1 1 1	warnings.filterwarnings('ignore') import re from nltk.corpus import stopwords from sklearn.metrics.pairwise import linear_kernel from sklearn.feature_extraction.text import CountVectorizer from sklearn.feature_extraction.text import TfidfVectorizer
221.	# Reading the dataset zomato_data = pd.read_csv("/Users/sachin/Downloads/NLP/zomato.csv") zomato_data.head(5)
93]:	url address name online_order book_table rate votes phone location rest_type dish_liked cuisines approx_cost(for two people) reviews_list Pasta, Lunch North North North (Pated 4.0)
0	https://www.zomato.com/bangalore/jalsa-banasha Road, 2nd Stage, Jalsa Yes Yes 4.1/5 775 42297555\r\n+91 Banashankari Dining Banashankari Dining Papad, Paneer Laja Horth Indian, Waghlai, Chinese Laja Horth Indian, Waghlai, Chinese Paneer Laja Womos, It (Pated 4.0), Rated
1	https://www.zomato.com/bangalore/spice- Feet Road, Spice Yes No 4.1/5 787 080 41714161 Banashankari Casual Buffet, North Bazaar, 6th Casual Buffet, North Had been Here for din 1112, Next to San
3	https://www.zomato.com/bangalore/addhuri-udupi https:
4	Banashankar 10, 3rd Floor, Lakshmi Grand Village Associates, Village Gandhi Baza 10, 3rd Floor, Lakshmi Grand Village No No 3.8/5 166 8026612447\r\n+91 Basavanagudi 9901210005 10, 3rd Floor, Lakshmi Grand No No 3.8/5 166 8026612447\r\n+91 Basavanagudi Panipuri, Dining Gol Gappe Rajasthani Casual Panipuri, North Indian, Rajasthani Very good restaurant
De pe	ow the next step is data cleaning and feature engineering for this step we need to do a lot of stuff with the data such as: eleting Unnecessary Columns\ Removing the Duplicates\ Remove the NaN values from the dataset\ Changing the column names\ Data Transformations\ Data Cleaning\ Adjust the column names Now, let's erform all the above steps in our data:
)5]: _#	# Deleting Unnecessary columns zomato = zomato_data.drop(["url","dish_liked","phone"],axis=1) # Removing the duplicates zomato.duplicated().sum() # check number of duplicated rows
96]:	zomato.drop_duplicates(inplace= True) # Removing NaN zomato.isnull().sum() # check number of null values in each columns zomato.dropna(how="any",inplace= True)
97]:	# Renaming Columns zomato.rename(columns={"approx_cost(for two people)":"cost","listed_in(type)":"type","listed_in(city)":"city"},inplace=True)
# 2 2	<pre># Transformations # Removing "," from cost zomato['cost'] = zomato['cost'].astype(str) zomato['cost'] = zomato['cost'].apply(lambda x : x.replace(",","")) zomato['cost'] = zomato['cost'].astype(float)</pre>
2 2	# Removing "NEW","-" and "/5" from rate zomato = zomato.loc[zomato['rate']!="NEW"] zomato = zomato.loc[zomato['rate']!='-'].reset_index(drop=True) remove_slash = lambda x: x.replace("/5","") if type(x)==np.str else x zomato['rate'] = zomato['rate'].apply(remove_slash).str.strip().astype(float)
L04 #	# Adjusting the column name zomato.name = zomato.name.apply(lambda x:x.title()) zomato['book_table'].replace(("Yes","No"),(True,False), inplace=True) zomato['online_order'].replace(("Yes","No"),(True,False), inplace=True)
114	# Computing mean Rating for each restaurant(Feature Engineering) ratings=pd.DataFrame(zomato.groupby('name')['rate'].mean().reset_index()) ratings
	name rate 0 #Feeltheroll 3.400 1 #L-81 Cafe 3.900 2 #Refuel 3.700
6	3 1000 B.C 3.200 4 100Ã Â Ã Â Â Â Â Â Â Â Â Â Â Â Â Â Â Â Â
6: 6:	567 Zoey'S 4.300 568 Zoroy Luxury Chocolate 4.000 569 Zu'S Doner Kebaps 3.700 570 Zyara 3.875 571 Zyksha 3.400
