	22\nSkipping line 500314: expected 15 fields, saw 22\nSkipping line 500339: expected 15 fields, saw 22\nSkipping line 507760: expected 15 fields, saw 22\nSkipping line 510760: expected 15 fields, saw 22\nSkipping line 527638: expected 15 fields, saw 22\nSkipping line 534209: expected 15 fields, saw 22\nSkipping line 535687: expected 15 fields, saw 22\nSkipping line 547671: expected 15 fields, saw 22\nSkipping line 549054: expected 15 fields, saw 22\nSkipping line 547671: expected 15 fields, saw 22\nSkipping line 604776: expected 15 fields, saw 22\nSkipping line 70217: expected 15 fields, saw 22\nSkipping line 604776: expected 15 fields, saw 22\nSkipping line 70217: expected 15 fields, saw 22\nSkipping line 804776: expected 15 fields, saw 22\nSkipping line 904776: expected 15 fields, saw 22\nSkipping line 904776: expected 15 fields, saw 22\nSkipping line 906776: expected 15 fields, saw 22\nSkipping line 906776: expected 15 fields, saw 22\nSkipping line 100747: expected 15 fields, saw 22\nSkipping line 1018984: expected 15 fields, saw 22\nSkipping line 1006478: expected 15 fields, saw 22\nSkipping line 1006478: expected 15 fields, saw 22\nSkipping line 1006478: expected 15 fields
	saw 22\nSkipping line 1066869: expected 15 fields, saw 22\nSkipping line 1068809: expected 15 fields, saw 22\nSkipping line 1087983: expected 15 fields, saw 22\nSkipping line 1108184: expected 15 fields, saw 22\nSkipping line 1118137: expected 15 fields, saw 22\nSkipping line 1142723: expected 15 fields, saw 22\nSkipping line 1156947: expected 15 fields, saw 22\nSkipping line 1172563: expected 15 fields, saw 22\nSkipping line 1209254: expected 15 fields, saw 22\nSkipping line 1217598: expected 15 fields, saw 22\nSkipping line 1237598: expected 15 fields, saw 22\nSkipping line 1273825: expected 15 fields, saw 22\nSkipping line 1273825: expected 15 fields, saw 22\nSkipping line 127898: expected 15 fields, saw 22\nSkipping line 1286023: expected 15 fields, saw 22\nSkipping line 1302038: expected 15 fields, saw 22\nSkipping line 138602: expected 15 fields, saw 22\nSkipping line 1338503: expected 15 fields, saw 22\nSkipping line 1341513: expected 15 fields, saw 22\nSkipping line 1346493: expected 15 fields, saw 22\nSkipping line 1373127: expected 15 fields, saw 22\nSkipping line 1341513: expected 15 fields, saw 22\nSkipping line 1433626: expected 15 fields, saw 22\nSkipping line 1442698: expected 15 fields, saw 22\nSkipping line 1443626: expected 15 fields, saw 22\nSkipping line 1442698: expected 15 fields, saw 22\nSkipping line 14500636: expected 15 fields, saw 22\nSkipping line 1537952: expected 15 fields, saw 22\nSkipping line 1594217: expected 15 fields, saw 22\nSkipping line 1541020: expected 15 fields, saw 22\nSkipping line 1615907: expected 15 fields, saw 22\nSkipping line 1541020: expected 15 fields, saw 22\nSkipping line 1615907: expected 15 fields, say 22\
	saw 22/nskipping line 162189; expected 15 fields, saw 22/nskipping line 1671837; expected 15 fields, saw 22/nskipping line 1673837; expected 15 saw 22/nskipping line 1674523; expected 15 fields, saw 22/nskipping line 1744482; expected 15 fields, saw 22/nskipping line 1744482; expected 15 fields, saw 22/nskipping line 1744482; expected 15 fields, saw 22/nskipping line 1831899; expected 15 fields, saw 22/nskipping line 1874790; expected 15 saw 22/nskipping line 1879952; expected 15 fields, saw 22/nskipping line 1879952; expected 15 fields, saw 22/nskipping line 1879952; expected 15 fields, saw 22/nskipping line 1879892; expected 15 fields, saw 22/nskipping line 1878969; expected 15 fields, saw 22/nskipping line 18786969; expected 15 fields, saw 22/nskipping line 1886989; expected 15 fields, saw 22/nskipping line 2886182; expected 15 fields, saw 22/nskipping line 2886183; expected 15 fields, saw 22/nskipping
