	Class: TE-10 Batch: K-10
	Roll No : 33241  Assn : 4  Problem Statement :
	Write a program for the Information Retrieval System using appropriate NLP tools (such as NLTK, Open NLP,) and perform following operations a. Text tokenization b. Count word frequency c. Remove stop words d. POS tagging  Neccessary Imports
In [8]:	<pre>import numpy as np import pandas as pd import re import nltk from nltk.corpus import stopwords import string from wordcloud import WordCloud import seaborn as sns import matplotlib.pyplot as plt %matplotlib inline</pre>
[n [9]:	<pre>nltk.download('wordnet') [nltk_data] Downloading package wordnet to</pre>
Out[9]:	<pre>[nltk_data]</pre>
n [10]:	<pre>df = pd.read_csv(r'Resume_Data.csv', encoding = 'utf-8') df['Cleaned_Resume'] = ''</pre> Exploratory Data Analysis
n [11]:	df.head()  Category Resume Cleaned_Resume
ı [12]:	<ul> <li>Data Science Skills * Programming Languages: Python (pandas</li> <li>Data Science Education Details \r\nMay 2013 to May 2017 B.E</li> <li>Data Science Areas of Interest Deep Learning, Control Syste</li> <li>Data Science Skills â□¢ R â□¢ Python â□¢ SAP HANA â□¢ Table</li> <li>Data Science Education Details \r\n MCA YMCAUST, Faridab</li> </ul>
	print("Resume Categories") print(df['Category'].value_counts())  Resume Categories  Java Developer 84  Testing 70  DevOps Engineer 55 Python Developer 48  Web Designing 45  HR 44  Hadoop 42  Sales 40  Data Science 40  Operations Manager 40  Mechanical Engineer 40  Blockchain 40  ETL Developer 40  Arts 36  Database 33  PMO 30
	Electrical Engineering 30 Health and fitness 30 Business Analyst 28 DotNet Developer 28 Automation Testing 26 Network Security Engineer 25 Civil Engineer 24 SAP Developer 24 Advocate 20 Name: Category, dtype: int64
n [13]:	Visualizing types of people who have given the resume  plt.figure(figsize = (10, 10))  # Setting size of plot plt.xticks(rotation = 90)  # Rotating plot to organize horizonta sns.countplot(y = 'Category', data = df)  # Deciding which column of Dataframe
it[13]:	Data Science -
	Advocate - Arts - Web Designing - Mechanical Engineer - Sales -
	Health and fitness -  Civil Engineer -  Java Developer -  Business Analyst -
	SAP Developer -  Automation Testing -  Electrical Engineering -  Operations Manager -
	Python Developer -  DevOps Engineer -  Network Security Engineer -  PMO -
	Database - Hadoop - ETL Developer - DotNet Developer -
	Blockchain - Testing - O C C C COunt
[14]:	<pre>Data Cleaning  def Clean_Resume(resumeText):     Removals = [</pre>
	<pre>'RT cc',</pre>
n [15]:	<pre>df['Cleaned_Resume'] = df.Resume.apply(lambda x: Clean_Resume(x)) df.head()</pre>
it[15]:	CategoryResumeCleaned_Resume0Data ScienceSkills * Programming Languages: Python (pandasSkills Programming Languages P thon pandas1Data ScienceEducation Details \r\nMay 2013 to May 2017 B.EEducation Details Ma 2013 to Ma 2017 B E UIT2Data ScienceAreas of Interest Deep Learning, Control SysteAreas of Interest Deep Learning Control S ste3Data ScienceSkills â□¢ R â□¢ Python â□¢ SAP HANA â□¢ TableSkills R P thon SAP HANA Table4Data ScienceEducation Details \r\n MCA YMCAUST, FaridabEducation Details MCA YMCAUST Faridabad Har
[16]:	<pre>for i in range(len(df)): corpus += df['Cleaned_Resume'][i] corpus[450:1000]</pre>
t[16]:	'ticSearch D3 js DC js Plot1 kibana matplotlib ggplot Tableau Others Regular Expression HTML CSS Angular 6 Logstash Kafka P thon Flask Git Docker computer vision Open CV and understanding of Deep leaning Education Details Data Science Assurance Associate Data Science Assurance Associate Ernst Young LLP Still Details JAVASCRIPT Exprience 24 months jQuer Exprience 24 months P thon Exprience 24 months Comparable Compan Ernst Young LLP description Fraud Investigations and Dispute Services Assurance TEC'
[17]:	Creating the Tokenizer and Tokenizing  tokenizer = nltk.tokenize.RegexpTokenizer('\w+') tokens = tokenizer.tokenize(corpus) # Tokenizing the text into individual
	<pre>words = [word.lower() for word in tokens]  # Transforming all words to lowercase print(len(words))</pre>
[18]:	<pre>Fetching English Stop Words  nltk.download('stopwords') stopwords = nltk.corpus.stopwords.words('english')</pre>
	<pre>[nltk_data] Downloading package stopwords to [nltk_data]</pre>
n [19]:	<pre>words_new = [     word     for word in words     if word not in stopwords</pre>
[20]:	len(words_new) 326374
[21]:	Lemmatization  from nltk.stem import WordNetLemmatizer
	<pre>wnl = WordNetLemmatizer()  lem_words = [     wnl.lemmatize(word)     for word in words_new ]</pre>
