

Configuration Manual

MSc Research Project
Data Analytics

Chethan Marigowda
StudentID:x21145377

School of Computing
National College of Ireland

Supervisor: Dr. Paul Stynes

**National College of Ireland
Project Submission Sheet
School of Computing**



Student Name:	Chethan Marigowda
Student ID:	x21145377
Programme:	Data Analytics
Year:	2022
Module:	MSc Research Project
Supervisor:	Dr. Paul Stynes
Submission Due Date:	15/12/2022
Project Title:	Configuration Manual
Word Count:	1309
Page Count:	15

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	Chethan Marigowda
Date:	15th December 2022

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project , both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Chethan Marigowda
x21145377

1 Introduction

This research experiment followed a specific implementation configuration, and the objective of this guide is to describe how that configuration was set up. This documentation provides comprehensive details of the software, hardware and library configurations utilized in the development of this project. In addition to this, it details the method of programming as well as the procedure that should be carried out in order to execute the code.

2 Local Computing Machine System Configuration

The Figure 1 shows system configuration used in this project.

Device specifications

Yoga 7 14ITL5

Device name	DESKTOP-THA8E1A
Processor	11th Gen Intel(R) Core(TM) i7-1165G7 @ 2.80GHz 2.80 GHz
Installed RAM	16.0 GB (15.8 GB usable)
Device ID	C65A0108-CB2E-40E6-BAF0-71B54D7597E1
Product ID	00327-35938-59332-AAOEM
System type	64-bit operating system, x64-based processor

Fig. 1. System Configuration

3 Software, Packages and Library.

Hybrid Machine Learning Framework to Recommend E-Commerce Products is implemented in open-source freeware python programming language on Google Colaboratory notebook. Following Fig 1 in the fig details important packages/libraries used for this project. In the first step, all the necessary libraries and modules are installed and imported in the notebook. Some of the major libraries that were imported include Keras, Pandas, NumPy among, genism and others. With help of these libraries, necessary functions to carry out data manipulation and machine learning related computations can be carried out

```
[ ] #import libraries
import json
import numpy as np
from collections import defaultdict
import warnings
import gensim
import pandas as pd
from gensim.models.doc2vec import Doc2Vec
import numpy as np
import keras
import math
import gzip
from keras import backend as K
from keras.models import Model
from keras.layers import Embedding, Input, Dense, Flatten, Concatenate, Multiply, Lambda, Reshape
from keras.layers.core import Dropout
import scipy.sparse as sp
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error, mean_absolute_error
from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint
import math as mt
```

Fig. 2. Libraries used

4 Notebook Execution Guidance.

The research was conducted using free software provided by Google called Colaboratory notebook, which does not require any static setup settings to be entered. The modest addition of this project is extensive code that enables it to execute by sourcing files directly from the source destination. This is the minor contribution of this research and useful in generating required product categories files.

The standard steps to be followed to execute the code from the fig 3 to 4 for the end-to-end execution.

