Amazon Customer Reviews (US Based Dataset) CIS 4130 Big Data Technologies Semester Project

Prepared By
Pranesh Chitrakar
pranesh.chitrakar@baruchmail.cuny.edu

Main Dataset

https://www.kaggle.com/datasets/cynthiare mpel/amazon-us-customer-reviews-dataset

Introduction:

Amazon Customer Reviews dataset is a collection of millions of Amazon's user reviews from 1995 to 2015. Over hundred million customers have conveyed their experiences regarding the merchandises on the Amazon.com website. The data helps us to study the properties and evolution of customer reviews on how the users express their views on a product/s.

This dataset consists of 15 columns:

- 1. marketplace- 2 letters country code (AU for Australia, US for United States)
- customer_id- random identifier that can be used to aggregate reviews by a single user
- 3. review_id- unique ID of the review
- 4. product_id- unique Product ID the review pertains to
- 5. product_parent- random identifier that can be used to aggregate reviews for the same product.
- 6. product_title- title of the product
- 7. product category- broad product category that can be used to group reviews (
- 8. star_rating- 1-5 star rating of the review
- 9. helpful_votes- no. of helpful votes
- 10. total_vote- no. of helpful votes received
- 11. vine- review was written as part of the Vine program.
- 12. verified purchase- the review is on a verified purchase.
- 13. review headline- the title of the review.
- 14. review body- the review text.
- 15. review_date- the date the review was written

This dataset has a bunch of opportunities to that come in had for Amazon Sellers, Customers as well as the site itself. By addressing the reviews, products can be recommended to certain parties and forecast(predict) the demand.

My focus for this project will be predicting the review ratings for respective products from the dataset.

For this project, I will be using Amazon Web Service (AWS) Management Console and it's platform like EC2, S3, and so on to manage and monitor the dataset. The Amazon Customer Reviews dataset is approximately 22GB, and to download and re-upload manually from Kaggle(kaggle.com/) would take me hours. I am using Amazon S3(Simple Storage Service) to create a bucket where the dataset can be extracted and stored in less amount of time.

The steps I followed are below:

- 1. Set up AWS Command Line Interface (CLI) in my EC2 instance using my Access Key ID and Secret Access Key.
- 2. In Amazon EC2 instance with Amazon Linux, I first created a bucket named 'mydata-bucket-pc' to store the dataset.

```
$ aws s3api create-bucket --bucket my-data-bucket-pc --region us-east-2
\ --create-bucket-configuration LocationConstraint=us-east-2
```

- Followed the steps for downloading data directly from Kaggle to Amazon S3 using Kaggle API
- 4. After the installation and configuration, I used \$ kaggle datasets list to search the data set

```
[ec2-user@ip-172-31-39-117 ~]$ kaggle datasets list

[ef teCount usabilityRating title size lastUpdated

iamsouravbanerjee/world-population-dataset World Population Dataset 17KB 2022-08-31 11:20:04
358 1.0

whenamancodes/hr-employee-attrition Employee Analysis | Attrition Report 50KB 2022-09-12 10:46:33
29 1.0

pantanjali/unemployment-dataset Unemployment dataset 17KB 2022-09-08 08:26:10
```

5. To download my data, I used \$ kaggle datasets download command to fetch the dataset and to pipe the output and direct it to my S3 bucket, I used aws s3 cp command.

```
$ kaggle datasets download --quiet -d cynthiarempel/amazon-us-customer-
reviews-dataset -p - | aws s3 cp - s3://my-data-bucket-pc/data.zip
```

6. To check my S3 bucket to see if the file was downloaded

```
$ aws s3 ls s3://my-data-bucket-pc/
```

7. Testing the Boto3 module in Python

As we can see above, 'amazon-customer-review', and 'my-data-bucket-pc' are the bucket I created for this project.

8. Unzipping files with S3

```
>>> import zipfile
>>> import boto3
>>> from io import BytesIO
>>> bucket="my-data-bucket-pc"
>>> zipfile_to_unzip="data.zip"
>>> s3 client= boto3.client('s3', use ssl=False)
>>> s3_resource= boto3.resource('s3')
>>> zip obj = s3 resource.Object(bucket name=bucket, key=zipfile to unzip)
>>> buffer = BytesIO(zip_obj.get() ["Body"].read())
>>> z = zipfile.ZipFile(buffer)
>>> for filename in z.namelist():
       print('Working on ' + filename)
       .
s3_resource.meta.client.upload_fileobj(z.open(filename), Bucket=bucket, Key=f'{filename}')
Working on amazon_reviews_multilingual_US_v1_00.tsv
Working on amazon_reviews_us_Apparel_v1_00.tsv
Working on amazon_reviews_us_Automotive_v1_00.tsv
Working on amazon_reviews_us_Baby_v1_00.tsv
Working on amazon_reviews_us_Beauty_v1_00.tsv
```

