

Resume Cleaning using NLP Techniques

▼ 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

▼ Necessary Imports

```
1 import numpy as np
2 import pandas as pd
3 import re
4 import nltk
5 from nltk.corpus import stopwords
6 import string
7 from wordcloud import WordCloud
8 import seaborn as sns
9 import matplotlib.pyplot as plt
10 %matplotlib inline
```

```
1 nltk.download('wordnet')
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Unzipping corpora/wordnet.zip.
True
```

```
1 nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
True
```

▼ Importing the dataset

```
1 df = pd.read_csv('/content/Resume_Data.csv', encoding = 'utf-8')
2 df['Cleaned_Resume'] = ''
```

▼ Exploratory Data Analysis

```
1 df.head()
```

	Category	Resume	Cleaned_Resume
0	Data Science	Skills * Programming Languages: Python (pandas...	
1	Data Science	Education Details \r\nMay 2013 to May 2017 B.E...	
2	Data Science	Areas of Interest Deep Learning, Control Syste...	
3	Data Science	Skills â¬ R â¬ Python â¬ SAP HANA â¬ Table...	
4	Data Science	Education Details \r\n MCA YMCAUST, Faridab...	

```
1 print("Resume Categories")
2 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
Blockchain	40
ETL Developer	40
Operations Manager	40
Data Science	40
Sales	40
Mechanical Engineer	40
Arts	36
Database	33
Electrical Engineering	30
Health and fitness	30
PMO	30
Business Analyst	28
DotNet Developer	28
Automation Testing	26
Network Security Engineer	25
SAP Developer	24
Civil Engineer	24
Advocate	20
Name: Category, dtype: int64	

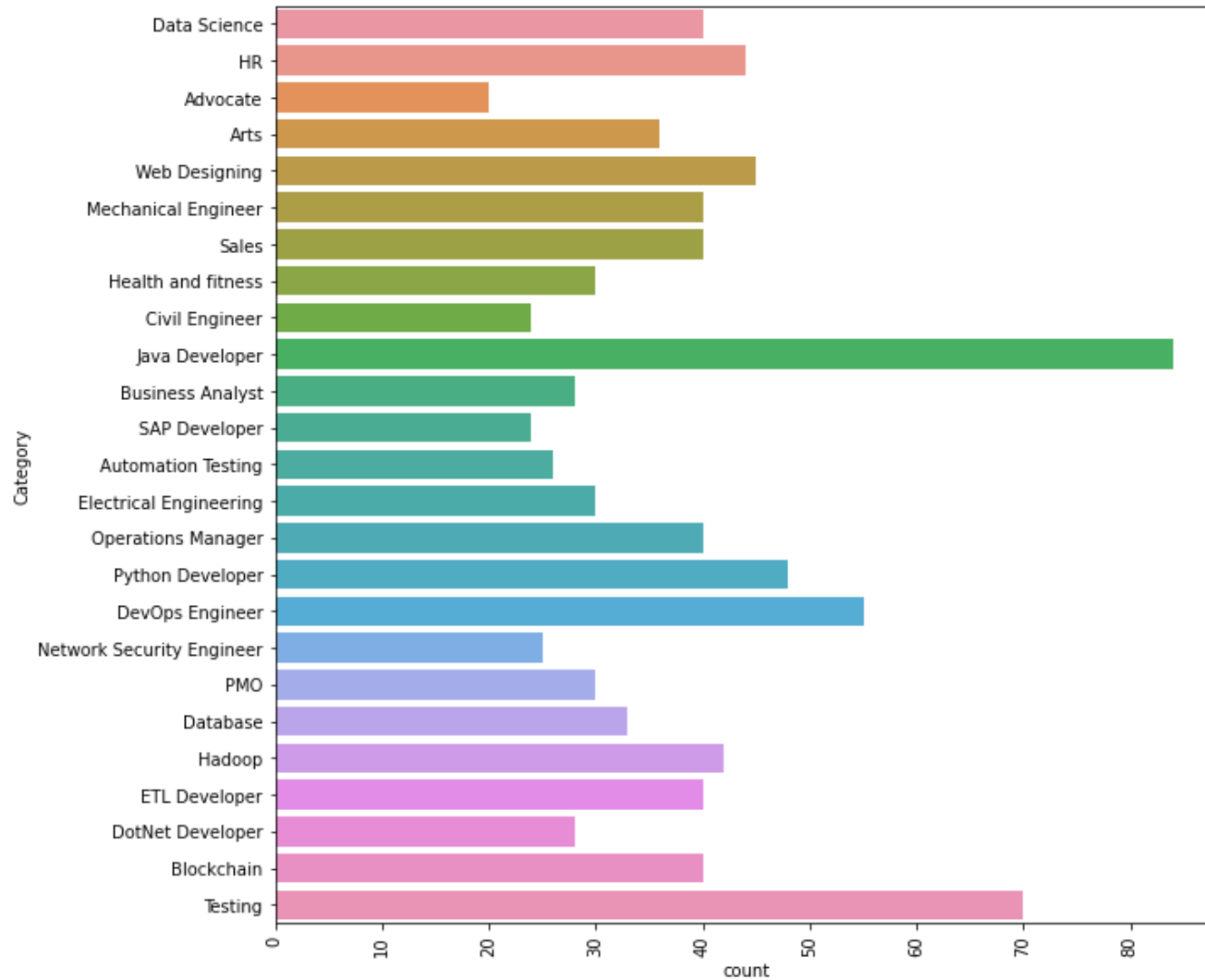
▼ Visualizing types of people who have given the resume