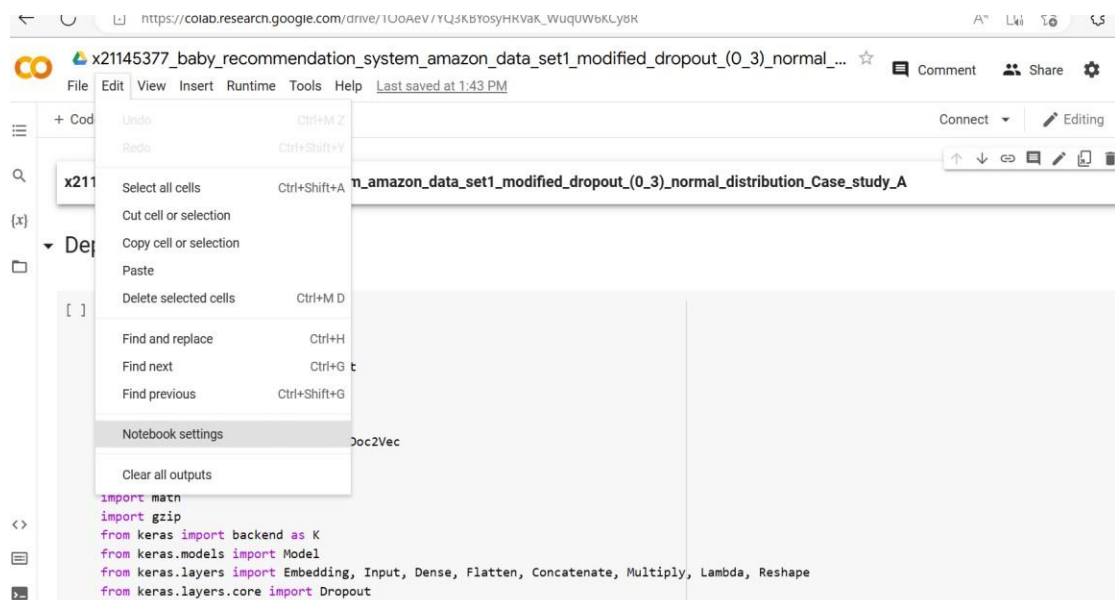


Fig. 3. Execution step1

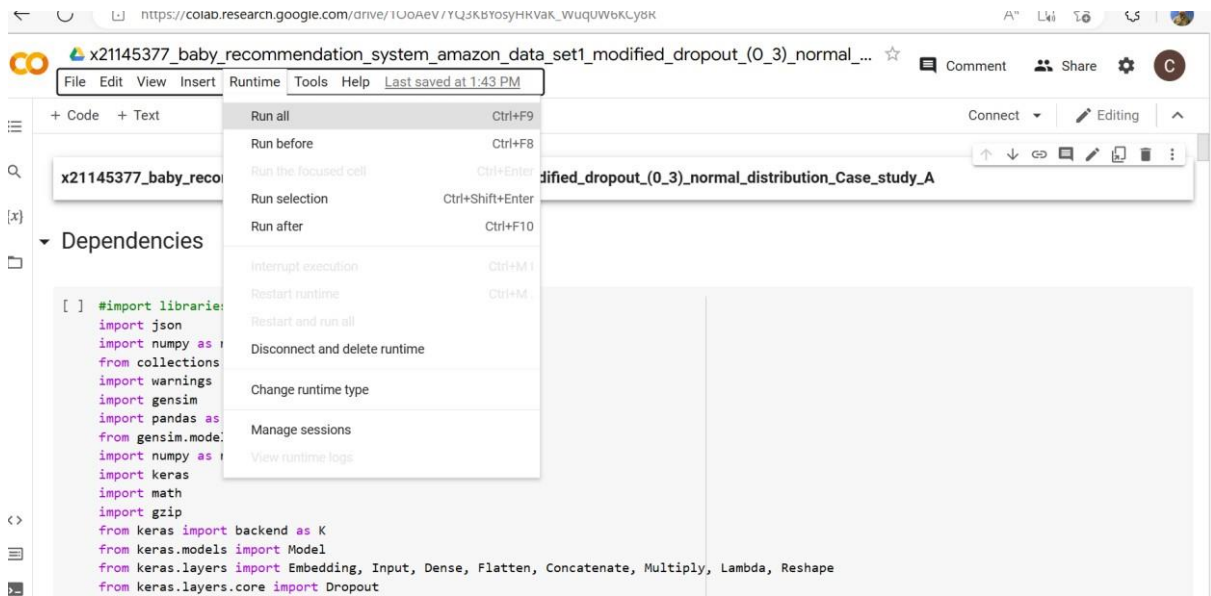


Fig. 4. Execution step2

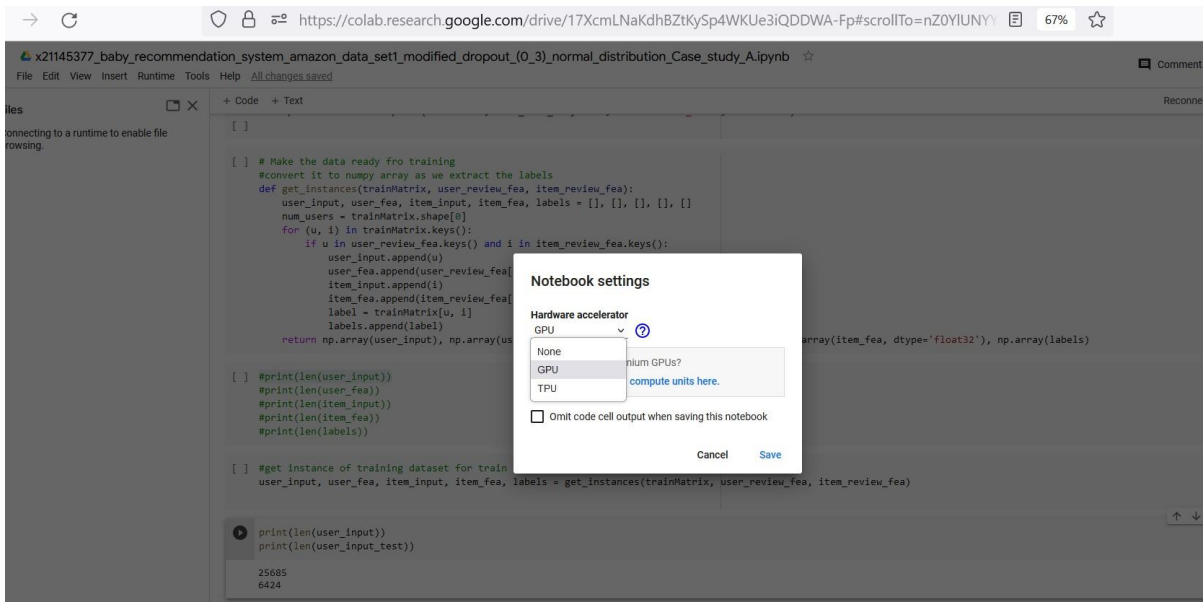


Fig. 5. Execution step3

5 Data Acquisition

The dataset collected from public open-source web server, after decompression contains 915445 consumer reviews and metadata of various baby products up to 71316 in semi structured json format for Data set 1[1]. Data set 2[2] contains 836005 consumer reviews and metadata of various digital music products up to 279898 in semi structured json format. Wget is a networking command-line tool that lets you download files and interact with REST APIs.

▼ Display the dataset

```
[ ] ## Convert data from gz to DataFrame
def parse(path):
    g = gzip.open(path, 'rb')
    for l in g:
        yield eval(l)

def getDF(path):
    i = 0
    df = {}
    for d in parse(path):
        df[i] = d
        i += 1
    return pd.DataFrame.from_dict(df, orient='index')
```

Fig. 6. Data acquisition illustration 1

▼ Download The Review and Meta Data

```
[ ] # Download Review Data
!wget http://snap.stanford.edu/data/amazon/productGraph/categoryFiles/reviews_Baby.json.gz

--2022-12-14 18:05:23-- http://snap.stanford.edu/data/amazon/productGraph/categoryFiles/reviews_Baby.json.gz
Resolving snap.stanford.edu (snap.stanford.edu)... 171.64.75.80
Connecting to snap.stanford.edu (snap.stanford.edu)|171.64.75.80|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 191601185 (183M) [application/x-gzip]
Saving to: 'reviews_Baby.json.gz.3'

reviews_Baby.json.gz 100%[=====] 182.72M  4.60MB/s   in 34s

2022-12-14 18:05:57 (5.42 MB/s) - 'reviews_Baby.json.gz.3' saved [191601185/191601185]
```

Fig. 7. Data acquisition illustration 2

```
[ ] # Download Meta Data
!wget http://snap.stanford.edu/data/amazon/productGraph/categoryFiles/meta_Baby.json.gz

--2022-12-14 18:05:57-- http://snap.stanford.edu/data/amazon/productGraph/categoryFiles/meta_Baby.json.gz
Resolving snap.stanford.edu (snap.stanford.edu)... 171.64.75.80
Connecting to snap.stanford.edu (snap.stanford.edu)|171.64.75.80|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 31328004 (30M) [application/x-gzip]
Saving to: 'meta_Baby.json.gz.3'

meta_Baby.json.gz.3 100%[=====] 29.88M  2.11MB/s   in 13s

2022-12-14 18:06:10 (2.30 MB/s) - 'meta_Baby.json.gz.3' saved [31328004/31328004]
```