In this step, I imported the ZIP file into my EC2 instance. While doing so, I had to change the instance type from t2.micro to t2.2xlarge(32GB).

```
import pandas as pd
df=pd.read_table('s3://my-data-bucket-
pc/amazon_reviews_multilingual_US_v1_00.tsv',lineterminator='\
n',error_bad_lines=False)
```

>>> df= pd.read_table('s3://my-data-bucket-pc/amazon_reviews_multilingual_US_v1_00.tsv',lineterminator='\n',error_bad_lines=False)
sys:1: FutureWarning: The error_bad_lines argument has been deprecated and will be removed in a future version.

Here, I had to use the "lineterminator='\n',error_bad_lines=False)" in order to skip the lines that were causing errors to load so that we can move with our pandas DataFrame.

```
>> print(df.info)
bound method DataFrame.info of
US 53096384
                                                                               marketplace customer_id ...
this is the first 8 issues of the series. it i...
I've always been partial to immutable laws. Th...
This is a book about first contact with aliens...
This is quite possibly *the* funniest book I h...
The story behind the book is almost better tha...
                                                                                                                                                                                                                                   review_body review_date
                                                                                                                                                                                          1995-08-13
                                                53096339
53096332
53096335
51747709
                                                                                                                                                                                          1995-08-17
1995-08-30
                                  US
                                                       52303
6900881
                                                                                                                                                           very fun game
                                                     565563
254421
                                                                              This is my first book from Jodi, and she has b...

Good movie!
6900883
                                                                              Amazing show that runs close to history. Love ...
6900885
                                                     146004
```

df.dtypes

Ten of the variables are categorical(which are labelled as 'object') while the remaining variables are numerical (which are labelled as 'int64'(2) and 'float64'(3)).

dtypes: float64(3), int64(2), object(10)

```
>>> df.dtypes
marketplace
                       object
customer id
                        int64
review id
                       object
product id
                       object
product parent
                        int64
product title
                       object
product_category
                       object
star_rating
                      float64
helpful votes
                      float64
total votes
                      float64
                       object
vine
verified purchase
                       object
review headline
                       object
review_body
                       object
review date
                       object
dtype: object
```

Grouping the data

```
results_star_rating =
df.groupby('customer_id').star_rating.agg(['count', 'min',
'max', 'mean', 'std'])
>>> print(results star rating)
```

```
>> results_star_rating = df.groupby('customer_id').star_rating.agg(['count', 'min', 'max', 'mean', 'std'])
>>> print(results_star_rating)

count min max n
                                mean
                                             std
customer_id
10001
                                       0.000000
                  2 5.0 5.0
                                  5.0
                                 4.0
5.0
10018
                     4.0 4.0
                                            NaN
10019
                          5.0
                     5.0
                                             NaN
                     5.0
10020
                          5.0
                                  5.0
                                             NaN
10022
                  2
                           5.0
                                       0.000000
                                  5.0
                           1.0
53096567
                          1.0
4.0
53096575
                     1.0
                                  1.0
                                             NaN
53096582
                     4.0
                                  4.0
                                       0.000000
                     2.0
                           5.0
53096584
                                       1.414214
```

```
>>> results_helpful_votes =
df.groupby('customer_id').helpful_votes.agg(['count', 'min',
'max', 'mean', 'std'])
```

