

Clustering

In [1]:

```
1 # from google.colab import drive
2 # drive.mount('/gdrive')
3 # %cd /gdrive
```

Importing the data and necessary libraries

In [2]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5 import seaborn as sns; sns.set()
6 import matplotlib.pyplot as plt
7
8 import warnings
9 warnings.filterwarnings('ignore')
10
```

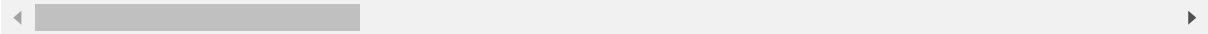
In [3]:

```
1 path = 'After_BERT(Embedded_Data).csv'
2 df = pd.read_csv(path, index_col = 0, keep_default_na=False)
3 df.head()
```

Out[3]:

	Industry	Company Location	New Job (90 Days)	Year Started	Profile Headline	Profile Summary	School	Degr
0	Mechanical or Industrial Engineering	Dublin, Ohio, United States	False	2020.0	Mechanical Design Engineer, System Integration...	In the ever-growing technological world where ...	Chalmers University of Technology	
1	Information Technology and Services	New York City Metropolitan Area	False	2018.0	Digital DevOps Engineer at HSBC	AWS Certified Cloud Engineer holding 3 AWS Ass...	Binghamton University	
2	Information Technology and Services	Chicago, Illinois, United States	False	2018.0	Leading Product + UX at Remedy (Two Point Conv...	http://aroonmathai.com	Carnegie Mellon University	
3	Information Technology and Services	Bangalore Urban, Karnataka, India	False	2018.0	Product Designer at udaan		Vellore Institute of Technology	
4	Information Technology and Services	Jaipur, Rajasthan, India	True	2021.0	Digital Technology Intern at General Electric ...		Vellore Institute of Technology	

5 rows × 30 columns



In [4]:

```
1 coords = pd.read_csv('Coordinates.csv')  
2 coords.head()
```

Out[4]:

	latitude	longitude
0	40.099229	-83.114077
1	44.870970	-0.547490
2	41.875562	-87.624421
3	12.945142	77.553645
4	26.915458	75.818982

In [5]:

```
1 df['latitude'] = coords['latitude']  
2 df['longitude'] = coords['longitude']
```

In [6]:

1 df.info()

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 4910 entries, 0 to 4909
Data columns (total 32 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Industry                             4910 non-null   object
1   Company Location                     4910 non-null   object
2   New Job (90 Days)                   4910 non-null   bool
3   Year Started                        4910 non-null   float64
4   Profile Headline                    4910 non-null   object
5   Profile Summary                     4910 non-null   object
6   School                             4910 non-null   object
7   Degree                             4910 non-null   object
8   Education End                       4910 non-null   object
9   Domain                             4910 non-null   object
10  CompanyName                         4910 non-null   object
11  JobTitle                           4910 non-null   object
12  My Network                         4910 non-null   object
13  Country                            4910 non-null   object
14  Continent                          4910 non-null   object
15  FieldOfStudy                       4910 non-null   object
16  Industry_embedding                 4910 non-null   float64
17  Company Location_embedding         4910 non-null   float64
18  Profile Headline_embedding         4910 non-null   float64
19  Profile Summary_embedding          4910 non-null   float64
20  School_embedding                   4910 non-null   float64
21  Degree_embedding                   4910 non-null   float64
22  Education End_embedding            4910 non-null   float64
23  Domain_embedding                   4910 non-null   float64
24  CompanyName_embedding              4910 non-null   float64
25  JobTitle_embedding                 4910 non-null   float64
26  My Network_embedding               4910 non-null   float64
27  Country_embedding                  4910 non-null   float64
28  Continent_embedding                4910 non-null   float64
29  FieldOfStudy_embedding             4910 non-null   float64
30  latitude                           4910 non-null   float64
31  longitude                           4910 non-null   float64
dtypes: bool(1), float64(17), object(14)
memory usage: 1.2+ MB

```

In [7]:

```
1 df.columns
```

Out[7]:

```
Index(['Industry', 'Company Location', 'New Job (90 Days)', 'Year Started',  
      'Profile Headline', 'Profile Summary', 'School', 'Degree',  
      'Education End', 'Domain', 'CompanyName', 'JobTitle', 'My Network',  
      'Country', 'Continent', 'FieldOfStudy', 'Industry_embedding',  
      'Company Location_embedding', 'Profile Headline_embedding',  
      'Profile Summary_embedding', 'School_embedding', 'Degree_embedding',  
      'Education End_embedding', 'Domain_embedding', 'CompanyName_embeddin  
g',  
      'JobTitle_embedding', 'My Network_embedding', 'Country_embedding',  
      'Continent_embedding', 'FieldOfStudy_embedding', 'latitude',  
      'longitude'],  
      dtype='object')
```

In [8]:

```
1 df1 = df[['Country_embedding',  
2         'Continent_embedding',  
3         'Domain_embedding',  
4         'Industry_embedding',  
5         'FieldOfStudy_embedding',  
6         'latitude', 'longitude']]
```

Standardisation and Normalisation of Data

In [9]:

```
1 # Standardize data  
2 from sklearn.preprocessing import StandardScaler  
3 scaler = StandardScaler()  
4 scaled_df = scaler.fit_transform(df1)
```

In [10]:

```
1 # Normalizing the Data  
2 from sklearn.preprocessing import normalize  
3 normalized_df = normalize(scaled_df)
```