Fig. 8. Data acquisition illustration 3

The screenshot shows a Jupyter Notebook with the following code and output:

```
[ ] #display the meta data as dataframe
metadata = getDF('./meta_Baby.json.gz')
metadata.head()
```

	asin	categories	description	title	price	imUrl	brand	related	salesRank
0	0188399313	[[Baby]]	Wee-Go Glass baby bottles by LifeFactory (Baby...	Lifefactory 4oz BPA Free Glass Baby Bottles - ...	69.99	http://ecx.images-amazon.com/images/I/41Swthpd...	Lifefactory	{'also_bought': ['B002SG7K7A', 'B003CJSXW8', '...	NaN
1	0188399518	[[Baby]]	The Planet Wise Flannel Wipes are 10 super sof...	Planetwise Flannel Wipes	15.95	http://ecx.images-amazon.com/images/I/41otjnA4...	Planet Wise	{'also_bought': ['B00G96N3YY', 'B003XSEV2O', '...	NaN
2	0188399399	[[Baby]]	The Planet Wise Wipe PouchTM features our pate...	Planetwise Wipe Pouch	10.95	http://ecx.images-amazon.com/images/I/61x8h9u6...	NaN	{'also_bought': ['B005WWI0DA', 'B005WWWIMGA', '...	NaN
3	0316967297	[[Baby]]	Hand crafted set includes 1 full quilt (76x86 ...	Annas Dream Full Quilt with 2 Shams	109.95	http://ecx.images-amazon.com/images/I/51%2BZ1%...	NaN	{'also_viewed': ['B009LTER3W', 'B00575T15Q', '...	NaN
4	0615447279	[[Baby]]	Thumbuddy To Love- The Binky Fairy helps child...	Stop Pacifier Sucking without tears with Thumb...	16.95	http://ecx.images-amazon.com/images/I/51RKKENI...		{'also_bought': ['0979670004', '1601310234', '...	NaN

Fig. 9. Data acquisition illustration 4

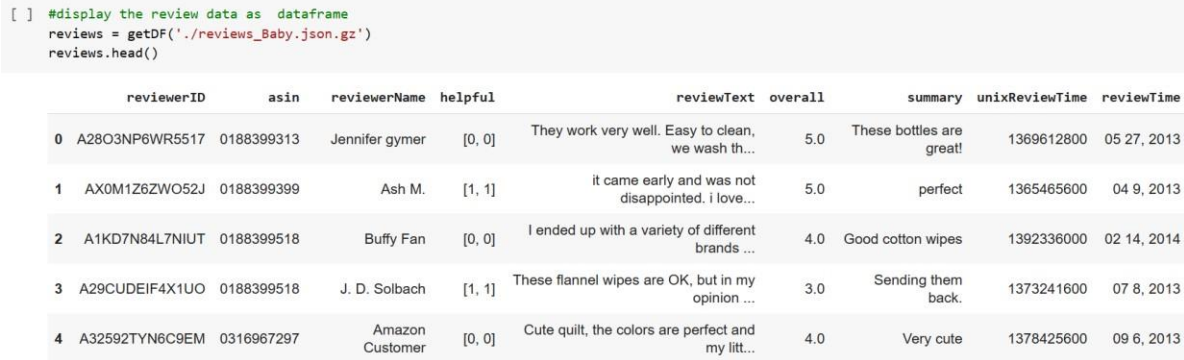


Fig. 10. Data acquisition illustration 5

6 Symbols and Notations

The important symbols used in the coding are shown in Table 1. For the convenience, consumer and product are used in programming as user and item for easier notation respectively and same is commented in code as shown in fig 11.

Table 1: Parameter conventions used with their function

Symbols	Meaning
vector_size	vector size of review text and item description
epoch_num	number of epochs to train for paragraph vectors.
num_ui_link	Min. no. of reviews required for a reviewer to be filtered
num_iu_link	Min. no. of reviews required for a product
max_num_item	max number of item
max_num_user	max number of users
split_ratio	80-20 split of train and test data
dataname	Name of the product eg: baby, digital music

```

#Parameters
dataName = "Baby"# name of your product
vector_size = 100 # vector size of reviewtext and item description
epoch_num = 100 # number of vector
reviews_data = [] # review list to insert review data inside
meta_data = [] # meta list to insert meta data inside
num_ui_link = 20 # the number of each user link to items
num_iu_link = 0 # the number of each item link to user

# user -> item -> for text
ui_dict = defaultdict(list)
#item -> user -> for text
iu_dict = defaultdict(list)
#rating training data list
reviews_train_data = []
#rating testing data list
reviews_test_data = []
# max number of item
max_num_item = 1
# max number of users
max_num_user = 1
#item -> user -> for rating
iu_dict2 = defaultdict(list)

```

Fig 11: Example of symbols

2 Implementation framework Specs

The input data pre-processing algorithm is described in Algorithm 1. The algorithm is capable of being trained in the same manner as other keras-based deep learning networks. On top of the recommendation model, this research add a regression layer, whose output serves as our loss function. This permits us to accomplish our rating prediction target. Figure 4 depicts the training procedure using the Python keras package. The dense network of consumer sentiment analysis model and product content model is shown by Algorithm