```

1 plt.figure(figsize = (10, 10))           # Setting size of plot
2 plt.xticks(rotation = 90)                 # Rotating plot to organize horizontally
3 sns.countplot(y = 'Category', data = df)  # Deciding which column of Dataframe will the

```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f96de65b810>
```



▼ Data Cleaning

```
1 def Clean_Resume(resumeText):
2     Removals = [                                     # Deciding weeds in resume
3         'http\S+\s*',                                # Web URLs
```

```

4      'RT|cc',                                # Regular characters
5      '#\S+',                                  # Hashtags
6      '@\S+',                                  # Emails
7      '\s+'
8  ]
9
10 for weed in Removals: resumeText = re.sub(weed, ' ', resumeText)    # Removing weeds using regular expression
11 resumeText = re.sub('%s'%re.escape("#$%&'_=-+()[];:.,/?^*@[{}|\~\""), ' ', resumeText)
12 resumeText = re.sub(r'[\x00-\x7f]', r' ', resumeText)
13
14 return resumeText

1 df['Cleaned_Resume'] = df.Resume.apply(lambda x: Clean_Resume(x))
2 df.head()

```

	Category	Resume	Cleaned_Resume
0	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	Skills R Python SAP HANA

```

1 corpus = ''
2 for i in range(len(df)): corpus += df['Cleaned_Resume'][i]
3 corpus[450:1000]

```

```

'ticSearch D3 js DC js Plotl 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 learning Education
Details Data Science Assurance Associate Data Science Assurance Associate Er

```

▼ Creating the Tokenizer and Tokenizing

```
1 tokenizer = nltk.tokenize.RegexpTokenizer('\w+')
2 tokens = tokenizer.tokenize(corpus)           # Tokenizing the text into individual words
3
4 words = [word.lower() for word in tokens]      # Transforming all words to lowercase
5 print(len(words))
```

423116

▼ Fetching English Stop Words

```
1 stopwords = nltk.corpus.stopwords.words('english')
```

▼ Removing Stop words

```
1 words_new = [
2     word
3     for word in words
4     if word not in stopwords
5 ]
```

```
1 len(words_new)
```