:	Labelling Reviews:  The reviews with rating 4,5 are labelled to be 1 and 1,2 are labelled as 0. Discard the reviews with ratin discarded_count = len(data[data['star_rating'] == 3.0])
	<pre>class 'pandas.core.frame.DataFrame'&gt; Int64Index: 4525105 entries, 0 to 4874889 Data columns (total 3 columns): # Column</pre>
	<pre>neg_data = data[data['label']==0]  # print("Positive Reviews: ", len(pos_data)) # print("Negative Reviews: ", len(neg_data))  pos = pos_data.sample(n=100000, random_state=1) neg = neg_data.sample(n=100000, random_state=1)  data_both = pd.concat([pos, neg], ignore_index=True) data_both = data_both.sample(frac=1).reset_index(drop=True) # data_both.head(10)  Positive Reviews: 3856296 Negative Reviews: 668809</pre>
	<pre>from statistics import mean  data_both['review_body'] = data_both['review_body'].astype(str)  avg_len1 = mean([len(s) for s in data_both["review_body"]])  #before cleaing and pre-processing # print("BEFORE CLEANING") # print(*[rev+", length: "+str(len(rev))+"\n" for rev in data_both["review_body"].head(3)]) # print("Average review length before processing: "+ str(avg_len1))  BEFORE CLEANING It is as pictured. It is cool, unique gift idea. I also purchased "The Best Ice Mold" (Amazon) and the round ice ball fits perfect his glass. This was a birthday themed idea that I put together for my husband. I also purchased a bottle of Markers Mark whiskey (local) M glass and coaster (Amazon). He really liked it. FYI: I also bought some real good cigars (locally) and a travel humidor (Amazon). To really liked it. I also purchased a place length but a func a current but the whole theme.</pre>
	the whole themebecause lets face it, buying a surprise birthday gift for my husband of 30 years is a challenge LoL but a fun one. I reced with this theme.:-), length: 612  The tags are satin finish and very hard to write on. Either the ink wipes off (sharpie) or is hard to write on (ball point)., length: 124  This was really cheap made. It serves the purpose but tore very easily, length: 70  Average review length before processing: 323.274835  Data Cleaning  #regexs for expressions to be removed  replace_url = re.compile(r'https?:\/\/.*[\r\n]*')  remove_href = re.compile(r'\ <a href')="" remove_amp="re.compile(r'\amp;')" remove_sym="re.compile(r'[_&quot;\-;%() +&amp;=*%.,~!?:#\$\[\]/]')&lt;/td"></a>
	<pre>remove_br = re.compile(r' ') remove_space = re.compile(r' +') replace_brackets = re.compile('[/(){}\[]\],;+]') remove_numbers = re.compile("[\d]")  def cleaning_text(text):     text = text.lower() # lowercase text     text = replace_url.sub("", text) # delete URLs from text     text = replace_url.sub("", text) # delete numbers tags from text     text = remove_numbers.sub("", text) # delete numbers tags from text     text = remove_href.sub('', text) # delete HTML tags from text     text = remove_amp.sub('', text) # delete "&amp;" from text     text = remove_sym.sub('', text) # replace "remove_sym" symbols by space in text     text = remove_br.sub('', text) # delete HTML tags from text     text = remove_brs.sub('', text) # replace "replace_brackets" symbols by space in text     text = replace_brackets.sub('', text) # replace_brackets" symbols by space in text     text = remove_space.sub('', text) # replace_brackets" symbols by space in text     text = remove_space.sub('', text)</pre>
:	return text  #applying all the above cleaning to the dataset data_both['review_body'] = data_both['review_body'].apply(cleaning_text)  # data_both.head(10)  review_id