2. The Rating Predictor is shown by Algorithm 3 in Fig. The Model Summary is depicted in the fig 13

Algorithm 1 Preprocessing of Amazon Review Data	
Input:	Reviewers data set
Output:	Preprocessed data for data transformation
1:	Preprocess()
2:	Initialize hyperparameters:
3:	vector_size ← 100
4:	epoch_num ← 100
5:	num_ui_link ← 20
6:	num_iu_link ← 0
7:	max_num_item ← 1
8:	max_num_user ← 1
9:	split_ratio ← 4
10:	Batch_size ← 64
11:	for each data_row, i ∈ review_data do
12:	user_id ← i["reviewerID"]
13:	item_id ← i["asin"]
14:	ui_dict[user_id] ← append(item_id)
15:	iu_dict[item_id] ← append(user_id)
16:	If ((num_ui_link < 20) (num_iu_link < 1))
17:	delete i
18:	end for
19:	Covert review_data.json to review_data.txt
20:	for each data_row, i ∈ meta_data do
21:	item_id ← i["asin"]
22:	if (iu_num != len(iu_dict[item_id]))
23:	if (iu_num < num_iu_link)
24:	delete i


```

25: end for
26: Covert meta_data.json to meta_data.txt
27: Gensim_doc2vec()
28: Use pretrained model from Gensim with parameters:
29:     min_count←1
30:     window←3
31:     vector_size←vector_size
32:     sample←10-3
33:     negative←5
34:     workers←4
35: Convert data.txt to vector format
36: Extract_rating()
37: for each data_row, i ∈ review_data do
38:     user_id ← i["reviewerID"]
39:     item_id ← i["asin"]
40:     if (user_id ! ∈ ui_dict)
41:         ui_dict[user_id] ← user_id
42:         if (user_id > max_num_user)
43:             max_num_user ← user_id
44:     ui_dict[user_id] ← item_id
46:     if (item_id ! ∈ iu_dict2)
47:         iu_dict[item_id] ← item_id
48:         if (item_id > max_num_item)
49:             max_num_item ← item_id
50:     iu_dict[item_id] ← user_id
51:
52: Test_Train_Split()
53:     num ← 0
54: for each data_row, i ∈ reviews_data do
55:     if (num % (split_ratio + 1) < split_ratio)
56:         reviews_train_data.append(i)
57:     else
58:         num ← num + 1
59:         reviews_test_data.append(i)
60: num ← num+1

```

Algorithm 2 User Sentiment and Item Content Dense Network

Input: User_Data, User_Latent, Item_vector, item_latent

Output: Users_Sentiment, Item_Content

```

1: User_Sentiment(User_data, User_Latent)
2: Initialize hyperparameters:
3:     input ← User_data
4:     activation() ← relu()
5:     Lamda() ← softmax()
6: User_Sentiment ← Multiply(User_data, Lamda(activation(User_Latent), input))
7: Item_Content(Item_vector, item_Latent)
8: Initialize hyperparameters:
9:     input ← Item_vector
10:    activation ← relu
11:    Lamda ← softmax()
12: Item_Content ← Multiply(Item_vector, Lamda(activation(item_latent), input))

```

Algorithm 3 Recommender Model

Input: User_Sentiment, Item_Content

Output: Prediction, Results

```

1: Recommender(User_Sentiment, Item_Content)
2: Initialize hyperparameters:
3:     optimizer() ← adam
4:     loss() ← MAE()
5:     metrics() ← RMSE
6: model() ← Compile(optimizer(), loss(), metrics())
7: Prediction, Results ← Compute(model())

