

Configuration Manual

MSc Research Project MSc in Cloud Computing

Vipin Yadav Student ID:19211791

School of Computing National College of Ireland

Supervisor : Prof. Divyaa Manimaran Elango



11 Club
Vipin Yadav
19211791
MSc in Cloud Computing
2021
MSc Research Project
Divyaa Manimaran Elango www.ncirl.ie
16-12-2021
Configuration Manual
767
11

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.

<u>ALL</u> 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:	Vipin Yadav
Date:	31-01-2022

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	
Attach a Moodle submission receipt of the online project submission, to each project (including multiple copies).	
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.	

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

Vipin Yadav x19211791

1 Introduction

This configuration manual will document the implementation of the current research project. It will explain how to run the implemented scripts. This documentation will help to execute the code without any problem. There will be recommendations of hardware and software versions required to run this. The project's results can be replicated if these methods are followed properly.

2 System Configuration

2.1 Hardware and Software Configuration

For this project implementation, we have used Google Colab IDE and Microsoft Office tools are being used for formatting and writing

Hardware (Configuration	Software Configuration			
Operating System	Microsoft Windows 10	Microsoft Office suite			
RAM	8 GB	IDE - Google collab			
Hard Disk	500 GB	Python 3.6.9			
Processor	Intel(R) Core(TM) i7- 10510U CPU @ 1.80GHz, 2304 Mhz, 4 Core(s), 8 Logical Processor(s)				

Figure 1: Software and Hardware configuration

2.2 Project Development

Amazon's product reviews dataset was created by Jianmo Ni and is publically presented on the below link:

http://deepyeti.ucsd.edu/jianmo/amazon/index.html

This Dataset of Electronic product needs to be uploaded on Google drive

Below figure 2 shows to mount the google collab with the drive so that the dataset can be easily available.