:	4 RIIP7DH28CEL72 unacceptable for my cooktop which requires tha 1.0 0  5 R2W5RCO3L4VBHH the only downside i see is that it takes up a 5.0 1  6 RINIVXSOO7SH3J i'm shocked this has good reviews i just got i 2.0 0  7 RIRSNTGPQMLKZU really like this pan perfect for large dishes 4.0 1  8 R3QFR9AAC4DGWQ i purchased these when they were advertised as 2.0 0  9 R256OUIH01QXW quite large and a little cheap looking i have 2.0 0  perform contractions on the reviews.  def contractionfunction(s):  s = re.sub(r"won\'t", "will not", s) s = re.sub(r"an\'t", "can not", s) s = re.sub(r"can\'t", "can not", s) s = re.sub(r"let\'s", "let us", s) s = re.sub(r"let\'s", "let us", s) s = re.sub(r"n\'t", "modam', s) s = re.sub(r"n\'t", "modam', s) s = re.sub(r"\'t", "i not", s) s = re.sub(r"\'t", "would", s) s = re.sub(r"\'t", "would", s) s = re.sub(r"\'t", "modam', s)
	data_both['review_body'] = data_both['review_body'].apply(contractionfunction)  review_id
	RINNYXSOOTSH33 iam shocked this has good reviews just got 2.0 0  RIGHTON TEACH TO THE PROPORTIES OF THE PROPORESSING IN THE PROPORTIES OF THE PROP
	<pre>remove the stop words  from nltk.corpus import stopwords  def remove_stopwords(text):     STOPWORDS = set(stopwords.words('english'))     text = ' '.join(word for word in text.split() if word not in STOPWORDS) # delete stopwors from text     return text  data_both['review_body'] = data_both['review_body'].apply(remove_stopwords) # data_both.head(10)  review_id</pre>
	1 R3BN73A6RT3LFP tags satin finish hard write either ink wipes 1.0 0 2 R1QRZ3W592MREX really cheap made serves purpose tore easily 2.0 0 3 R1T1KAXBW2HJEC product shipped present carry overseas paid da 1.0 0 4 R1IP7DH28CEL72 unacceptable cooktop requires cookware totally 1.0 0 5 R2W5RCO3L4VBHH downside see takes little room kitchen larger 5.0 1 6 R1NIVXSOO7SH3J shocked good reviews got mail today excited us 2.0 0 7 R1RSNTGPQMLK2U really like pan perfect large dishes used calp 4.0 1 8 R3QFR9AAC4DGWQ purchased advertised popsicle holders pain kno 2.0 0 9 R256OUIHOI1QXW quite large little cheap looking seen much nic 2.0 0  perform tokenization and lemmatization  from nltk.stem import WordNetLemmatizer from nltk.tokenize import TweetTokenizer from nltk.stem import PorterStemmer  #creating an tokenizer object text_tokenizer = TweetTokenizer(strip_handles=True, reduce_len=True)  def tokenize text  #tokenize text
: .	return review_tokens  #tokenizing each review in the dataset data_both['review_body'] = data_both['review_body'].apply(tokenize_text)  # data_both.head(10)  review_id review_body star_rating label  0 RZMWGGF37ISTO [pictured, cool, unique, gift, idea, also, pur 4.0 1  1 R3BN73A6RT3LFP [tags, satin, finish, hard, write, either, ink 1.0 0  2 R1QRZ3W59ZMREX [really, cheap, made, serves, purpose, tore, e 2.0 0  3 R1T1KAXBW2HJEC [product, shipped, present, carry, overseas, p 1.0 0
:	4 R1IP7DH28CEL72 [unacceptable, cooktop, requires, cookware, to 1.0 0  5 R2W5RCO3L4VBHH [downside, see, takes, little, room, kitchen, 5.0 1  6 R1NIVXSOO7SH3J [shocked, good, reviews, got, mail, today, exc 2.0 0  7 R1RSNTGPQMLK2U [really, like, pan, perfect, large, dishes, us 4.0 1  8 R3QFR9AAC4DGWQ [purchased, advertised, popsicle, holders, pai 2.0 0  9 R256OUIH0I1QXW [quite, large, little, cheap, looking, seen, m 2.0 0
:	<pre>#Lemmatizer object lemmatizer=WordNetLemmatizer()  def lemmatize_tokens(tokens):     #stemming     stemmed_tokens=[stemmer.stem(t) for t in tokens]  #Stemming with lemmatization lemmatized_tokens = [lemmmatizer.lemmatize(st) for st in stemmed_tokens] lemmatize_review = ' '.join(lemmatized_tokens)  return lemmatize_review  data_both['review_body'] = data_both['review_body'].apply(lemmatize_tokens)  from statistics import mean  avg_len3 = mean([len(s) for s in data_both["review_body"]])  # #after cleaing and pre-processing print("AFTER PREPROCESSING") # print("*[rev+", length: "+str(len(rev))+"\n" for rev in data_both["review_body"].head(3)])</pre>
