laugh loud recit car...
```

```
1 X_Data_Text.index=[x for x in range(0,10**5)]
```

## Get the list of all the words from the "Cleaned Text"

```
1 list_of_words = []
2 for i in tqdm(X_Data_Text):
3     list_of_words += i
4 list_of_words = ''.join(list_of_words)
5 list_of_all_words=list_of_words.split()
```

100%|████████████| 100000/100000 [00:00<00:00, 318290.67it/s]

```
1 list_of_all_words[0]
```

'witti'

## Convert the List of words into SET and List

```
1 vocab=set(list_of_all_words)
2 vocab_list=list(vocab)
```

```
1 type(vocab)
```

set

```
1 vocab_list[0]
```

'lisa'

```
1 dict_freq={}
2 for x in tqdm(vocab_list):
3     dict_freq[x]=list_of_all_words.count(x)
```

100%|████████████| 103077/103077 [1:34:36<00:00, 18.16it/s]

```
1 dict_freq
```

SHOW HIDDEN OUTPUT

```
1 (pd.DataFrame.from_dict(data=dict_freq, orient='index')
2     .to_csv('dict_file.csv', header=False))
```

## Use Counter function to count the number of Unique words and pick the top 5000 words and add in the dictionary

```
1 print("Number of sentences in complete dataset : ",len(list_of_all_words))
2
3 counts_words = Counter(list_of_all_words)
4 print("Number of unique words present : ",len(counts_words.most_common()))
5 vocab_size = len(counts_words.most_common()) + 1
6 top_words_count = 5000
```

```
 7 common_words = counts_words.most_common(top_words_count)
 8
 9 word_index = dict()
10 i = 1
11 for word,frequency in common_words:
12     word_index[word] = i
13     i += 1
14
15 print()
16 print("Top 25 words with their frequencies:")
17 print(counts_words.most_common(25))
18 print()
19 print("Top 25 words with their index:")
20 print(list(islice(word_index.items(), 25)))
```

```
Number of sentences in complete dataset :   3488783
Number of unique words present :   103077

Top 25 words with their frequencies:
[('like', 38877), ('tast', 37698), ('tea', 33156), ('good', 29934), ('product', 29213), ('use', 29058), ('flavor

Top 25 words with their index:
[('like', 1), ('tast', 2), ('tea', 3), ('good', 4), ('product', 5), ('use', 6), ('flavor', 7), ('one', 8), ('gre
```

```
1 type(X_Data_Text)
```

```
pandas.core.series.Series
```

## Create a new Column called "CleanedText_Index" and add the index of each word occured in the "Cleaned_text"

```
 1 def use_index(row):
 2     holder = []
 3     for word in row['Cleaned_text'].split():
 4         if word in word_index:
 5             holder.append(word_index[word])
 6         else:
 7             holder.append(0)
 8     return holder
 9
10
11 final_data_100K['CleanedText_Index'] = final_data_100K.apply(lambda row: use_index(row),axis=1)
12 final_data_100K.head(5)
```

| Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Scor |
|---|---|---|---|---|---|---|
| **138706** | 150524 | 0006641040 | ACITT7DI6IDDL | shari zychinski | 0 | 0 | positiv |
| **138688** | 150506 | 0006641040 | A2IW4PEEKO2R0U | Tracy | 1 | 1 | positiv |
| **138689** | 150507 | 0006641040 | A1S4A3IQ2MU7V4 | sally sue "sally sue" | 1 | 1 | positiv |

▾ Convert the Score to 1 & 0.

positive - 1

negative - 0

```
1 final_data_100K['Score'] = final_data_100K['Score'].map(lambda x : 1 if x == 'positive' else 0)
2 final_data_100K.head(5)
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score |
|---|---|---|---|---|---|---|---|
| **138706** | 150524 | 0006641040 | ACITT7DI6IDDL | shari zychinski | 0 | 0 | 1 |

```
1 final_data_100K['CleanedText_Index'].head(2)
```

```
138706    [0, 27, 932, 11, 384, 1976, 2578, 0, 1196, 123...
138688    [995, 247, 0, 932, 551, 23, 0, 988, 2594, 10, ...
Name: CleanedText_Index, dtype: object
```

## Split the Data into TEST and Train

| 138689 | 150507 | 0006641040 | ATS4ASIQZMU7V4 | "sally sue" | 1 | 1 | 1 |

```
1 X_Train, X_Test, Y_Train, Y_Test = train_test_split(final_data_100K['CleanedText_Index'].values,final_data_100K['S
```

```
1 X_Train[6]
```

```
    [365,
     1773,
     266,
```

## Apply Padding

```
     384,
```

```
1 from keras.preprocessing import sequence
2
3 max_review_length = 600
4 X_Train = sequence.pad_sequences(X_Train, maxlen=max_review_length)
5 X_Test = sequence.pad_sequences(X_Test, maxlen=max_review_length)
6
7 print(X_Train.shape)
8 print(X_Train[1])
```

```
(70000, 600)
[   0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0    0
    0    0    0    0    0    0    0    0    0    0    0    0    0  995
  247    0  932  551   23    0  988 2594   10  384   10  149  558  111
  562  366    0  113  211 3563   82   63  252   68 1751  166]
```

```
1 from tensorflow.python.client import device_lib
2 print(device_lib.list_local_devices())
```

```
[name: "/device:CPU:0"
device_type: "CPU"
memory_limit: 268435456
locality {
}
incarnation: 6202525998069733733
, name: "/device:XLA_CPU:0"
device_type: "XLA_CPU"
memory_limit: 17179869184
locality {
}
incarnation: 15694142288273005724
physical_device_desc: "device: XLA_CPU device"
, name: "/device:GPU:0"
device_type: "GPU"
memory_limit: 15695549568
locality {
  bus_id: 1
  links {
  }
}
incarnation: 16992963004178590259
physical_device_desc: "device: 0, name: Tesla P100-PCIE-16GB, pci bus id: 0000:00:04.0, compute capability: 6.0'
, name: "/device:XLA_GPU:0"
device_type: "XLA_GPU"
memory_limit: 17179869184
```