```

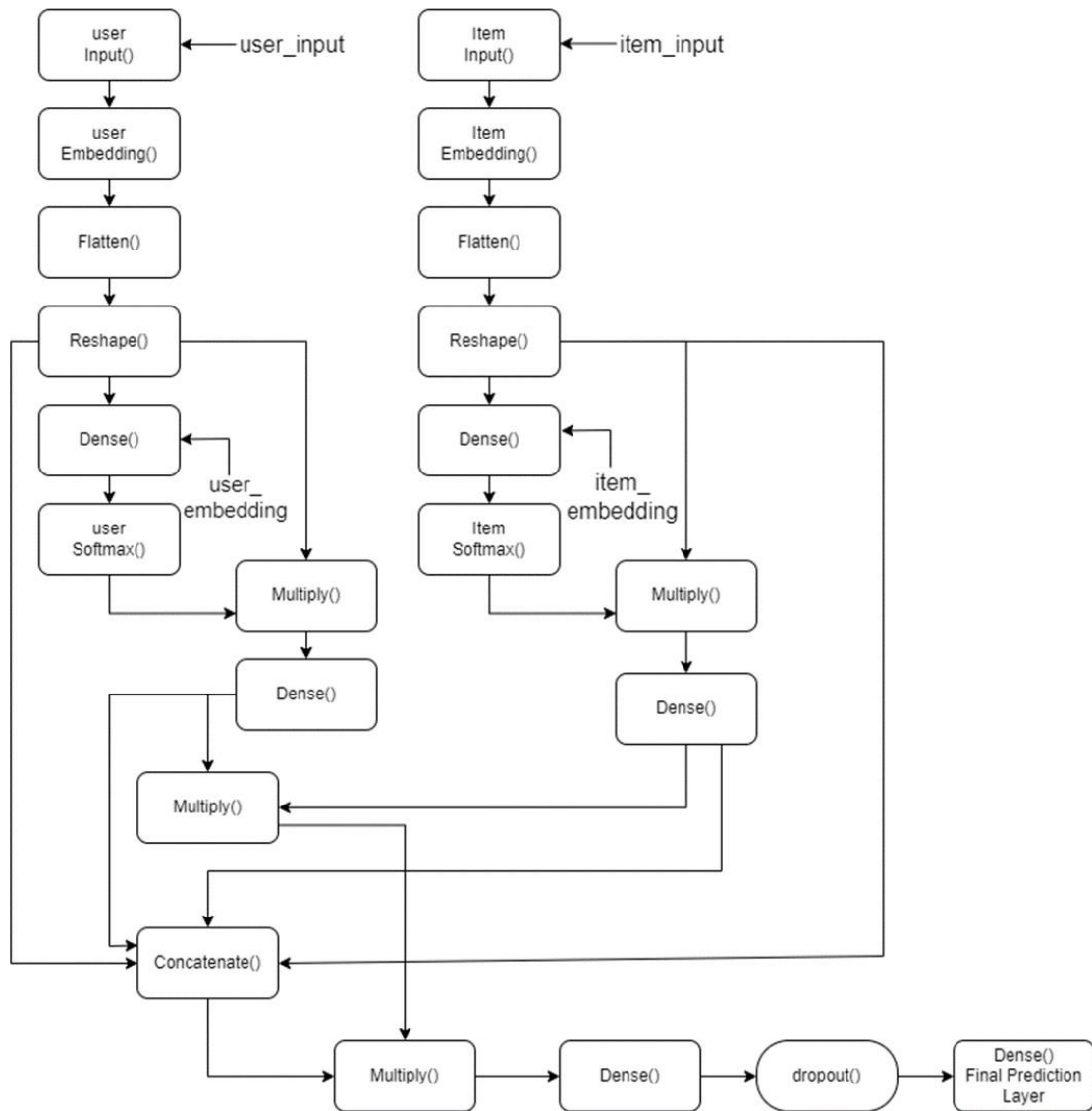


Fig. 12. Keras Training Flowchart

```

x2145377_baby_recommendation_system_amazon_data_set1_modified_dropout_0_3_normal_distribution_Case_study_A.ipynb
File Edit View Insert Runtime Tools Help All changes saved
Code + Text

model.summary()

user_sentiment (InputLayer) [(None, 100)] 0 []
item_embedding (Embedding) (None, 1, 100) 10049000 ['item_input[0][0]']
item_content (InputLayer) [(None, 100)] 0 []
flatten (flatten) (None, 100) 0 ['user_embedding[0][0]']
user_attention_layer (Dense) (None, 100) 10100 ['user_sentiment[0][0]']
flatten_1 (flatten) (None, 100) 0 ['item_embedding[0][0]']
item_attention_layer (Dense) (None, 100) 10100 ['item_content[0][0]']
reshape (Reshape) (None, 100) 0 ['flatten[0][0]']
user_sentiment_softmax (Lambda (None, 100)) 0 ['user_attention_layer[0][0]']
reshape_1 (Reshape) (None, 100) 0 ['flatten_1[0][0]']
item_content_softmax (Lambda) (None, 100) 0 ['item_attention_layer[0][0]']
multiply (Multiply) (None, 100) 0 ['reshape[0][0]', 'user_sentiment_softmax[0][0]']
multiply_1 (Multiply) (None, 100) 0 ['reshape_1[0][0]', 'item_content_softmax[0][0]']
dense (Dense) (None, 100) 10100 ['multiply[0][0]']
dense_1 (Dense) (None, 100) 10100 ['multiply_1[0][0]']
concatenate (Concatenate) (None, 400) 0 ['dense[0][0]', 'item_content[0][0]', 'dense_1[0][0]', 'dense_1[0][0]']
multiply_2 (Multiply) (None, 100) 0 ['dense[0][0]', 'dense_1[0][0]']
dense_2 (Dense) (None, 100) 40100 ['concatenate[0][0]']
multiply_3 (Multiply) (None, 100) 0 ['multiply_2[0][0]', 'dense_2[0][0]']
dense_3 (Dense) (None, 100) 10100 ['multiply_3[0][0]']
dropout (Dropout) (None, 100) 0 ['dense_3[0][0]']
prediction (Dense) (None, 1) 101 ['dropout[0][0]']

=====
Total params: 1,246,701
Trainable params: 1,246,701
Non-trainable params: 0

```

Fig. 13. Model Summary

3 Evaluation and Results

In all of our experiments and case studies, this research evaluated the performance of our proposed sentiment model by using the Mean Squared Error (MSE) and Root Mean Squared Error (RMSE). This is because the Mean Squared Error is the only metric that has been utilized for assessment in the majority of the relevant works [3].

Let's say that the total number of datapoints in the test set is denoted by the letter n. The MSE may be determined using the formulas below.

$$MAE = \frac{\sum_{k=0}^{n-1} |Y(k) - y(k)|}{n}$$

In a more similar manner, RMSE () is computed by:

$$RMSE = \sqrt{\frac{\sum_{k=0}^{n-1} (Y(k) - y(k))^2}{n}}$$

Where Y(k) and y(k) represent the actual and predicted value respectively

6.1 State of Art

State of Art-evaluation Metrics Capture are illustrated from the Fig14 to 15.