120	C 🔒 colab.r	research.go	ogle.com/drive/18)	KK5nH4	QQeDANKJvki_gJ1222Mlsx7qo?auth	user=1#s	scrollTo=gf5sAJg8LZ07			12	Q	Ċ	☆	*	=1	Y)
G Fi	inal.ipynb 💠 Edit View Insert Runtii	me Tools He	lp All changes saved									Q c	omment	**	Share 🗘	2
+ Code	+ Text												Reconnec	t •	🖍 Editio	g
[]	from google.colab impo drive.mount(' <u>/content/</u>	ont drive (drive')														
D	Drive already mounted	at /content/	drive; to attempt to for	cibly rem	mount, call drive.mount("/content/drive", fo	rce_remount	-True).									
	import pandas as pd import torch import numpy as np import matplotlib.pypl import seaborn as sns from wordcloud import from sklearn.feature.e from textlibel import 1	lot as plt WordCloud, S Extraction.te	TOPWORDS xt import CountVectorize	tr.												
			A Line Inc. Bulley Description		Continue to the second second sectors to the		-									
[];	review_df - pd.read_ review_df.head()	_json(' <u>/conte</u>	nt/drive/My Drive/Amazon	Product	Sentiment_IE/Electronics_5.json', orient-'re	cords', lin	es-True)	unixReviewTite	reviewTine							
	review_df = pd.read_ review_df.head() reviewerID 0 A094DHGC771SJ	json(' <u>/conte</u> asin 0528881469	nt/drive/My Drive/Amazon reviewerName amazdnu	helpful	Sentiment IE/Electronics 6. json', orient-'re reviewText We got this GPS for my husband who is an (OTR)	cords', lin overall 5	es-True) summary Gotta have GPSI	unixReviewTime 1370131200	reviewTime 06.2, 2013							
11 ;	review_df = pd.read_ review_df.head() reviewerID 0 A094DHGC771SJ 1 AM0214LNFCEI4	_json(' <u>/conte</u> asin 0528881469 0528881469	nt/drive/My Drive/Amazon reviewerName amazonu Amazon Customer	helpful [0,0] [12,15]	Sentiment_IE/Electronics_5.json',orient-'re reviewText We got this GPB for my husband who is an (OTR). Trm a professional OTR fuck driver, and I bou.	cords', lin overall 5 1	es=True) summary Gotta have GPSI Very Disappointed	unixReviewTime 1370131200 1290643200	reviewTine 06 2, 2013 11 25, 2010							
C1 ;	review_df - pd.read_ review_df.head() reviewerID 0 A094DHGC771SJ 1 AM0214LNFCEI4 2 A3N7T0DY83Y4IG	_json('/conte asin 0528881469 0528881469 0528881469	nt/drive/My Drive/Amazon reviewerName amazonu Amazon Customer C. A. Freeman	helpful [0, 0] [12, 15] [43, 45]	Sentiment_IF/Electronics_6.json*, erient=*rev reviewText We got this OPS for my husband who is an (OTR) Tra a professional OTR truck driver, and Tolo Weld, what can I say. We had this unfit in m	cords', lin overall 5 1 3	es-True) Sumary Gota have GPSI Very Disappointed 1st impression	unixReviewTime 1370131200 1290643200 1283990400	reviewTine 06.2, 2013 11.25, 2010 09.9, 2010							
[];	review_df - pd.read_ review_df.head() 0 A094DHGC77ISJ 1 AMO214LNFCEI4 2 A3N7T0DY83Y4IG 3 A1H8PY3QHIAQQA0	_json(' <u>/conte</u> asin 0528681469 0528681469 0528681469 0528681469	nt/drive/My Drive/Amazon reviewerNane amazonu Amazon Customer C. A. Freeman Dave M. Shaw "mack dave"	helpful [0, 0] [12, 15] [43, 45] [9, 10]	Sentiment_IF/Electronics_5.json*.orient='re reviewText We got this GPS for my husband who is an (OTR). I'm a professional OTR hurux drives, and 10ou. Welk, what can I say. The had this unit in m. Not going to write a long review, even thought.	cords', lin overall 5 1 3 2	es-True) Sumary Oota have GPSI Veny Disappointed 151 impression Great grafics, POCR GPS	unixReviewTime 1370131200 1290643200 1283990400 1290556800	reviewTine 06 2, 2013 11 25, 2010 09 9, 2010 11 24, 2010							
	review_df = pd.read_ reviewedf = pd.read_ reviewerID 0 A0940HGC71SJ 1 AM0214LINFCEH 2 A3N7T00Y83Y4IG 3 A1H8PY3GHMQGA0 4 A24EV6RXELQ263	jsen('/conte asin 0528681469 0528681469 0528681469 0528681469	nt/drive/My Drive/Amazon reviewerName amazonu Amazon Customer C. A. Freeman Dave M. Shaw "mack dave" Wayne Smith	helpful [0,0] [12,15] [43,45] [9,10] [0,0]	Sentiment_IF/Electronics_5_ison*, erlant* re reveloation We got this GPS for my husband who is an (OTN). I'm a professional GTR fluck drive; and I bou. Well, what can i say, be had the unit in , hot grang to ware a long reverse, we not thought. Due had mine for a year and here's what we go	cords', lin overall 5 1 3 2 1 h	es=Trut) Sommary Ootta have GPSI Very Disappointed 15t impression Great grafics, POOR GPS tepor issues, only excuses for support	unixReviewTime 1370131200 1290643200 1283990400 1290556800 1317254400	revies/Tine 06 2, 2013 11 25, 2010 09 9, 2010 11 24, 2010 09 29, 2011							
	review_df - pd.read review_df.head) reviewerID 0 A094DHGC715J 1 AN074LNFCEH 2 A3NT00753VIG 3 A1HBPY30H/AGA0 4 A34EVERXEL2253 # Charget the overall review_df("classe") - review_df("classe").	json('/conte asin 0528081469 0528081409 05280 050000000000	nt/drive/My Orive/Amazon reviewerName amacdnu Amazon Cushomer C. A. Freeman Dave M. Shaw "mack dave" Wayne Smith www.smith categorial from numeric coverpail of the second second "positive"_4.8:"positiv	helpful [0,0] [12,15] [43,45] [9,10] [0,0] (al.	Sentiment Tri/Electronics, 5, 5, 5, 5, 7, 4, 74 and * /re received and We got this GPS for my husband who is an (OTR). I'm a professional OTR link of where, and Ibou. West, what can I say. The had this will in m. Not going to white a long inview, even thought. The had mine for a year and here's what we go.	overall 5 1 2 1 h lace=True)	sex-True) Gota have GPBI Very Disquored Issi impression Great grades, Booch GPB bipor issues, only excuses for support	unixReviewTime 1370131200 1290643200 1283990400 1290056800 1317254400	reviesTine 06 2, 2013 11 25, 2010 09 9, 2010 11 24, 2010 09 29, 2011							
	review_df - pd.read, review_df.read() reviewrD 0 A040HGC7181 2 A002HNCEH 2 A3N7D0YS3Y4G 3 A148PY2GHAD2A5 4 A24EVERDEL0255 4 A24EVERDEL0255 review_df.rease1,	json('/conte asin 0528881469 0528881469 0528881469 0528881469 0528881469 0528881469 0528881469 0528881469 0528881469 category to = review_f('.etc) replace((5.0: mique()) agative', 'nep	nt/driw/hy Grive/Amazon reviewerKame amazon Armazon Coustoner C A Freeman Dave M Shaw 'mack daw' Wayne Semth Dave M Shaw 'mack daw' Wayne Semth Topolitier ', dripolitis' 'positive'', dripolitis' tral'], drype=object)	helpful [0, 0] [12, 15] [43, 45] [9, 10] [0, 0] (al.	Sentiment Tri/Electronics, 5, 5, 5, 5, 7, 4, 74 and * /e recallent We got this GPS for my husband who is an (OTR). I'm a professional OTR link of where and bou- West, what can I say. Ne had this will n in Not going to where a long network, when thought. Ne had mine for a year and here's what we go.	overall 6 1 3 2 1 h	sex-True) Gota have GPBI Very Disquored Issi impression Great grades, Booch GPB tepri issues, only excuses for support	un1skev1ewTine 1370131200 1290643200 1283990400 1290556800 1317254400	reviewTine 06 2, 2013 11 25, 2010 09 9, 2010 11 24, 2010 09 29, 2011							
	review_df - pd.read, review_df.read() review_df.read() evelwerID AOQ4HGG7T81 AOQ5HGG7T81 AOQ5HG7T82 AOD5HG7C84 AOA0000 AOA0000 AOA00000 AOA000000 AOA0000000 AOA00000000	jion('/conte asin 0528881469 0528881469 0528881469 0528881469 0528881469 0528881469 0528881469 0528881469 0528881469 0528881469 scategory to replace((5.0) influe() influe() influe() influe() influe() influe()	et/driv/hy Drive/Amazon residentman Amazon Customer C A Freeman Dave M Sham Yang Wayne Smth comportional from numeric reseal[1] "positive",4.0:"positiv stral"], dtyp=object) lamore datast c[rescalingtam]	helpful [0.0] [12.15] [43.45] [9.10] [0.0] (al. (e ⁺ ,3.0: ⁺)	Sentiment Tri/Electronics 5.300°, erdent - re- revelations We got this GPS for my husband who is an (OTR). This a professional OTR into diver, and hou. Well, what can i say. I've had the unit in m. Not grang to wise a long review, even thought. Due had mine for a year and here's what we go.	overall 6 1 3 2 1 h	es+True) Sumary Octa have GP81 Very Diagoniete 1st Impression Gree grade, POCH Gree tepr Issues, only excuses for support	un1x8ev1ew11xe 1370131200 1290643200 1283990400 129065600 1317254400	reviewTine 06 2, 2013 11 25, 2010 09 9, 2010 11 24, 2010 09 29, 2011							