65	zomato=pd.merge(zomato, ratings,how='left', on='name')
	zomato.rename(columns={"rate_y":"mean_rate"},inplace=True) zomato address name online_order book_table rate_x votes location rest_type cuisines cost reviews_list menu_item type city rate_y r 942, 21st Main Road, 0 2nd Stage lalsa True True 4.1 775 Banashankari Casual Mughlai 800.0 ('RATED\n A II Buffet Banashankari 4.118182 4.1183)
	2nd Floor, 80 Feet Road, Near Big Bazaar, 6th 2nd Floor, 80 Feet Bazaar, 6th 3nd Floor, 80 Feet Bazaar, 8th 3nd Floor, 8th
	1112, Next to KIMS 2 Medical College, 17th Cross San Churro Cafe True False 3.8 918 Banashankari Cafe, Casual Dining Italian Dining Cafe, Mexican, Ambience is not that 1st Floor, Annakuteera, 3rd Addhuri Udupi Rhoiana False False 3.7 88 Banashankari Quick South Indian, North Indian, North Indian, North Indian, North Indian, North Indian, Sites North Indian, North Indian, Sites S
	Stage, Banashankar 10, 3rd Floor, Lakshmi 4 Associates, Gandhi Baza Bhojana False False S.7 66 Banashankar Bites North Indian Bites North Indian False False S.7 66 Banashankar Bites North Indian False False False False S.7 66 Banashankar Bites North Indian False Fals
	136, SAP Labs India, KIADB Export Promotion In The Farm House Bar N Grill False False
	Table 139/CI, Next 10 GR 1233 Tech Park, Pattandur Agraha Four Points by Sheraton Bengaluru, Points By Sheraton Bengaluru False False 2.5 81 Whitefield Dining, Indian, Chinese, 800.0 'RATED\n A fine place to chill bars [('Rated 4.0', Pubs and Whitefield 2.283333 2.2833) [('Rated 5.0', "RATED\n Food and service are bars and service are bars whitefield 3.600000 3.6000 and service are bars whitefield 3.600000 3.60000 3.60000 3.60000 3.60000 3.60000 3.60000 3.60000 3.60000 3.60000 3.60000 3.600000 3.600000 3.600000 3.60000 3.60000 3.60000 3.60000 3.60000 3.60000 3.60000 3.60000 3.6
	Sheraton Grand Chime - Sheraton Grand Bengaluru Whitefield Hotel & Co Sheraton Grand Chime - Sheraton Grand Bengaluru False True 4.3 236 Road, Bar Finger Food 2500.0 'RATED\n Nice and friendly pla Sheraton Grand Grand Bengaluru Whitefield Hotel & Co Whitefield Hotel & Co The Nest - The Den Bengaluru False False S.4 13 Road, Casual North Indian, 1500.0 'RATED\n Great [] and Whitefield S.400000 3.40000 A.3000 A.
L27 1	1237 rows × 16 columns from sklearn.preprocessing import MinMaxScaler
No.	scaler = MinMaxScaler(feature_range = (1,5)) zomato[['mean_rate']] = scaler.fit_transform(zomato[['mean_rate']]).round(2) ow the next step is to perform some text preprocessing steps which include: ower casing\ Removal of Punctuations\ Removal of Stopwords\ Removal of URLs\ Spelling correction\ Now let's perform the above text preprocessing steps on the data:
L32 #	# Lower casing zomato['reviews_list'] = zomato['reviews_list'].str.lower() # Removal of Punctuations
Ç	<pre>import string punc_to_remove = string.punctuation def remove_punctuations(text): return text.translate(str.maketrans('', '', punc_to_remove)) zomato['reviews_list'] =zomato['reviews_list'].apply(lambda text: remove_punctuations(text))</pre>
(<pre># Removal of Stopwords stop_words = set(stopwords.words('english')) def remove_stopwords(text): return " ".join([word for word in str(text).split() if word not in stop_words]) zomato['reviews_list'] = zomato['reviews_list'].apply(lambda text: remove_stopwords(text))</pre>
L40 #	<pre># Removal of URL def remove_url(text): url_pattern = re.compile(r'https?://\S+ www\.\S+') return url_pattern.sub(r'',text)</pre>