n [22]:	<pre>same=0 diff=0 for i in range(0,1832):     if(lem_words[i]==words_new[i]):         same=same+1     elif(lem_words[i]!=words_new[i]):         diff=diff+1 print('Number of words Lemmatized=', diff)</pre>
n [23]:	<pre>print('Number of words not Lemmatized=', same)  Number of words Lemmatized= 311 Number of words not Lemmatized= 1521  freq_dist = nltk.FreqDist(lem_words)</pre>
	<pre>plt.subplots(figsize=(20,12)) freq_dist.plot(30)</pre>
	3500
	3000
	2500
	2000
	stem nent hon on the near and hirran and lie test lie test which here here here here here here here he
t[23]:	exprire pro exprire maharas) descrip de ser exprire de maharas) de developmen exprire de maharas) de developmen exprire de maharas de developmen exprire de descrip de maharas de descrip d
[24]: t[24]:	<pre>mostcommon = freq_dist.most_common(50) mostcommon  [('project', 4071),</pre>
	('exprience', 3829), ('compan', 3578), ('month', 3344), ('detail', 3132), ('description', 3122), ('team', 2159),
	('1', 2142), ('data', 2138), ('management', 2024), ('skill', 1998), ('stem', 1960), ('b', 1696),
	('sql', 1664), ('database', 1533), ('6', 1499), ('client', 1466), ('maharashtra', 1449),
	<pre>('anal', 1435), ('ear', 1418), ('application', 1394), ('service', 1375), ('testing', 1349), ('test', 1297),</pre>
	('requirement', 1274), ('business', 1273), ('e', 1256), ('le', 1237), ('report', 1229), ('development', 1204),
	('server', 1196), ('developer', 1194), ('customer', 1178), ('ltd', 1177), ('process', 1163), ('using', 1124),
	('c', 1088), ('januar', 1086), ('java', 1076), ('engineering', 1055), ('work', 1038), ('pune', 1026),
	('role', 969), ('ing', 925), ('user', 916), ('operation', 895), ('software', 886), ('pvt', 879),
[25]:	<pre>('responsibility', 866), ('sale', 845)]  res=' '.join([i for i in lem_words if not i.isdigit()])</pre>
[26]:	<pre>plt.subplots(figsize=(16,10)) wordcloud = WordCloud(background_color='black', max_words=200, width=1400, height=1200).generate(res) plt.imshow(wordcloud) plt.title('Resume Text WordCloud (100 Words)') plt.axis('off')</pre>
	Resume Text WordCloud (100 Words)  organi ation detail janetwork securit description project januar bachelor detail janetwork securit description project
	ear monthscompanioftere development of an all studio and single part of the constraint of the constrai
	The conducted by the complete distance administrator completed internal compan description and end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge team member unit testing testing exprience database sql end end good knowledge testing end end good knowledge end end good kn
	education detail nagpur maharashtra duration month technical skill operating stem operating stem operating stem operating stem operation
	Si e shell script school script ing web application team lead home websides till date description responsibility home series strong and school series are school series and school series and school series are school series and school series and school series are school series and school series and school series are school series are school series and school series are school series are school school series are school series
	business anall to description working to business anall to business anall to business anall project or description working to business anall to business ana
	anal ing well korea to see the distribution of the second part of the
[27]: t[27]:	Category Resume Cleaned_Resume  Data Science Skills * Programming Languages: Python (pandas Skills Programming Languages P thon pandas
	1 Data Science       Education Details \r\nMay 2013 to May 2017 B.E       Education Details Ma 2013 to Ma 2017 B E UIT         2 Data Science       Areas of Interest Deep Learning, Control Syste       Areas of Interest Deep Learning Control S ste         3 Data Science       Skills â□¢ R â□¢ Python â□¢ SAP HANA â□¢ Table       Skills R P thon SAP HANA Table
	4 Data Science Education Details \r\n MCA YMCAUST, Faridab Education Details MCA YMCAUST Faridabad Har  957 Testing Computer Skills: â□¢ Proficient in MS office ( Computer Skills Proficient in MS office
	958 Testing â□□ Willingness to accept the challenges. â□□ Willingness to a ept the challenges P  959 Testing PERSONAL SKILLS â□¢ Quick learner, â□¢ Eagerne PERSONAL SKILLS Quick learner Eagerne  960 Testing COMPUTER SKILLS & SOFTWARE KNOWLEDGE MS-Power COMPUTER SKILLS SOFTWARE KNOWLEDGE MS Power  961 Testing Skill Set OS Windows XP/7/8/8.1/10 Database MY Skill Set OS Windows XP 7 8 8 1 10 Database MY
n [ 1 ·	962 rows × 3 columns  Skill Set OS Windows XP/7/8/8.1/10 Database MY  Skill Set OS Windows XP / 8 8 1 10 Database MY
in [ ]:	