326374

▼ Lemmatization

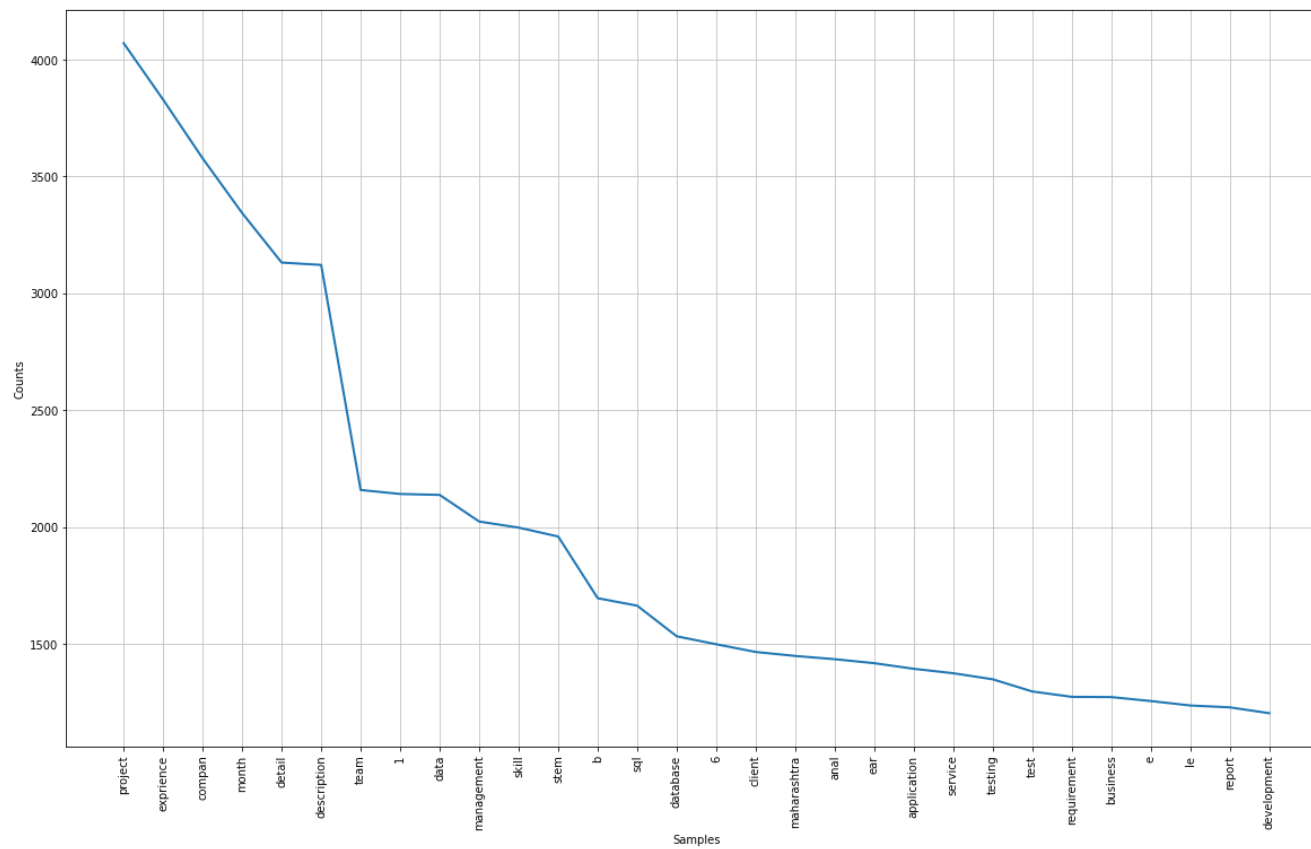
```
1 from nltk.stem import WordNetLemmatizer
2 wnl = WordNetLemmatizer()
3
```

```
4 lem_words = [  
5     wnl.lemmatize(word)  
6     for word in words_new  
7 ]
```

```
1 same=0  
2 diff=0  
3 for i in range(0,1832):  
4     if(lem_words[i]==words_new[i]):  
5         same=same+1  
6     elif(lem_words[i]!=words_new[i]):  
7         diff=diff+1  
8 print('Number of words Lemmatized=', diff)  
9 print('Number of words not Lemmatized=', same)
```

```
Number of words Lemmatized= 311  
Number of words not Lemmatized= 1521
```

```
1 freq_dist = nltk.FreqDist(lem_words)  
2 plt.subplots(figsize=(20,12))  
3 freq_dist.plot(30)
```



```
1 mostcommon = freq_dist.most_common(50)
2 mostcommon
```

```
[('project', 4071),
 ('expreience', 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),
('anal', 1435),
('ear', 1418),
('application', 1394),
('service', 1375),
('testing', 1349),
('test', 1297),
('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),
```

```
('responsibility', 866),  
( 'sale', 845)]
```

```
1 res= ' '.join([i for i in lem_words if not i.isdigit()])
```

```
1 import os  
2 os.system('pip install wordcloud')
```

```
0
```

```
1 plt.subplots(figsize=(16,10))  
2 wordcloud = WordCloud(  
3     background_color='black',  
4     max_words=200,  
5     width=1400,  
6     height=1200  
7     ).generate(res)  
8 plt.imshow(wordcloud)  
9 plt.title('Resume Text WordCloud (100 Words)')  
10 plt.axis('off')  
11 plt.show()
```

[illegible]

1 df

	Category	Resume	Cleaned_Resume
0	Data Science	Skills * Programming Languages: Python	Skills Programming Languages P thon
1	Data Science	Education Details \n\n MCA YMCAUST, Faridab...	Education Details MCA YMCAUST Faridab...
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 \n\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...
		PERSONAL SKILLS & QUALITIES	PERSONAL SKILLS & QUALITIES