141	<pre>zomato['reviews_list'] = zomato['reviews_list'].apply(lambda text: remove_url(text)) zomato['reviews_list']</pre>
	rated 40 ratedn good restaurant neighbourhood 1232 rated 40 ratedn ambience big spacious lawn use
4: 4: 4: N:	rated 40 ratedn fine place chill office hours 1234 rated 50 ratedn food service incomparably exce 1235 rated 40 ratedn nice friendly place staff awes 1236 rated 50 ratedn great ambience looking nice go 1236 ame: reviews_list, Length: 41237, dtype: object
143	IDF Vectorizer zomato.drop(columns=['address','rest_type', 'type', 'menu_item', 'votes'],inplace= True) df_percent.rename(columns={'mean_rate':'Mean_Rating'},inplace= True)
L49 (# Randomly sampling 60% OF dataset df_percent = zomato.sample(frac=0.6) df_percent.set_index('name',inplace=True) indices = pd.Series(df_percent.index)
L53 #	# Creating TF-IDF Matrix tfidf = TfidfVectorizer(analyzer='word',ngram_range=(1,2),min_df=0, stop_words='english') tfidf_matrix = tfidf.fit_transform(df_percent['reviews_list']) cosine_similarities = linear_kernel(tfidf_matrix, tfidf_matrix)
L56 (cosine_similarities = linear_kernel(tridr_matrix, tridr_matrix) cosine_similarities rray([[1.00000000e+00, 2.77366910e-03, 1.07108322e-02,,
	5.30357196e-03, 2.36064866e-02, 6.79103198e-03], [2.77366910e-03, 1.00000000e+00, 4.93025119e-03,, 6.34095547e-04, 1.13662745e-02, 2.69431239e-03], [1.07108322e-02, 4.93025119e-03, 1.0000000e+00,, 4.72002043e-03, 2.75998475e-02, 1.98241294e-02],, [5.30357196e-03, 6.34095547e-04, 4.72002043e-03,, 1.00000000e+00, 7.74911196e-03, 3.07933572e-03], [2.36064866e-02, 1.13662745e-02, 2.75998475e-02,, 7.74911196e-03, 1.00000000e+00, 2.52733138e-02], [6.79103198e-03, 2.69431239e-03, 1.98241294e-02,, 3.07933572e-03, 2.52733138e-02, 1.00000000e+00]])
	<pre>def recommend(name, cosine_similarities=cosine_similarities): recommend_restaurant = [] # list to append recommended restaurant idx = indices[indices==name].index[0] # Getting the index of restaurant score_series = pd.Series(cosine_similarities[idx]).sort_values(ascending=False) # Getting all similar restaurants and sorting by highest similarity</pre>
	<pre>top_30indices = list(score_series.iloc[0:31].index) # extracting top 30 similar restaurants for val in top_30indices: recommend_restaurant.append(list(df_percent.index)[val]) # Creating a new DF to shown similary restaurants</pre>
	<pre>df_new = pd.DataFrame(columns=['cuisines', 'Mean Rating', 'cost']) for each in recommend_restaurant: df_new = df_new.append(pd.DataFrame(df_percent[['cuisines', 'Mean Rating', 'cost']][df_percent.index == each].sample())) # Dropping Duplicates</pre>
	<pre>df_new = df_new.drop_duplicates(subset=['cuisines','Mean Rating', 'cost'],keep=False) # sorting and keeping only top 10 values df_new = df_new.sort_values(by='Mean Rating',ascending=False).head(10) # Printing print('TOP %s RESTAURANTS LIKE %s WITH SIMILAR REVIEWS: ' % (str(len(df_new)), name))</pre>
	return df_new recommend('Pai Vihar')
.70	OP 10 RESTAURANTS LIKE Pai Vihar WITH SIMILAR REVIEWS: Cuisines Mean Rating cost Burma Burma Asian, Burmese 4.74 1500.0 Lavonne Cafe, Desserts 4.35 800.0 Caffi Å Å Å Å Å Å Å Å Å Å Å Å Å Å Å Å Å Å
C	Cafã Â Ã Â Ã Â Ã Â Ã Â Ã Â Ã Â Ã Â Ã Â Ã Â
	Ilyazsab The House Of ChickenRolls, Kebab3.84250.0FoodhallItalian, Bakery, Fast Food3.801000.01992 Chats - SpaceStreet Food3.45200.0Karavali GrandMangalorean, Seafood, North Indian, Chinese3.45600.0