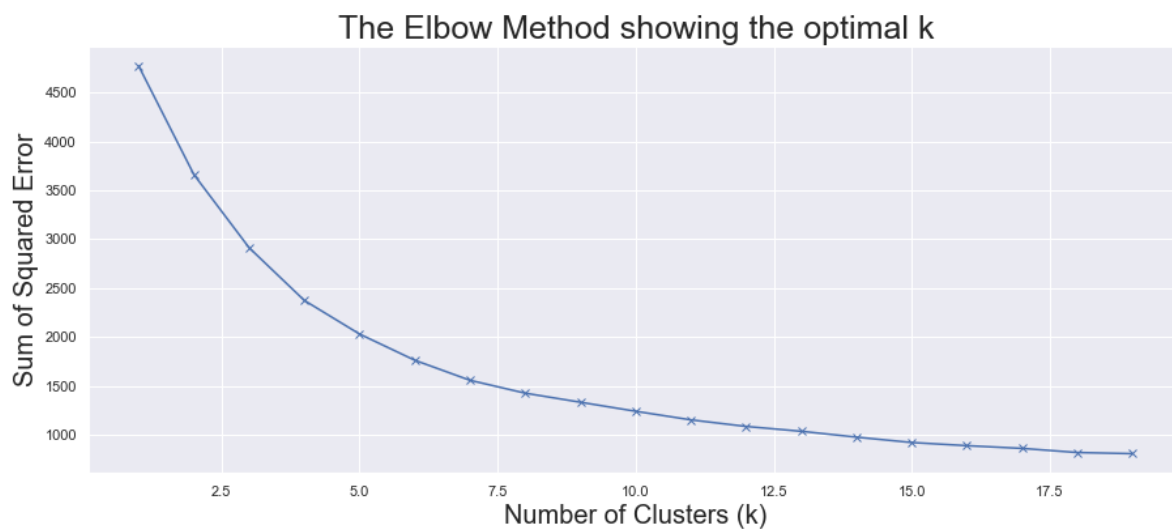
In [11]:

```
1 # Converting the numpy array into a pandas DataFrame  
2 normalized_df = pd.DataFrame(normalized_df)
```

Using Elbow Method to determine the best value of k in KMeans clustering

In [12]:

```
1 from sklearn.cluster import KMeans
2
3 sse = []
4 K = range(1,20)
5 for k in K:
6     kmeanModel = KMeans(n_clusters=k)
7     kmeanModel.fit(normalized_df)
8     sse.append(kmeanModel.inertia_)
9 plt.figure(figsize=(15,6))
10 plt.plot(K, sse, 'bx-')
11 plt.xlabel('Number of Clusters (k)', fontsize = 20)
12 plt.ylabel('Sum of Squared Error', fontsize = 20)
13 plt.title('The Elbow Method showing the optimal k', fontsize = 25)
14 plt.show()
```



Here we see that there is no appropriate value of k. Thus we use DBSCAN method to cluster data.

Visualising the data using PCA

In [13]:

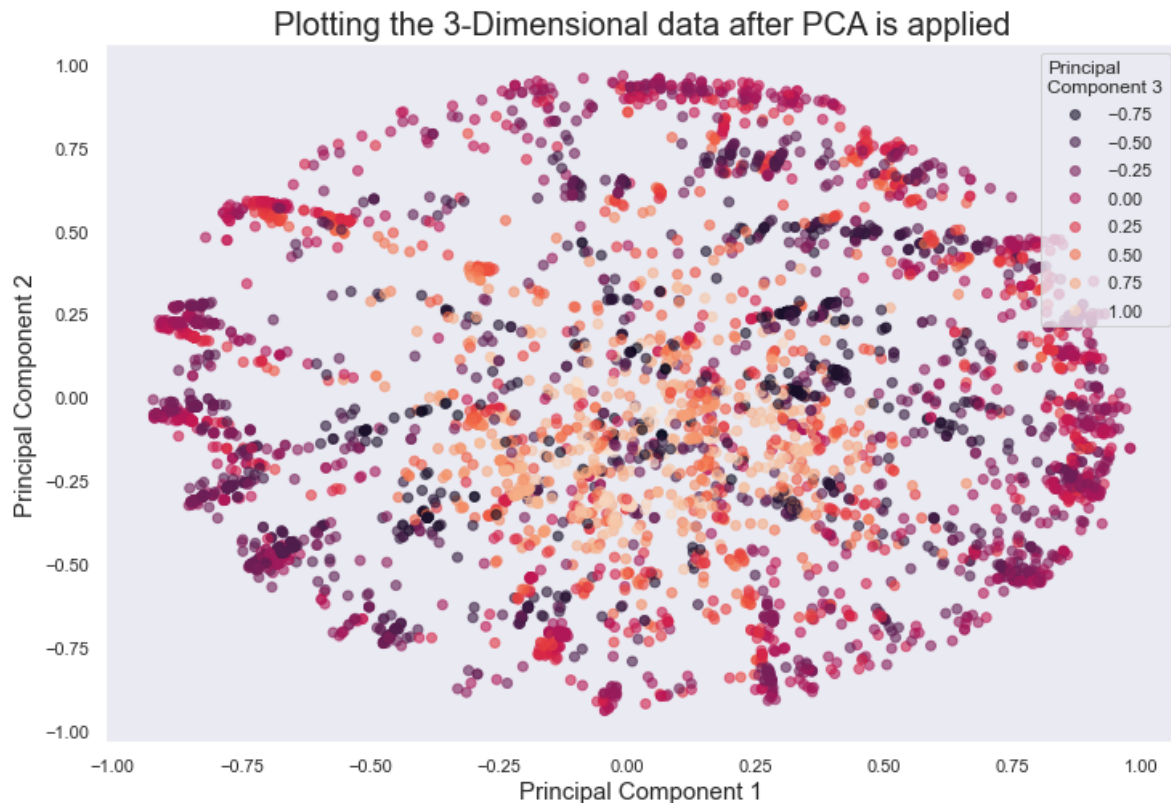
```
1 # Reducing the dimensions of the data
2 from sklearn.decomposition import PCA
3 pca = PCA(n_components = 3)
4 pcadf = pca.fit_transform(normalized_df)
5 pcadf = pd.DataFrame(pcadf)
6 pcadf.columns = ['Principal Component 1', 'Principal Component 2', 'Principal Component 3']
7
8 pcadf.head(10)
```

Out[13]:

	Principal Component 1	Principal Component 2	Principal Component 3
0	0.375353	-0.061732	0.885181
1	-0.285591	-0.137249	0.815803
2	-0.244175	-0.140513	0.851572
3	-0.840738	-0.028420	-0.292160
4	-0.872485	-0.068228	0.046006
5	0.373151	0.525823	-0.481568
6	-0.840738	-0.028420	-0.292160
7	0.297989	-0.033309	0.666742
8	0.820311	0.211774	-0.202186
9	0.871196	0.125964	-0.024124

In [14]:

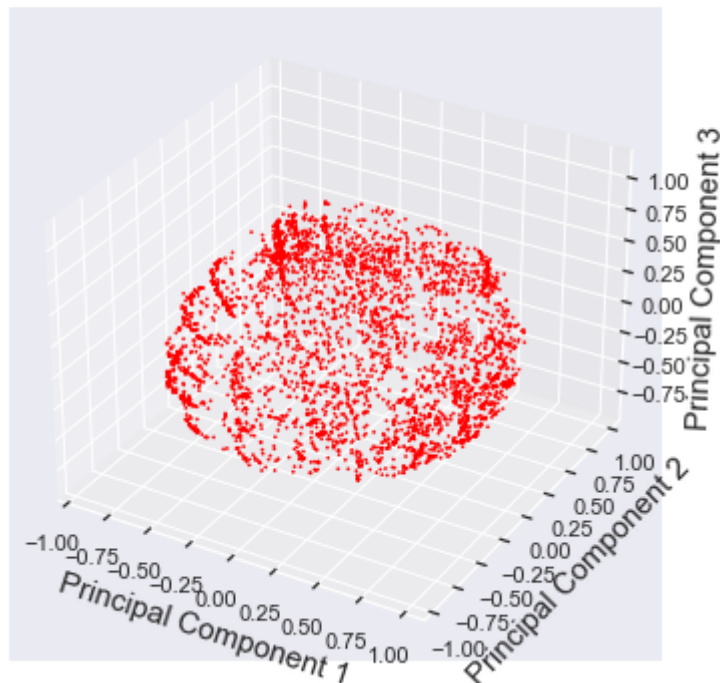
```
1 fig, ax = plt.subplots(figsize=(12, 8))
2 scatter = ax.scatter(pcadf['Principal Component 1'],
3                      pcadf['Principal Component 2'],
4                      c = pcadf['Principal Component 3'],
5                      alpha=0.6)
6 plt.title('Plotting the 3-Dimensional data after PCA is applied', fontsize = 20)
7 plt.xlabel('Principal Component 1', fontsize = 15)
8 plt.ylabel('Principal Component 2', fontsize = 15)
9 plt.legend(*scatter.legend_elements(), loc="best", title="Principal\nComponent 3")
10 ax.plot([])
11 ax.grid()
12 plt.show()
```



In [15]:

```
1 fig = plt.figure(figsize = (15, 6))
2 plt.suptitle("Plotting the 3 component with colors representing clusters\n'x axis : P1
3
4 x, s1 = pcadf['Principal Component 1'], "Principal Component 1"
5 y, s2 = pcadf['Principal Component 2'], "Principal Component 2"
6 z, s3 = pcadf['Principal Component 3'], "Principal Component 3"
7
8 ax = fig.add_subplot(111, projection = '3d')
9 ax.scatter(x, y, z, c = "red", s=0.5, alpha = 1)
10 ax.set_xlabel(s1, fontsize = 15)
11 ax.set_ylabel(s2, fontsize = 15)
12 ax.set_zlabel(s3, fontsize = 15)
13
14 plt.show()
```

Plotting the 3 component with colors representing clusters
'x axis : P1 | y axis : P2 | z axis : P3'



In [16]:

```

1 import plotly as py
2 import plotly.graph_objs as go
3 import numpy as np
4 import pandas as pd
5
6 trace = go.Scatter3d(
7     x=pcadf['Principal Component 1'],
8     y=pcadf['Principal Component 2'],
9     z=pcadf['Principal Component 3'],
10
11     mode='markers',
12     marker=dict(
13         size=5,
14         color=pcadf['Principal Component 3'],
15         colorscale='Viridis',
16     ),
17     name= 'PCA Visualization',
18
19     # List comprehension to add text on hover
20     text= [f"PCA 1: {a}<br>PCA 2: {b}<br>PCA 3: {c}" for a,b,c in list(zip(pcadf['Principal Component 1'], pcadf['Principal Component 2'], pcadf['Principal Component 3']))],
21     # if you do not want to display x,y,z
22     hoverinfo='text'
23 )
24
25
26 layout = dict(title = 'PCA Visualization',)
27
28 data = [trace]
29 fig = dict(data=data, layout=layout)
30
31 py.offline.plot(fig, filename = 'Test.html')

```

Out[16]:

'Test.html'

Since there are no distinguishable clusters, we go for DBSCAN clustering and Heirarchial Clustering algorithms

DBSCAN Clustering (World)

In [17]:

```

1 # df1 = df[['Country_embedding',
2 #           'Continent_embedding',
3 #           'Domain_embedding',
4 #           'Industry_embedding',
5 #           'FieldOfStudy_embedding',
6 #           'Latitude', 'Longitude']]

```

In [18]:

```

1 from sklearn.cluster import DBSCAN
2
3 db = DBSCAN(eps=0.3, min_samples=10).fit(df1.drop(['latitude', 'longitude'], axis = 1))
4 labels = db.labels_
5 print(labels, len(labels))

```

[0 1 1 ... 2 2 2] 4910

In [19]:

```

1 df1['label'] = labels
2 df1['label'].value_counts()

```

Out[19]:

```

2    4341
1     204
0     204
-1    118
3      17
4      16
5      10

```

Name: label, dtype: int64

In [20]:

```
1 print("Cluster types: ",set(labels))
```

Cluster types: {0, 1, 2, 3, 4, 5, -1}

In [21]:

```
1 df1.sample(5)
```

Out[21]:

	Country_embedding	Continent_embedding	Domain_embedding	Industry_embedding	Field
1161	-0.156677	-0.227953	-0.112251	0.394512	
206	-0.156677	-0.227953	0.195954	0.026966	
3847	0.812406	-0.227953	0.116207	0.389670	
2046	-0.156677	-0.227953	0.484690	-0.245237	
3864	-0.156677	-0.227953	0.116207	-0.245237	

Evaluation Metrics

Silhouette Score

In [22]:

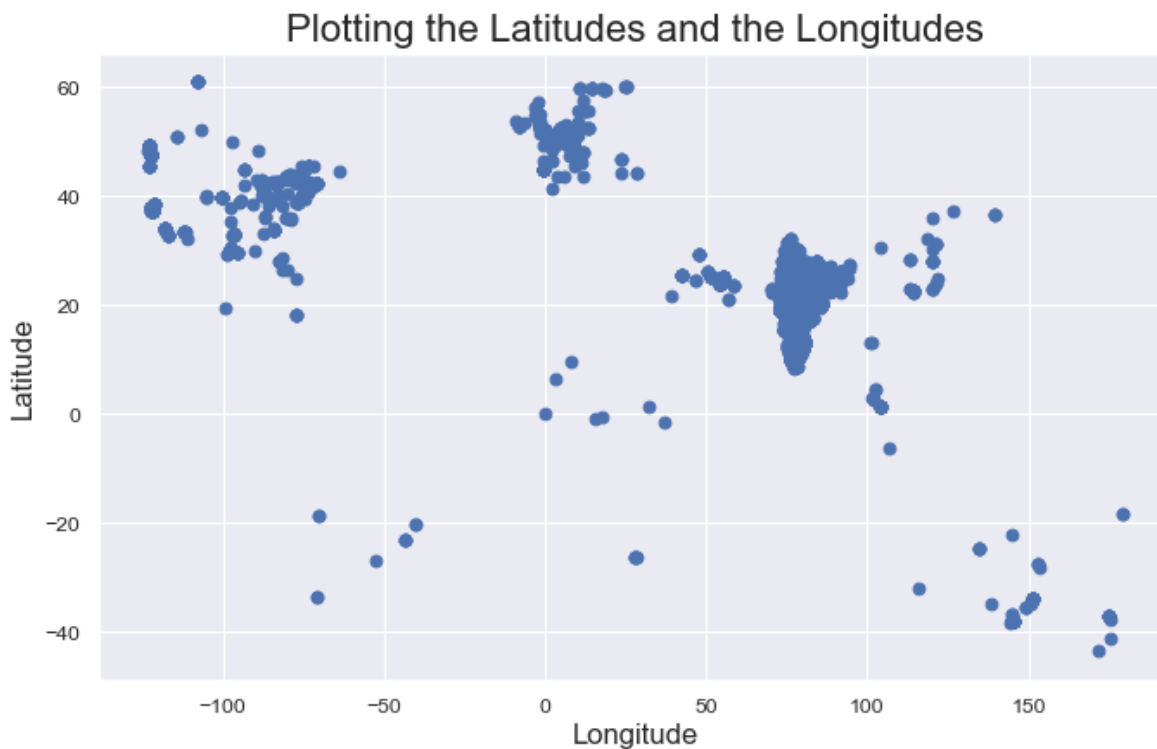
```
1 from sklearn.cluster import KMeans
2 from sklearn.metrics import silhouette_score
3
4 sil_coeff = silhouette_score(df1, labels, metric='euclidean')
5 print("Silhoutte Score: ", sil_coeff)
```

Silhoutte Score: 0.6633363585442195

Plotting the world-wide profile clusters

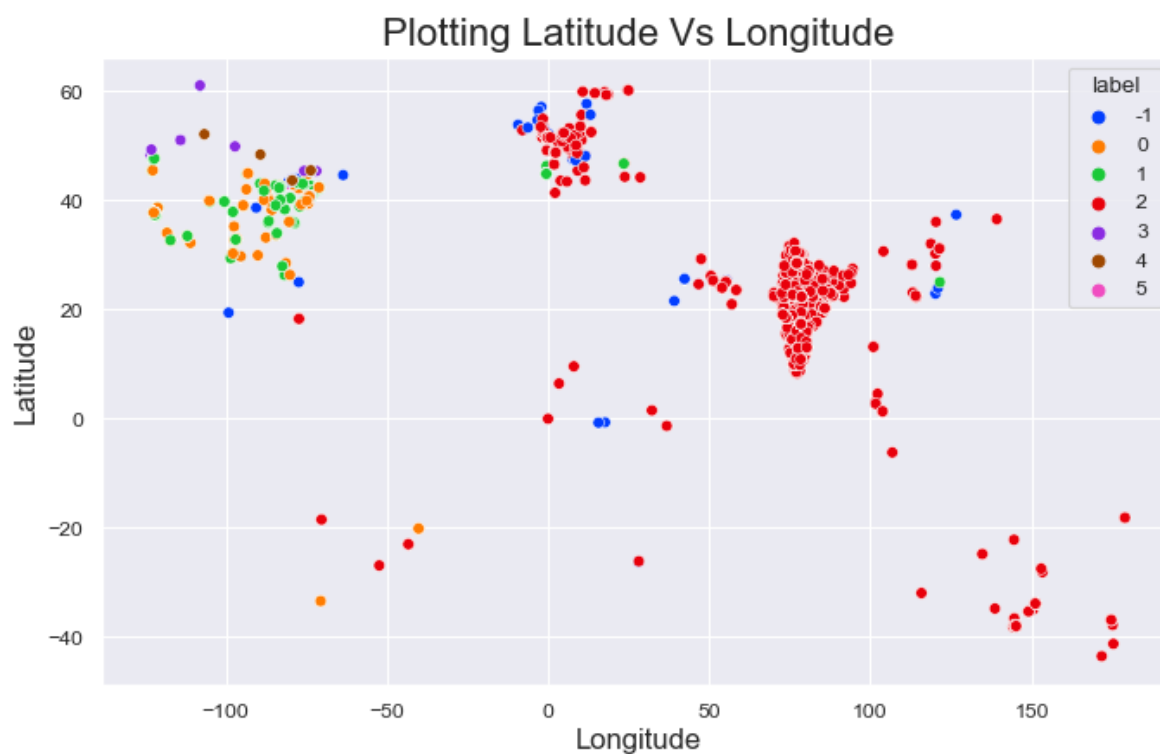
In [25]:

```
1 import matplotlib.pyplot as plt
2 from matplotlib.pyplot import figure
3
4 plt.figure(figsize=(10, 6), dpi=80)
5 plt.scatter(x=df1['longitude'], y=df1['latitude'])
6 plt.title('Plotting the Latitudes and the Longitudes', fontsize = 20)
7 plt.xlabel('Longitude', fontsize = 15)
8 plt.ylabel('Latitude', fontsize = 15)
9 plt.show()
```



In [26]:

```
1 #markers=['.',',','o','v','^','<','>','1','2','3','4','8']
2 import seaborn as sns
3 sns.set_theme(color_codes=True)
4 fig, ax = plt.subplots(figsize = (10, 6),dpi=80)
5 markers=['o','v','^','<','x','.',',']
6 sns.scatterplot(ax = ax , x = "longitude" , y = "latitude" , data = df1, hue = "label")
7 ax.set_xlabel( "Longitude" , size = 15)
8 ax.set_ylabel( "Latitude" , size = 15)
9 ax.set_title( "Plotting Latitude Vs Longitude" , size = 20)
10 plt.show()
```