	# print("Average review lentgh after processing: "+ str(avg_len3))  AFTER PREPROCESSING pictur cool uniqu gift idea also purchas best ice mold amazon round ice ball fit perfectli glass birthday theme idea put togeth husband als s bottl marker mark whiskey local mm glass coaster amazon realli like fyi also bought real good cigar local travel humidor amazon round who let face buy surpris birthday gift husband year challeng lol fun one realli score theme, length: 382 tag satin finish hard write either ink wipe sharpi hard write ball point, length: 72 realli cheap made serv purpos tore easili, length: 41  Average review lentgh after processing: 176.841055  from sklearn.model_selection import train_test_split  #splitting the dataset into train and test sets X = data_both.review_body
Ξ.	<pre>y = data_both.label X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)  Perceptron  from sklearn.pipeline import Pipeline     from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer     from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score     from sklearn.linear_model import Perceptron  #perceptron for reviews     perceptron = Pipeline([('vect', CountVectorizer()),</pre>
	<pre># print('Train Accuracy: %s' % accuracy_score(y_pred_perceptron_train, y_train)) # print('Train Precision: %s' % precision_score(y_pred_perceptron_train, y_train)) # print('Train Recall: %s' % recall_score(y_pred_perceptron_train, y_train)) # print('Train F1: %s \n' % f1_score(y_pred_perceptron_train, y_train))  #predicting the classes on test y_pred_perceptron = perceptron.predict(X_test)  #metrics # print('Test Accuracy: %s' % accuracy_score(y_pred_perceptron, y_test)) # print('Test Precision: %s' % precision_score(y_pred_perceptron, y_test)) # print('Test Recall: %s' % recall_score(y_pred_perceptron, y_test)) # print('Test F1: %s' % f1_score(y_pred_perceptron, y_test))</pre>
	Train Accuracy: 0.88214375 Train Precision: 0.910613269396471 Train Recall: 0.8617466320018908 Train F1: 0.885506287227002  Test Accuracy: 0.84735 Test Precision: 0.8781185683449626 Test Recall: 0.8262717868782768 Test F1: 0.8514065998247835  SVM  #SVM classifier for reviews from sklearn.linear_model import SGDClassifier
	<pre>svm_classifier = Pipeline([('vect', CountVectorizer()),</pre>
	Train Accuracy: 0.91041875 Train Precision: 0.905355609885238 Train Recall: 0.9147803265450324 Train F1: 0.9100437448613908  Test Accuracy: 0.893275 Test Precision: 0.8891119923698609 Test Recall: 0.8958122597612785 Test F1: 0.8924495502985413  Logistic Regression #logistic regression model for reviews from sklearn.linear_model import LogisticRegression
	<pre>log_reg = Pipeline([('vect', CountVectorizer()),</pre>
	<pre># print('Test F1: %s' % f1_score(y_pred_ir, y_test)) LR Train Accuracy: 0.9239125 Train Precision: 0.9226264064236567 Train Recall: 0.9251449393320895 Train F1: 0.9238839564836815  Test Accuracy: 0.881875 Test Precision: 0.8807790773555544 Test Recall: 0.8818858061921995 Test F1: 0.8813320943315669  Naive Bayes  #Multinomial Naive Bayes Classifier for reviews from sklearn.naive_bayes import MultinomialNB  multi_naive_bayes = Pipeline([('vect', CountVectorizer()),</pre>
	<pre>multi_naive_bayes = Pipeline((('vect', countvectorizer()),</pre>
	NB Train Accuracy: 0.87893125 Train Precision: 0.8714269658712022 Train Recall: 0.884919729133379 Train F1: 0.8781215198474868  Test Accuracy: 0.867075 Test Precision: 0.8590934190050701 Test Recall: 0.8720953933958419 Test F1: 0.8655455809836895  #final output  print("\nNeutral Reviews:", discarded_count, "Positive Reviews:", len(pos_data), "Negative Reviews:", len(neg_data), "\n")
	<pre>print(str(avg_len1)+","+ str(avg_len2)) print(str(avg_len2)+","+ str(avg_len3)+"\n")  print(str(accuracy_score(y_pred_perceptron_train, y_train))+","+str(precision_score(y_pred_perceptron_train, y_train))+","+ str(recall_score) print(str(accuracy_score(y_pred_svm_train, y_train))+","+str(precision_score(y_pred_svm_train, y_train))+","+ str(recall_score(y_pred_svm_train, y_train))+","+ str(recall_score(y_pred_svm_train, y_train))+","+ str(recall_score(y_pred_lr_train, y_train))+","+ str(recall_score(y_pred_lr_train, y_train))+","+ str(recall_score(y_pred_nb_train, y</pre>