>>> print(results helpful votes)

```
>>> results helpful_votes = df.groupby('customer_id').helpful_votes.agg(['count', 'min',
>>> print(results_helpful_votes)
                          min
                                                           std
              count
                                   max
                                            mean
customer id
10001
                          0.0
                   2
                                    1.0
                                             0.5
                                                     0.707107
10018
                          0.0
                                    0.0
                                             0.0
                                                           NaN
                          0.0
10019
                                   0.0
                                             0.0
                                                           NaN
10020
                          0.0
                                    0.0
                                             0.0
                                                           NaN
                   2
                                                     0.000000
10022
                          0.0
                                    0.0
                                             0.0
                          0.0
                 ...
53096567
                                    0.0
                                             0.0
                                                           NaN
                      4054.0
53096575
                                4054.0
                                          4054.0
                                                           NaN
                                                   746.704761
53096582
                   2
                          6.0
                                1062.0
                                           534.0
53096584
                                            52.5
3.0
                          0.0
                                 209.0
                                                   104.334398
                          3.0
                                    3.0
53096589
                                                           NaN
[4097784 rows x 5 columns]
```

```
>>> results_total_votes =
df.groupby('customer_id').total_votes.agg(['count', 'min',
'max', 'mean', 'std'])
>>> print(results total votes)
```

```
>>> results_total_votes = df.groupby('customer_id').total_votes.agg(['count', 'min', 'max', 'mean', 'std'])
>>> print(results_total_votes)
              count
                                            mean
                         min
customer_id
10001
                                            1.00
                                                     1.414214
                                   2.0
                                            0.00
10018
                          0.0
                                   0.0
                                                           NaN
10019
                          0.0
                                   0.0
                                            0.00
                                                           NaN
10020
                          0.0
                                   0.0
                                            0.00
                                                           NaN
10022
                   2
                                                     0.000000
                          0.0
                                   0.0
                                            0.00
53096567
                   1
                          0.0
                                   0.0
                                            0.00
                                                           NaN
                      4756.0
53096575
                               4756.0
                                         4756.00
                                                           NaN
                   2
53096582
                               1114.0
                                                   783.474314
                          6.0
                                          560.00
53096584
                                           61.75
                          0.0
                                 246.0
                                                   122.834238
53096589
                          3.0
                                   3.0
                                            3.00
                                                           NaN
[4097784 rows x 5 columns]
```

Finding out the total NaN values throughout the dataframe

```
>>> total_nan_value = df.isna().sum()
>>> print(total nan value)
```

```
>>> total_nan_value = df.isna().sum().sum()
>>> print(total_nan_value)
480
```

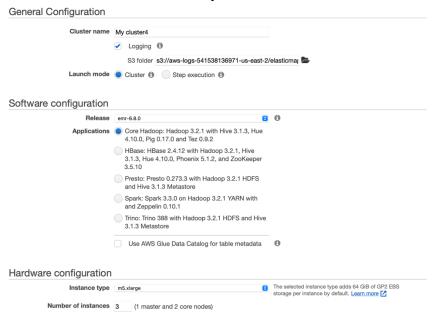
For the review_date variable, since it is categorical, I still need to convert it to datetime format to work with finding the min and max date using python function. The review dates in my dataframe are in ascending order, therefore in the image below, we can see the min and max date.

```
>>> df.loc[:,"review_date"]
           1995-08-13
           1995-08-17
           1995-08-30
           1995-09-11
           1995-10-17
6900881
           2015-08-31
6900882
           2015-08-31
6900883
           2015-08-31
6900884
           2015-08-31
6900885
           2015-08-31
Name: review_date, Length: 6900886, dtype: object
```

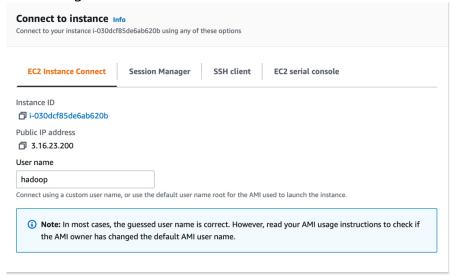
For this milestone, I used Amazon EMR to create cluster that implements Apache Hadoop and Spark.

Launching the EMR Cluster

I created a Cluster named "My cluster4" on Amazon EMR



Then navigated the instances from EC2



Writing Hadoop MapReduce program in Python

```
[hadoop@ip-172-31-23-236 ~]$ ls -l
total 0
[hadoop@ip-172-31-23-236 ~]$ nano mapper.py
[hadoop@ip-172-31-23-236 ~]$ nano reducer.py
[hadoop@ip-172-31-23-236 ~]$ ls -l
total 8
-rw-rw-r-- 1 hadoop hadoop 423 Dec 16 20:43 mapper.py
-rw-rw-r-- 1 hadoop hadoop 921 Dec 16 20:47 reducer.py
```

Making a new directory folder on HDFS

[hadoop@ip-172-31-23-236 ~]\$ hdfs dfs -mkdir hdfs:///bigdata Now we are going to use the dataset from S3 bucket