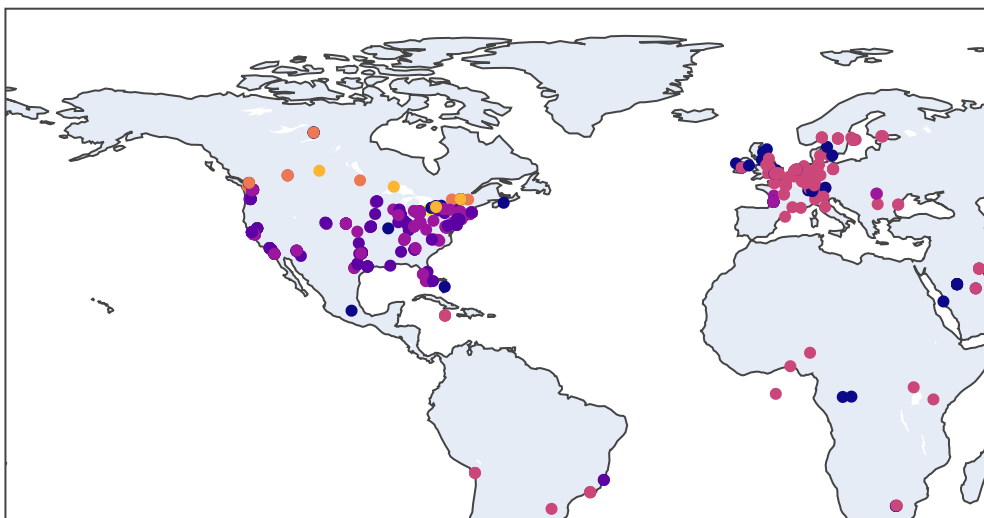
In [27]:

```

1 import plotly.io as pio; pio.renderers.default='notebook'
2 import plotly.express as px
3 import pandas as pd
4 fig = px.scatter_geo(df1,lat='latitude',lon='longitude', hover_name="label", color = d-
5 fig.update_layout(title = 'World map', title_x=0.5)
6 fig.show()

```

World map



In [52]:

```

1 import pandas as pd
2 from shapely.geometry import Point
3 import geopandas as gpd
4 from geopandas import GeoDataFrame
5
6 geometry = [Point(xy) for xy in zip(india_df['longitude'], df['latitude'])]
7 gdf = GeoDataFrame(india_df, geometry=geometry)
8
9 #this is a simple map that goes with geopandas
10 world = gpd.read_file(gpd.datasets.get_path('naturalearth_lowres'))
11 gdf.plot(ax=world.plot(figsize=(10, 6)), marker='o', color='red', markersize=15)

```

1

In []:

1

Word Clouds

In [29]:

```
1 df.columns
```

Out[29]:

```
Index(['Industry', 'Company Location', 'New Job (90 Days)', 'Year Started',
      'Profile Headline', 'Profile Summary', 'School', 'Degree',
      'Education End', 'Domain', 'CompanyName', 'JobTitle', 'My Network',
      'Country', 'Continent', 'FieldOfStudy', 'Industry_embedding',
      'Company Location_embedding', 'Profile Headline_embedding',
      'Profile Summary_embedding', 'School_embedding', 'Degree_embedding',
      'Education End_embedding', 'Domain_embedding', 'CompanyName_embeddin
g',
      'JobTitle_embedding', 'My Network_embedding', 'Country_embedding',
      'Continent_embedding', 'FieldOfStudy_embedding', 'latitude',
      'longitude'],
      dtype='object')
```

In [30]:

```
1 word_cloud_df = df[['Industry',
2                     'Degree',
3                     'Domain',
4                     'CompanyName',
5                     'JobTitle',
6                     'FieldOfStudy']]
```

In [31]:

```
1 word_cloud_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 4910 entries, 0 to 4909
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Industry        4910 non-null   object
1   Degree          4910 non-null   object
2   Domain          4910 non-null   object
3   CompanyName     4910 non-null   object
4   JobTitle        4910 non-null   object
5   FieldOfStudy    4910 non-null   object
dtypes: object(6)
memory usage: 268.5+ KB
```

In [32]:

```
1 def join_text(df):
2     k = ""
3     for col in df.columns:
4         try:
5             k = (' '.join(df[col].str.lower()))
6         except:
7             print(col)
8     return k
```

In [33]:

```
1 word_cloud_df['label'] = labels
2 word_cloud_df['label'].value_counts()
```

Out[33]:

```
2    4341
1     204
0     204
-1    118
3      17
4      16
5      10
```

Name: label, dtype: int64

In [34]:

```
1 clus1_df = word_cloud_df[word_cloud_df['label'] == -1]
```

In [35]:

```
1 def plot_wordCloud(df):
2     from textblob import TextBlob
3     from collections import Counter
4     from wordcloud import WordCloud, STOPWORDS
5     import matplotlib.pyplot as plt
6
7     k = join_text(df.drop(['label'], axis = 1))
8     wordcloud = WordCloud(width = 2000, height = 1000, prefer_horizontal = 0.3).generate(k)
9     plt.figure(figsize=(15,15))
10    plt.imshow(wordcloud)
11    plt.axis('off')
```


In [36]:

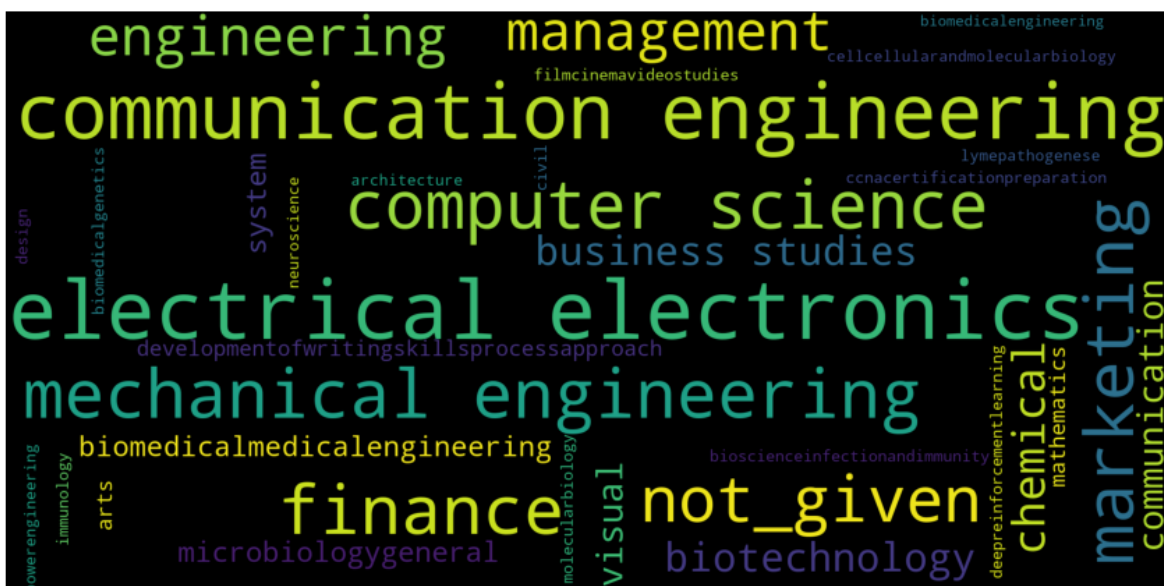
```

1 import collections
2 import numpy as np
3 import pandas as pd
4 import matplotlib.cm as cm
5 import matplotlib.pyplot as plt
6 from matplotlib import rcParams
7 from wordcloud import WordCloud, STOPWORDS
8 %matplotlib inline
9 stopwords = STOPWORDS
10 stopwords.add('&')
11 stopwords.add('-')
12
13 def analyse_word_cloud(df, num):
14     all_words = join_text(df)
15
16     filtered_words = [word for word in all_words.split() if word not in stopwords]
17     counted_words = collections.Counter(filtered_words)
18
19     words = []
20     counts = []
21     for letter, count in counted_words.most_common(10):
22         words.append(letter)
23         counts.append(count)
24
25     colors = cm.rainbow(np.linspace(0, 1, 10))
26     rcParams['figure.figsize'] = 20, 10
27     plt.title('Top words in the cluster (' + str(num) + ') vs their count', fontsize =
28     plt.xlabel('Count', fontsize = 20)
29     plt.ylabel('Words', fontsize = 20)
30     plt.xticks(fontsize = 20)
31     plt.yticks(fontsize = 20)
32     plt.barh(words, counts, color=colors)
33     plt.show()

```

In [37]:

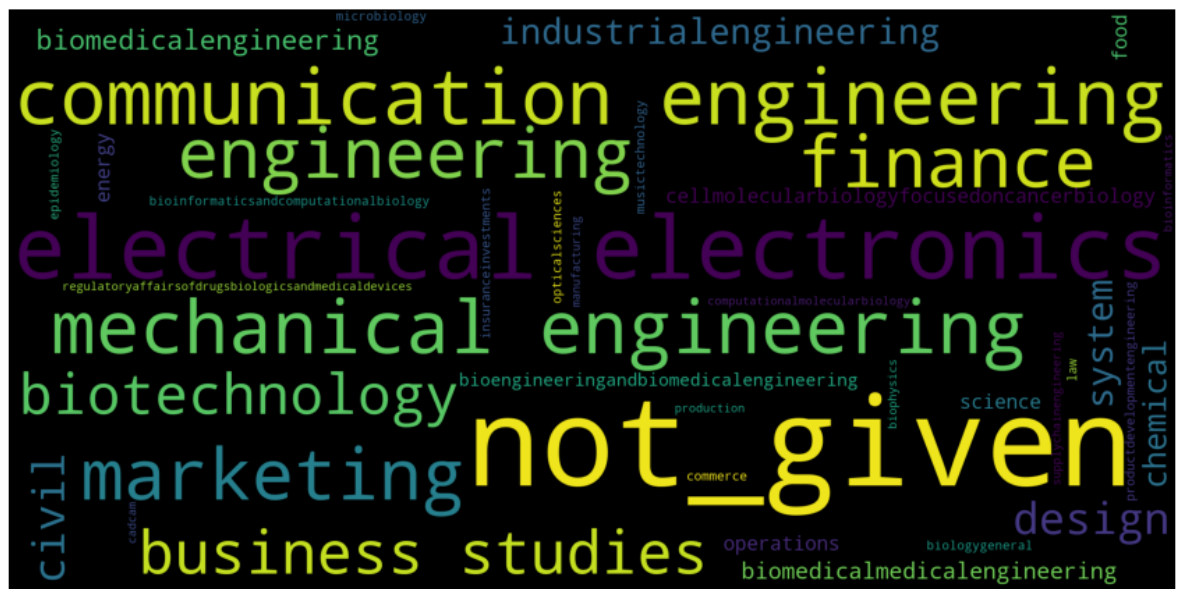
```
1 plot_wordCloud(clus1_df)
```



```
1 analyse_word_cloud(clus1_df, -1)
```

Words	Count
not_given	11
finance	12
marketing	12
science	12
computer	12
mechanical	12
electronics	16
electrical	16
communication	19
engineering	36

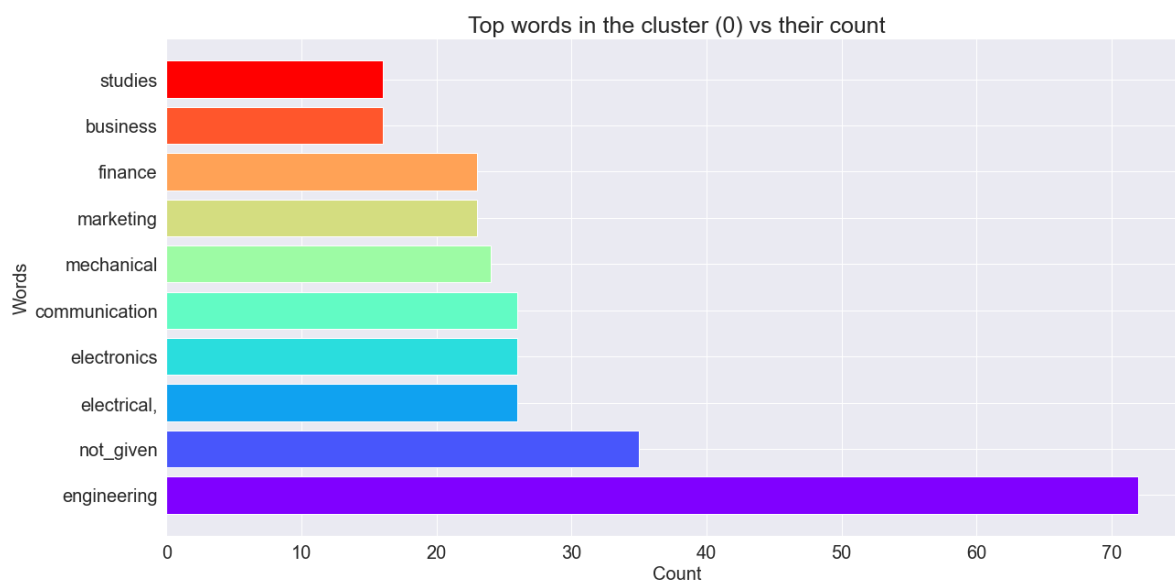
```
1 clus2_df = word_cloud_df[word_cloud_df['label'] == 0]
2 plot_wordCloud(clus2_df)
```