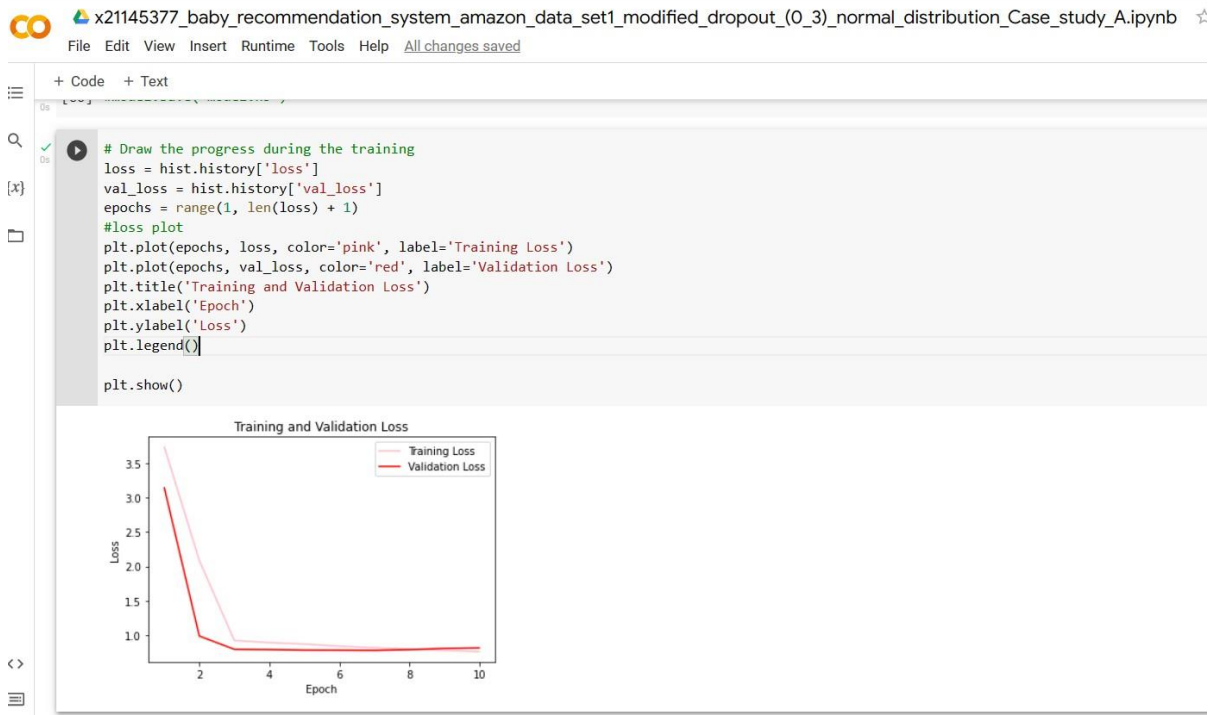


Fig. 14. Training and validation loss

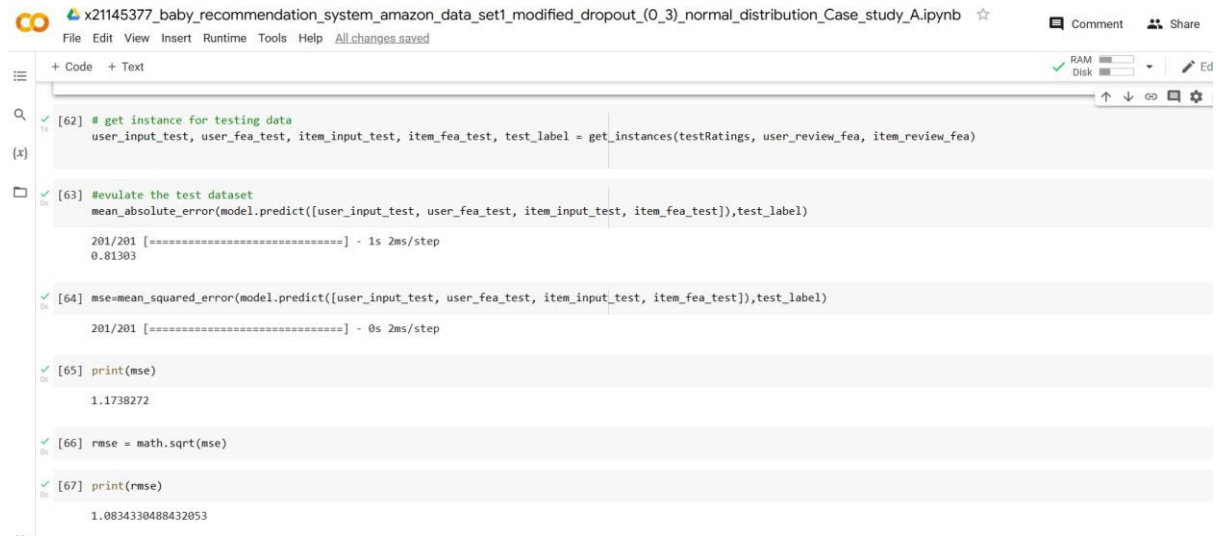


Fig. 15. Baby evaluation accuracy Metrics Capture

6.2 Replicating state of Art with different product category data set- :

The accuracy evaluation Metrics Capture is illustrated from the Fig16

```
x21145377_digital_music_recommendation_system_amazon_data_set_2_Case_study_B.ipynb ☆
File Edit View Insert Runtime Tools Help Last saved at 3:18 PM

+ Code + Text Conn

[ ] # get instance for testing data
user_input_test, user_fea_test, item_input_test, item_fea_test, test_label = get_instances(testRatings, user_review_fea, item_review_fea)

[ ] #evaluate the test dataset
mean_absolute_error(model.predict([user_input_test, user_fea_test, item_input_test, item_fea_test]),test_label)
475/475 [=====] - 1s 2ms/step
0.55482465

[ ] mean_squared_error(model.predict([user_input_test, user_fea_test, item_input_test, item_fea_test]),test_label)
475/475 [=====] - 1s 2ms/step
0.65298164

[ ] prediction = model.predict([user_input_test, user_fea_test, item_input_test, item_fea_test])
475/475 [=====] - 1s 2ms/step
```

Fig. 16. Digital music dataset evaluation accuracy Metrics Capture

6.3 Sentiment analysis

Sentiment analysis accuracy evaluation Metrics Capture is illustrated from the Fig 17 to 19