s3-dist-cp --src s3://my-data-bucket-pc/amazon_reviews_multilingual_US_v1_00.tsv --dest hdfs:///bigdata

```
[hadoop@ip-172-31-23-236 -]$ s3-dist-cp --src s3://my-data-bucket-pc/amazon reviews multilingual US v1 00.tsv --dest hdfs:///bigdata
2022-12-16 21:18:38,827 INFO s3distcp,Main: Running with args: -libjars /usr/share/aws/emr/s3-dist-cp/lib/aopalliance-1.0.jar,/usr/share/aws/emr/s3-dist-cp/lib/aopalliance-1.0.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.10.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-ervlet-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.jar,/usr/share/aws/emr/s3-dist-cp/lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/guice-4.2.lib/gui
```

```
2022-12-16 21:18:45,018 INFO mapreduce.Job: Running job: job_1671223184898_0001
2022-12-16 21:18:52,338 INFO mapreduce.Job: Job job_1671223184898_0001 running in uber mode: false
2022-12-16 21:18:52,339 INFO mapreduce.Job: map 0% reduce 0%
2022-12-16 21:18:58,386 INFO mapreduce.Job: map 100% reduce 0%
2022-12-16 21:19:04,414 INFO mapreduce.Job: map 100% reduce 33%
2022-12-16 21:19:05,419 INFO mapreduce.Job: map 100% reduce 67%
2022-12-16 21:19:16,458 INFO mapreduce.Job: map 100% reduce 100%
2022-12-16 21:20:45,730 INFO mapreduce.Job: Job job_1671223184898_0001 completed successfully
2022-12-16 21:20:45,814 INFO mapreduce.Job: Counters: 59
```

Below, I have pasted the screenshot of the counters and Shuffle Errors that it gave me as a result

```
Map-Reduce Framework
File System Counters
                                                                            Map input records=1
         FILE: Number of bytes read=163
                                                                            Map output records=1
         FILE: Number of bytes written=979031
                                                                            Map output bytes=177
         FILE: Number of read operations=0
                                                                            Map output materialized bytes=151
         FILE: Number of large read operations=0
                                                                            Input split bytes=169
         FILE: Number of write operations=0
                                                                            Combine input records=0
         HDFS: Number of bytes read=426
                                                                            Combine output records=0
         HDFS: Number of bytes written=3629753164
                                                                            Reduce input groups=1
         HDFS: Number of read operations=20
                                                                            Reduce shuffle bytes=151
                                                                            Reduce input records=1
         HDFS: Number of large read operations=0
         HDFS: Number of write operations=7
                                                                            Reduce output records=0
                                                                            Spilled Records=2
         HDFS: Number of bytes read erasure-coded=0
                                                                            Shuffled Maps =3
         S3: Number of bytes read=3629753164
         S3: Number of bytes written=0
                                                                            Failed Shuffles=0
                                                                            Merged Map outputs=3
         S3: Number of read operations=0
                                                                            GC time elapsed (ms)=650
         S3: Number of large read operations=0
                                                                            CPU time elapsed (ms)=650

CPU time spent (ms)=44560

Physical memory (bytes) snapshot=2159296512

Virtual memory (bytes) snapshot=25683906560

Total committed heap usage (bytes)=2070937600

Peak Map Physical memory (bytes)=498593792

Peak Map Virtual memory (bytes)=4397223936
         S3: Number of write operations=0
Shuffle Errors
            BAD ID=0
            CONNECTION=0
            IO_ERROR=0
                                                                            Peak Reduce Physical memory (bytes)=1107648512
Peak Reduce Virtual memory (bytes)=7113846784
            WRONG LENGTH=0
            WRONG_MAP=0
            WRONG_REDUCE=0
                                                  Job Counters
File Input Format Counters
                                                            Launched map tasks=1
            Bytes Read=257
                                                            Launched reduce tasks=3
File Output Format Counters
                                                            Data-local map tasks=1
                                                            Total time spent by all maps in occupied slots (ms)=386976
Total time spent by all reduces in occupied slots (ms)=21051648
Total time spent by all map tasks (ms)=4031
            Bytes Written=0
                                                            Total time spent by all reduce tasks (ms)=109644
                                                            Total vcore-milliseconds taken by all map tasks=4031
Total vcore-milliseconds taken by all reduce tasks=109644
                                                            Total megabyte-milliseconds taken by all map tasks=12383232
                                                            Total megabyte-milliseconds taken by all reduce tasks=673652736
```

Checking the resulting file on HDFS:

```
[hadoop@ip-172-31-23-236 ~]$ hdfs dfs -ls hdfs:///bigdata/amazon_reviews_multilingual_US_v1_00.tsv
-rw-r--r-- 1 hadoop hdfsadmingroup 3629753164 2022-12-16 21:20 hdfs:///bigdata/amazon reviews_multilingual_US_v1_00.tsv
```