[] from google.colab import drive drive.mount('/content/drive')

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

Figure 2: Mounting drive with Google colab

from tqdm.notebook import tqdm #used as a progress bar import pandas as pd import torch import numpy as np import matplotlib.pyplot as plt import seaborn as sns from wordcloud import WordCloud, STOPWORDS from sklearn.feature extraction.text import CountVectorizer from textblob import TextBlob

Figure 3: Import the import libraries

To run the code successfully required packages will be imported. Word cloud, Stopwords CountVectorisation (Pankaj et al. 2019) are important libraries for data cleaning. Figure 2a shows the list of required libraries that have been imported. from tqdm.notebook import tqdm #used as a progress bar import pandas as pd import torch import numpy as np import matplotlib.pyplot as plt import seaborn as sns from wordcloud import WordCloud, STOPWORDS from sklearn.feature_extraction.text import CountVectorizer from textblob import TextBlob

Figure 2a: List of libraries imported

Data Read and Load: Dataset is in JSON format and uploaded in Drive location of /content/drive/My Drive/Amazon_Product_Sentiment_IE/Electronics_5.json. As shown in below figure 3, A data frame was created to read the dataset.



Fig A: Google Drive.

ev.	1ew_dt.head()								
	reviewerID	asin	reviewerName	helpful	reviewText	overall	summary	unixReviewTime	reviewTime
D	AO94DHGC771SJ	0528881469	amazdnu	[0, 0]	We got this GPS for my husband who is an (OTR)	5	Gotta have GPSI	1370131200	06 2, 2013
ľ	AMO214LNFCEI4	0528881469	Amazon Customer	[12, 15]	I'm a professional OTR truck driver, and I bou	1	Very Disappointed	1290643200	11 25, 2010
2	A3N7T0DY83Y4IG	0528881469	C. A. Freeman	[43, 45]	Well, what can I say. I've had this unit in m	3	1st impression	1283990400	09 9, 2010
	A1H8PY3QHMQQA0	0528881469	Dave M. Shaw "mack dave"	[9, 10]	Not going to write a long review, even thought	2	Great grafics, POOR GPS	1290556800	11 24, 2010
ł	A24EV6RXELQZ63	0528881469	Wayne Smith	[0, 0]	I've had mine for a year and here's what we go	1	Major issues, only excuses for support	1317254400	09 29, 2011