In [40]:

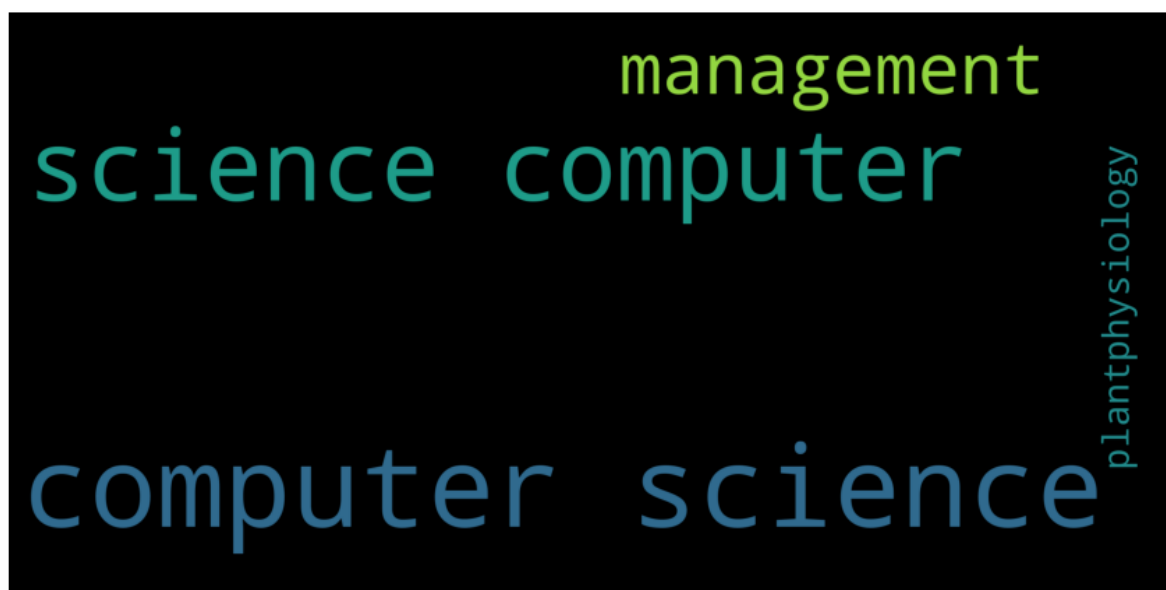
```
1 analyse_word_cloud(clus2_df, 0)
```

label



In [41]:

```
1 clus3_df = word_cloud_df[word_cloud_df['label'] == 1]  
2 plot_wordCloud(clus3_df)
```

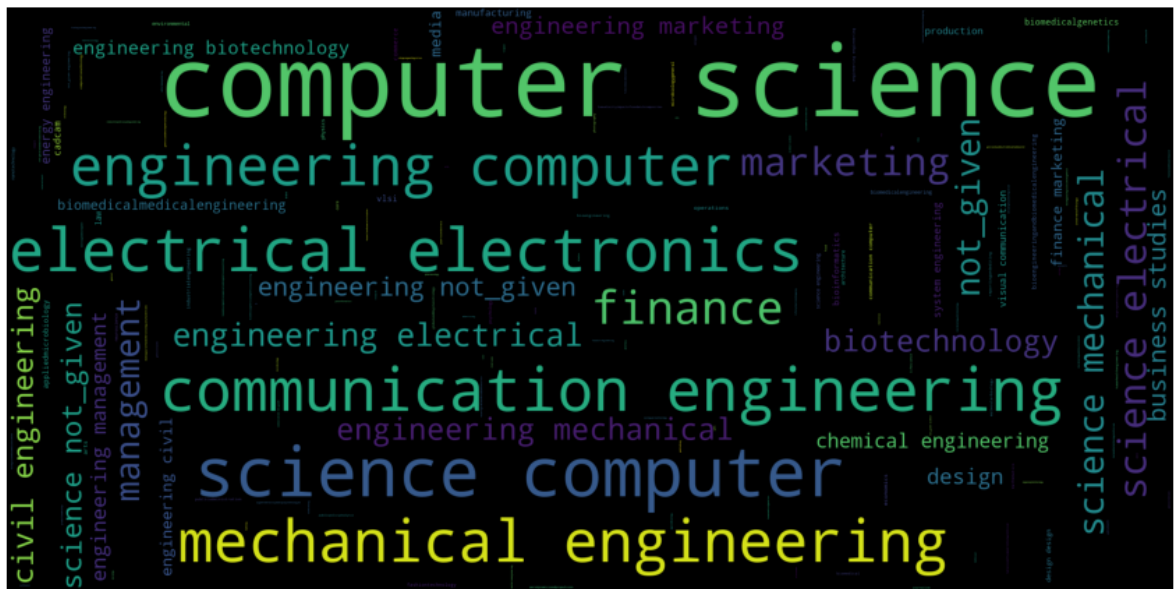


```
1 analyse_word_cloud(clus3_df, 1)
```