\$ hdfs fsck hdfs:///bigdata/amazon_reviews_multilingual_US_v1_00.tsv -files -blocks - locations

```
Status: HEALTHY
Number of data-nodes: 2
                                1
Number of racks:
Total dirs:
Total symlinks:
Replicated Blocks:
Total size:
              3629753164 B
Total files:
Total blocks (validated):
                                28 (avg. block size 129634041 B)
Minimally replicated blocks:
                                28 (100.0 %)
Over-replicated blocks:
                                0 (0.0 %)
                                0 (0.0 %)
Under-replicated blocks:
Mis-replicated blocks:
                                0 (0.0 %)
Default replication factor:
Average block replication:
                               1.0
Missing blocks:
Corrupt blocks:
Missing replicas:
                                0 (0.0 %)
Erasure Coded Block Groups:
Total size:
               0 B
Total files:
               0
Total block groups (validated):
                                        0
Minimally erasure-coded block groups:
                                        0
                                        0
Over-erasure-coded block groups:
Under-erasure-coded block groups:
Unsatisfactory placement block groups: 0
Average block group size:
                                0.0
Missing block groups:
                                0
Corrupt block groups:
                                0
Missing internal blocks:
                                0
```

```
>>> from pyspark.sql.functions import col, isnan, when, count, udf
>>> sc.setLogLevel("ERROR")
>>> bucket= 'my-data-bucket-pc/'
>>> filename='amazon_reviews_multilingual_US_v1_00.tsv'
>>> file path = 's3a://' + bucket + filename
```

```
>>> print(file path)
s3a://my-data-bucket-pc/amazon reviews multilingual US v1 00.tsv
>>> sdf = spark.read.csv(file path, sep='\t', header=True, inferSchema=True)
>>> sdf.printSchema()
root
  - marketplace: string (nullable = true)
    customer id: integer (nullable = true)
    review_id: string (nullable = true)
     product id: string (nullable = true)
     product parent: integer (nullable = true)
     product_title: string (nullable = true)
     product_category: string (nullable = true)
     star rating: integer (nullable = true)
    helpful votes: integer (nullable = true)
     total_votes: integer (nullable = true)
     vine: string (nullable = true)
     verified_purchase: string (nullable = true)
    review_headline: string (nullable = true)
  -- review_body: string (nullable = true)
    review_date: timestamp (nullable = true)
```

Below are some summaries of the data after dropping some of the records where the column are empty(null)

```
>>> sdf.summary().show()
summary | marketplace |
                                customer id
                                                   review id
                                                                         product id
                                                                                            product parent
                                                                                                                  product
ting
           helpful_votes
                                   total_votes
                                                   vine | verified_purchase |
                                                                                  review_headline
                                                                                                             review_body
                                     6930482
  count
             6930482
                                                     6930482
                                                                            6930482
                                                                                                    6930482
30482
                  6930482
                                       6930482 | 6930482 |
                                                                    6930482
                                                                                           6930471
                                                                                                                  6930482
                 null | 2.917439358716003E7 |
                                                         null|1.0721688256266124E9| 4.932823608636339E8| 376979.1262
41819|2.0379774451473938|3.2418290675886614|
| stddev| null|1.5654661763964374E7|
                                                                                                                 Infinity
                                                   null|
                                                                       null
                                                                                               NaN
                                                         null|1.4497338047699845E9|2.8615311643958396E8|1902259.3995
48847 | 31.779313582603148 | 36.26553224008541 |
                                                   nul1
                                                                       null|
                                                                                               NaN
                                                                                                                      NaN
     min|
                   US
                                       10001|R10000WAG70YS2|
                                                                         0001046314|
                                                                                                     225472
                        0 |
                                             0 |
                                                      и|
                                                                                                  \bPlease make oth...
     25%
                 null
                                    15005223
                                                         null|
                                                                       3.45803485E8
                                                                                                 249510863
                                                   null
                                                                       null|
                                                                                               4.5
                                                                                                                      1.0
                        0 |
     50%
                 null|
                                    28788910
                                                         null|
                                                                       6.34049011E8
                                                                                                 496520461
    5|
                        0 |
                                             0 |
                                                   null|
                                                                       null|
                                                                                              11.0
                                                                                                                      1.0
     75%|
                 null|
                                    44131558
                                                         null|
                                                                      1.419553429E9
                                                                                                 744807844
                                                   null|
                         1 |
                                                                       null|
                                                                                             802.0
                   US
                                    53096589
                                              RZZZZKQ8A3VTL
                                                                         B06XXNYHRQ
                                                                                                 999988076
                                                                                                                 風立ちぬ
                    27550
                                         28727
                                                      Υ|
```

As we see that, there are some emotes/characters we don't need. Therefore to ignore these character (ascii characters), I am using the user-defined function(udf) below.