Figure 3

Figure 4 shows the categorization of data into positive, negative, and neutral categories with the help of the overall rating column. A positive review will be for 4 and 5 ratings, Negative (1 and 2) and Neutral review (3)



array(['positive', 'negative', 'neutral'], dtype=object)

Data Cleaning:

Below figure 5 shows that null value rows are being dropped.

```
[ ] # Dropping null values to have a cleaner dataset
    review_df = review_df.dropna(subset=['reviewText'])
```

Figure 5

Once data is divided into positive, negative, and natural ratings. A bar graph Figure 6 is plotted to visualize the ratio of Positive, negative, and neutral ratings





Percentage of neutral, negative, positive words in train and test data def pert_count(data, category): return (len(data[data["classes"] == category])/len(data)) * 100 print(f"Percentage of neutral words in train --> {pert_count(review_df, 'neutral")} %") print(f"Percentage of negative words in train --> {pert_count(review_df, 'negative')} %") print(f"Percentage of positive words in train --> {pert_count(review_df, 'negative')} %")

□ Percentage of neutral words in train --> 8.421620328820712 % Percentage of negative words in train --> 11.299156754606356 % Percentage of positive words in train --> 80.27922291657293 %





Figure 7a : Most repeated words

Below diagram 8 shows different data cleaning methods for removing punctuation



marks, Changing text to lower case.

Tokenization of the text data needs to be down Figure 9.



Figure 9

For implementing machine learning models below libraries needed to be imported as shown in Fig 10.



Figure 10

Dataset is divided into test and train set in the ratio of 20: 80. This will help to analyze the ML algorithms easily. Below Fig 11 show the representation



Splitting Complete

✤ <u>Implementation of</u> <u>Machine Learning Models:</u>

1- MultinomialNB model (MNB):



Figure 12

2- MLP classifier (Multilayer Perceptron model)

3- SGD Classifier Model (sklearn.linear model):

```
from sklearn.linear_model import SGDClassifier
  def SGD_Count_Vec(X_train,y_train,X_test,y_test):
    global acc3
    clf = SGDClassifier(alpha=0.00001)
    clf.fit(X_train,y_train)
    open('/content/drive/My Drive/Amazon_Product_Sentiment_IE/classifier/sgd_cvec.pkl', 'wb').write(pickle.dumps(clf))
print("SGDClassifier :train set")
    y_pred = clf.predict(X_train)
    #pred=clf.predict_proba(X_test)
    print("SGDClassifier using Count Vectorizer :Confusion Matrix: ", confusion_matrix(y_train, y_pred))
    print ("SGDClassifier using Count Vectorizer :Accuracy : ", accuracy_score(y_train,y_pred)*100)
print("SGDClassifier :Test set")
    v pred = clf.predict(X test)
    print("SGDClassifier using Count Vectorizer :Confusion Matrix: ", confusion_matrix(y_test, y_pred))
    print ("SGDClassifier using Count Vectorizer :Accuracy : ", accuracy_score(y_test,y_pred)*100)
    acc3 =accuracy_score(y_test,y_pred)*100
    #confusion Matrix
    matrix =confusion_matrix(y_test, y_pred)
    class_names=['Negative', 'Neutral', 'Possitive']
fig, ax = plt.subplots()
    tick_marks = np.arange(len(class_names))
    plt.xticks(tick_marks, class_names)
    plt.yticks(tick_marks, class_names)
    sns.heatmap(pd.DataFrame(matrix), annot=True, cmap="YlGnBu" ,fmt='g')
    ax.xaxis.set_label_position("top")
    plt.tight_layout()
    plt.title('Confusion matrix', y=1.1)
    plt.ylabel('Actual label')
    plt.xlabel('Predicted label')
    plt.show()
    #Classification Report
    target_names = ['Negative', 'Neutral', 'Possitive']
    prediction=clf.predict(X_test)
    print(classification_report(y_test, prediction, target_names=target_names))
    classes = ['Negative', 'Neutral','Possitive']
    visualizer = ClassificationReport(clf, classes=classes, support=True)
    visualizer.fit(X_train, y_train)
    visualizer.score(X_test, y_test)
    g = visualizer.poof()
  SGD_Count_Vec(X_train_cVec,y_train_cVec,X_test_cVec,y_test_cVec)
```

Results:

The visualization of output is done by confusion matrix (Fig 14) and classifier classification matrix (Fig 15) MLPClassifier :train set



Figure 14: Confusion Matrix



Figure: 15: Classification report

References:

Pankaj, Pandey, P., Muskan and Soni, N. (2019). Sentiment analysis on customer feedback data: Amazon product reviews, *2019 International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COMITCon)*, pp. 320–322.