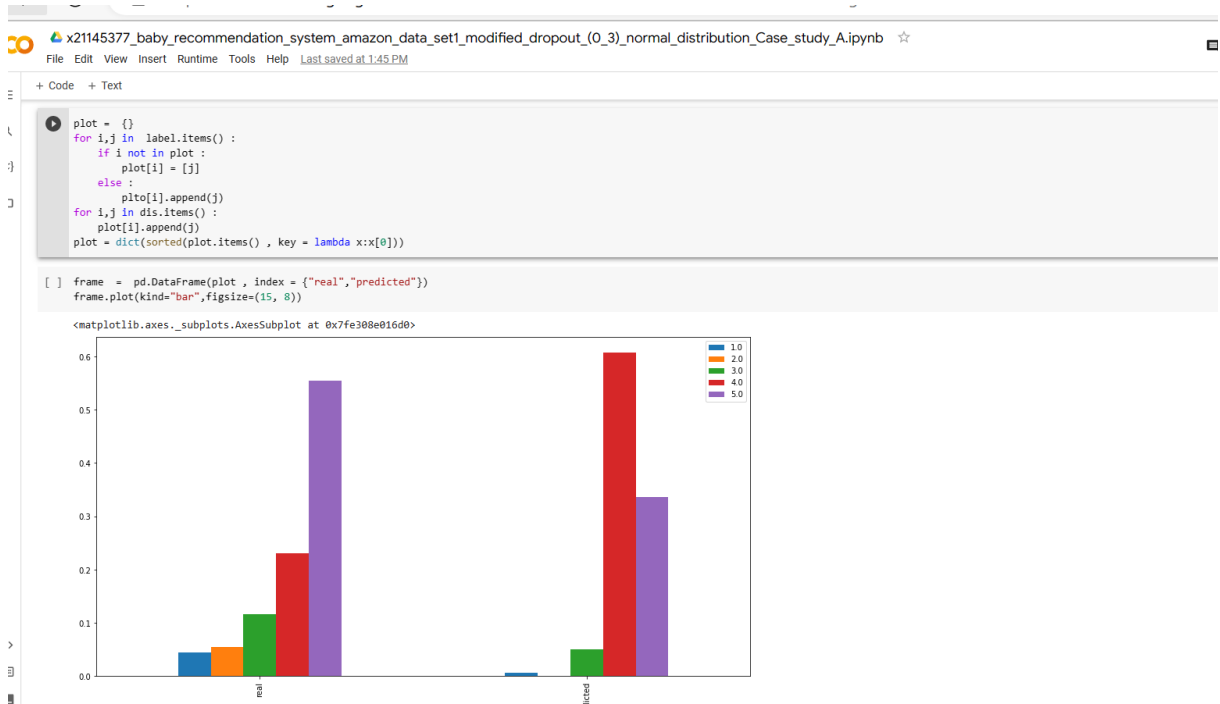


Fig17. Baby data set ratings distribution Metrics Capture

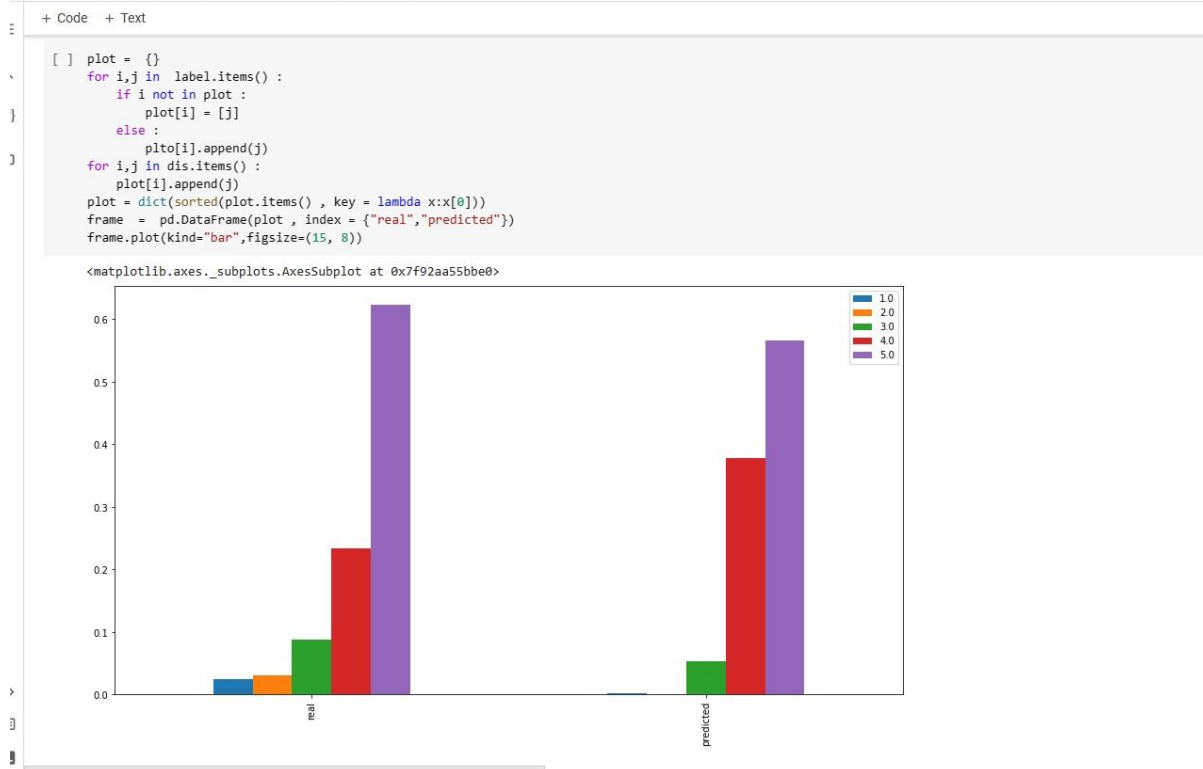


Fig. 18. Digital Music data set ratings distribution Metrics Capture

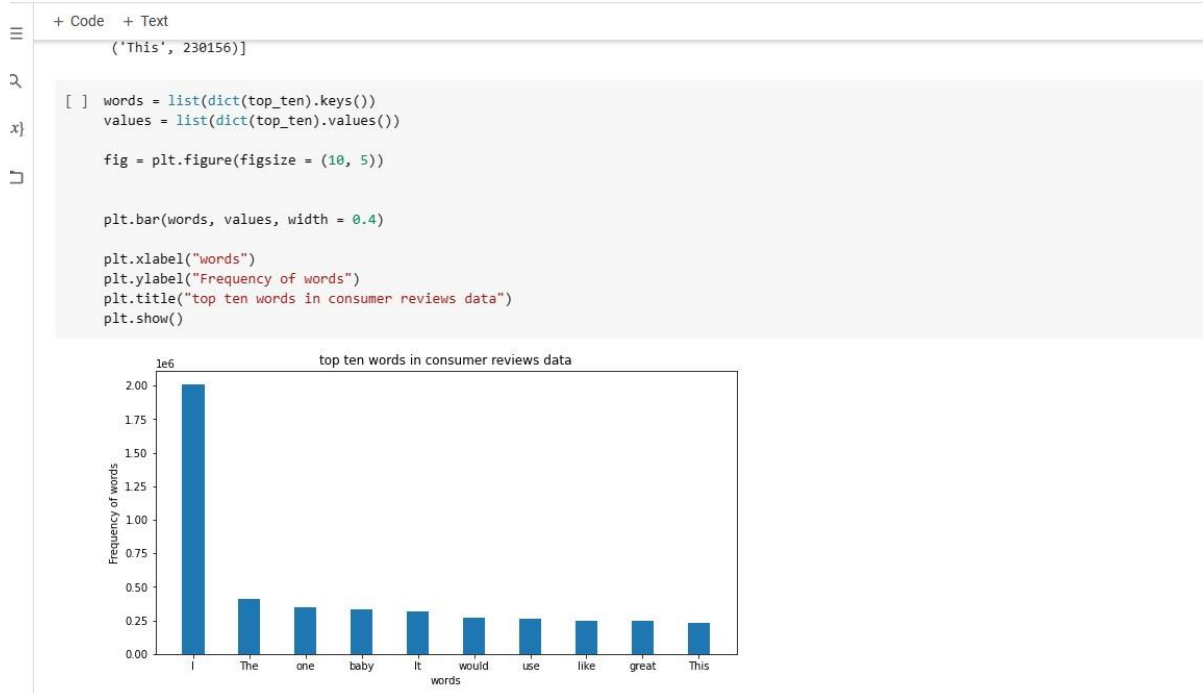
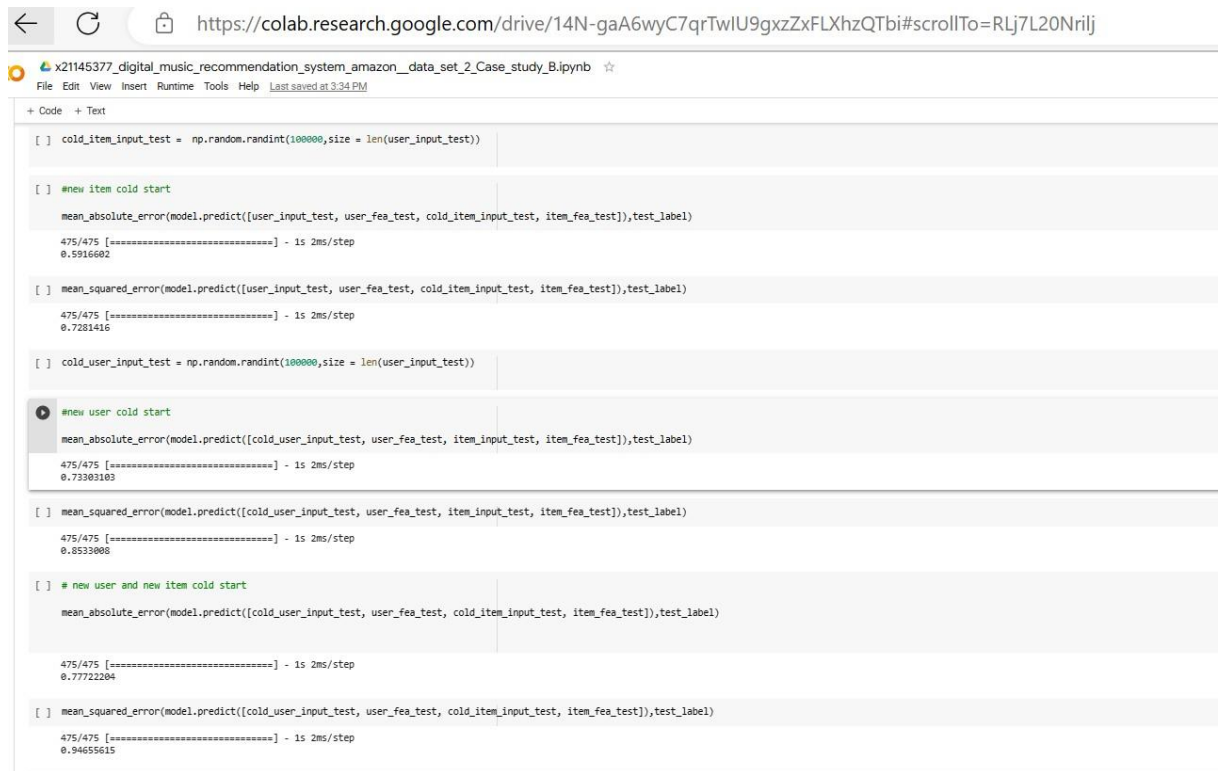


Fig. 19. Words distribution illustration

6.4 Cold Start Analysis

The Cold start accuracy Evaluation Metrics is Captured in Fig 19.



```
← ↻ 🔒 https://colab.research.google.com/drive/14N-gaA6wyC7qrTwiU9gxZxFLXhzQTbi#scrollTo=RLj7L20Nrij
x21145377_digital_music_recommendation_system_amazon_data_set_2_Case_study_B.ipynb ☆
File Edit View Insert Runtime Tools Help Last saved at 3:34 PM
+ Code + Text
[ ] cold_item_input_test = np.random.randint(100000,size = len(user_input_test))

[ ] #new item cold start
mean_absolute_error(model.predict([user_input_test, user_fea_test, cold_item_input_test, item_fea_test]),test_label)
475/475 [=====] - 1s 2ms/step
0.5916602

[ ] mean_squared_error(model.predict([user_input_test, user_fea_test, cold_item_input_test, item_fea_test]),test_label)
475/475 [=====] - 1s 2ms/step
0.7281416

[ ] cold_user_input_test = np.random.randint(100000,size = len(user_input_test))