## ▾ Apply LSTM 1

```
physical_device_desc:  device. XLA_GPU device
```

```python
1  import numpy
2  numpy.random.seed(7)
3
4  embedding_vecor_length = 32
5  model = Sequential()
6  model.add(Embedding(top_words_count+1, embedding_vecor_length, input_length=max_review_length))
7  model.add(LSTM(100))
8  model.add(Dense(1, activation='sigmoid'))
9  model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
10 print(model.summary())
```

```
Model: "sequential"
_____
Layer (type)                 Output Shape              Param #
=================================================================
embedding (Embedding)        (None, 600, 32)           160032
_____
lstm (LSTM)                  (None, 100)               53200
_____
dense (Dense)                (None, 1)                 101
=================================================================
Total params: 213,333
Trainable params: 213,333
Non-trainable params: 0
_____
None
```

```python
1  model.fit(X_Train, Y_Train, epochs=10, batch_size=64)
2  # Final evaluation of the model
3  scores = model.evaluate(X_Test, Y_Test, verbose=0)
4  print("Accuracy: %.2f%%" % (scores[1]*100))
```

```
Epoch 1/10
1094/1094 [==============================] - 35s 32ms/step - loss: 0.2506 - accuracy: 0.9027
Epoch 2/10
1094/1094 [==============================] - 34s 31ms/step - loss: 0.1884 - accuracy: 0.9275
Epoch 3/10
1094/1094 [==============================] - 35s 32ms/step - loss: 0.1729 - accuracy: 0.9331
Epoch 4/10
1094/1094 [==============================] - 34s 31ms/step - loss: 0.1582 - accuracy: 0.9391
Epoch 5/10
1094/1094 [==============================] - 35s 32ms/step - loss: 0.1433 - accuracy: 0.9451
Epoch 6/10
1094/1094 [==============================] - 34s 32ms/step - loss: 0.1301 - accuracy: 0.9502
Epoch 7/10
1094/1094 [==============================] - 34s 31ms/step - loss: 0.1176 - accuracy: 0.9563
Epoch 8/10
```

```
1 print('Test score: ',scores[0])
2 print('Test accuracy: ',scores[1])
```

```
Test score:  0.3073323667049408
Test accuracy:  0.9080333113670349
```

```
 1 embedding_vecor_length = 32
 2 model2 = Sequential()
 3 model2.add(Embedding(top_words_count+1, embedding_vecor_length, input_length=max_review_length))
 4 model2.add(LSTM(100,return_sequences=True))
 5 model2.add(Dropout(0.25))
 6 model2.add(LSTM(80))
 7 model2.add(Dropout(0.5))
 8 model2.add(Dense(1, activation='sigmoid'))
 9 model2.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
10 print(model2.summary())
```

```
Model: "sequential_3"
_____
Layer (type)                 Output Shape              Param #
=================================================================
embedding_2 (Embedding)      (None, 600, 32)           160032

lstm_3 (LSTM)                (None, 600, 100)          53200

dropout_1 (Dropout)          (None, 600, 100)          0

lstm_4 (LSTM)                (None, 80)                57920

dropout_2 (Dropout)          (None, 80)                0

dense_1 (Dense)              (None, 1)                 81
=================================================================
Total params: 271,233
Trainable params: 271,233
Non-trainable params: 0
_____
None
```

```
1 model2.fit(X_Train, Y_Train, epochs=10, batch_size=64)
2 # Final evaluation of the model
3 scores = model2.evaluate(X_Test, Y_Test, verbose=0)
4 print("Accuracy: %.2f%%" % (scores[1]*100))
```

```
Epoch 1/10
1094/1094 [==============================] - 62s 56ms/step - loss: 0.2522 - accuracy: 0.9031
Epoch 2/10
1094/1094 [==============================] - 61s 56ms/step - loss: 0.1947 - accuracy: 0.9255
Epoch 3/10
1094/1094 [==============================] - 62s 56ms/step - loss: 0.1744 - accuracy: 0.9338
Epoch 4/10
1094/1094 [==============================] - 62s 57ms/step - loss: 0.1582 - accuracy: 0.9408
Epoch 5/10
1094/1094 [==============================] - 63s 57ms/step - loss: 0.1411 - accuracy: 0.9471
Epoch 6/10
1094/1094 [==============================] - 62s 57ms/step - loss: 0.1289 - accuracy: 0.9530
Epoch 7/10
```

```
1 print('Test score: ',scores[0])
2 print('Test accuracy: ',scores[1])
```

```
Test score:  0.3317498564720154
Test accuracy:  0.9067999720573425
1094/1094 [------------------------------] - 62s 56ms/step - loss: 0.0800 - accuracy: 0.9733
```

# CONCLUSION

## After using LSTM Models , below are the conlusion made

1. LSTM without Dropouts --> Accuracy is 90.80%
2. LSTM with Dropout --> Accuracy is 90.67%