[ ] #new user cold start
mean_absolute_error(model.predict([cold_user_input_test, user_fea_test, item_input_test, item_fea_test]),test_label)
475/475 [=====] - 1s 2ms/step
0.73303103

[ ] mean_squared_error(model.predict([cold_user_input_test, user_fea_test, item_input_test, item_fea_test]),test_label)
475/475 [=====] - 1s 2ms/step
0.8533008

[ ] # new user and new item cold start
mean_absolute_error(model.predict([cold_user_input_test, user_fea_test, cold_item_input_test, item_fea_test]),test_label)
475/475 [=====] - 1s 2ms/step
0.77722204

[ ] mean_squared_error(model.predict([cold_user_input_test, user_fea_test, cold_item_input_test, item_fea_test]),test_label)
475/475 [=====] - 1s 2ms/step
0.94655615
```

Fig. 20. Cold start Evaluation Metrics Capture

References

- [1] Dataset1 (Baby Product Dataset) – Source: <http://jmcauley.ucsd.edu/data/amazon/links.html>, last accessed 2022/12/13.
- [2] Dataset2 (Digital Music Product Dataset) – Source: <http://jmcauley.ucsd.edu/data/amazon/links.html>, last accessed 2020/12/13.
- [3] Dezfouli, P. A. B., Momtazi, S. and Dehghan, M. [2020]. Deep neural review text interaction for recommendation systems, Appl. Soft Comput. 100: 106985.