Top words in the cluster (1) vs their count

Words	Count
plantphysiology	2
management	50
science	150
computer	150

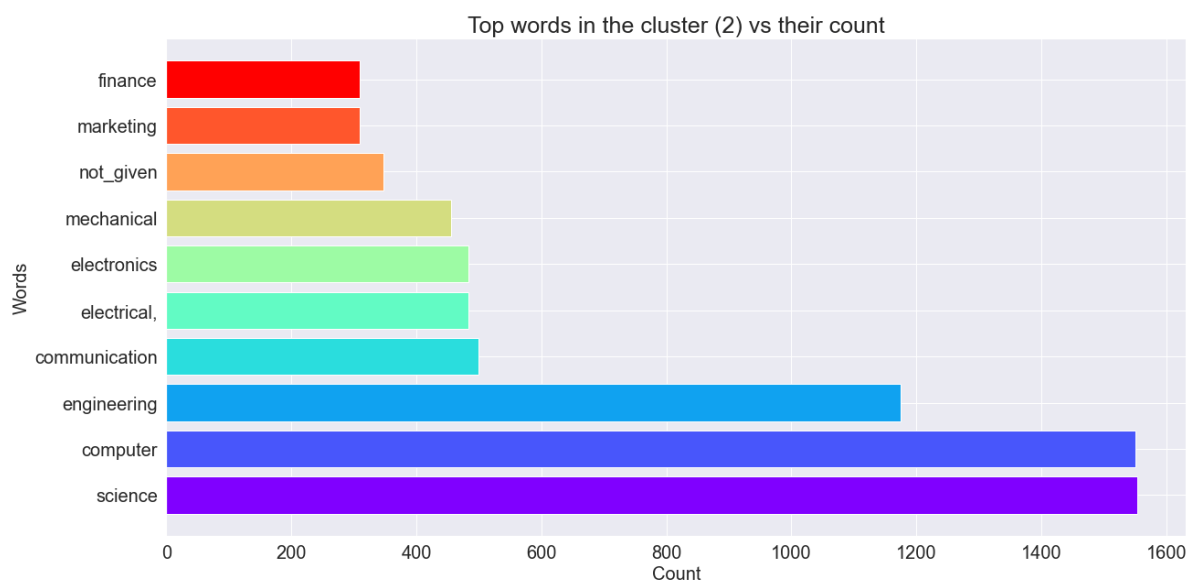
```
1 clus4_df = word_cloud_df[word_cloud_df['label'] == 2]
2 plot_wordCloud(clus4_df)
```



In [44]:

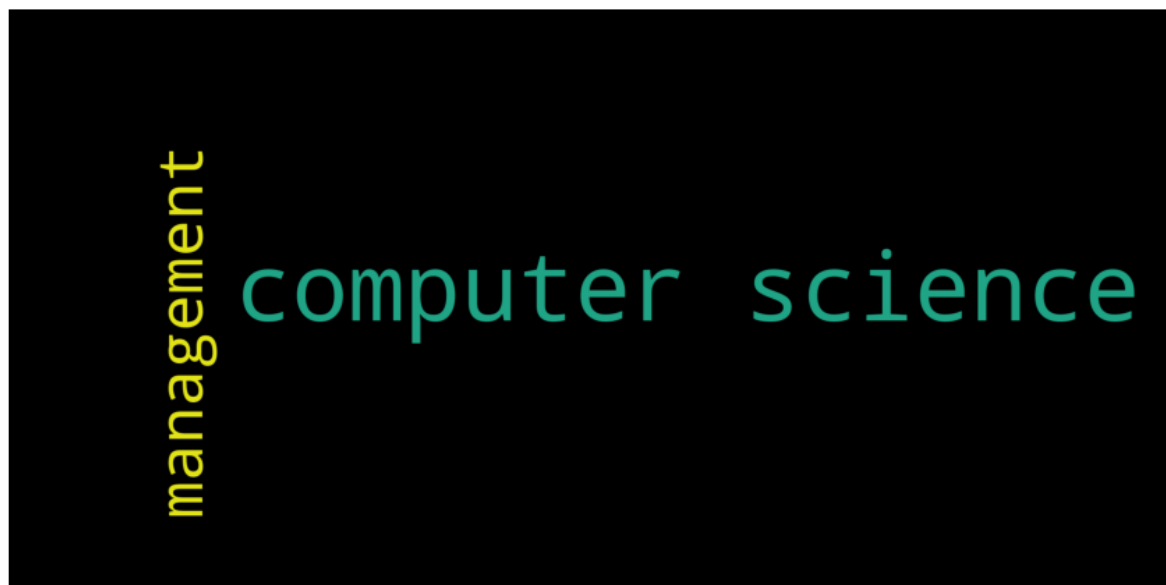
```
1 analyse_word_cloud(clus4_df, 2)
```

label



In [45]:

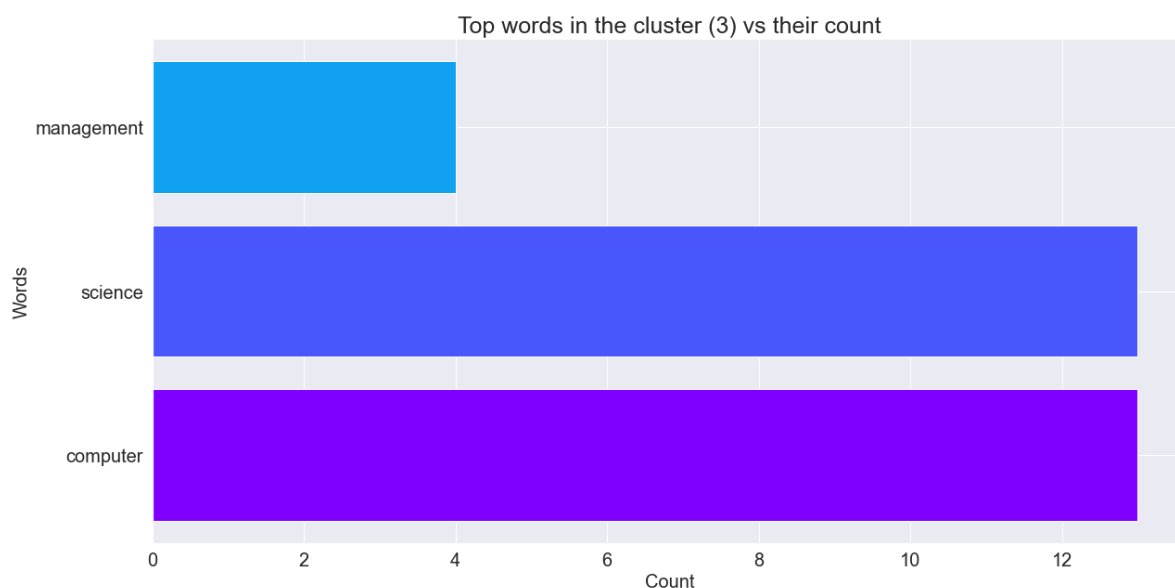
```
1 clus5_df = word_cloud_df[word_cloud_df['label'] == 3]  
2 plot_wordCloud(clus5_df)
```



In [46]:

```
1 analyse_word_cloud(clus5_df, 3)
```

label



In [47]:

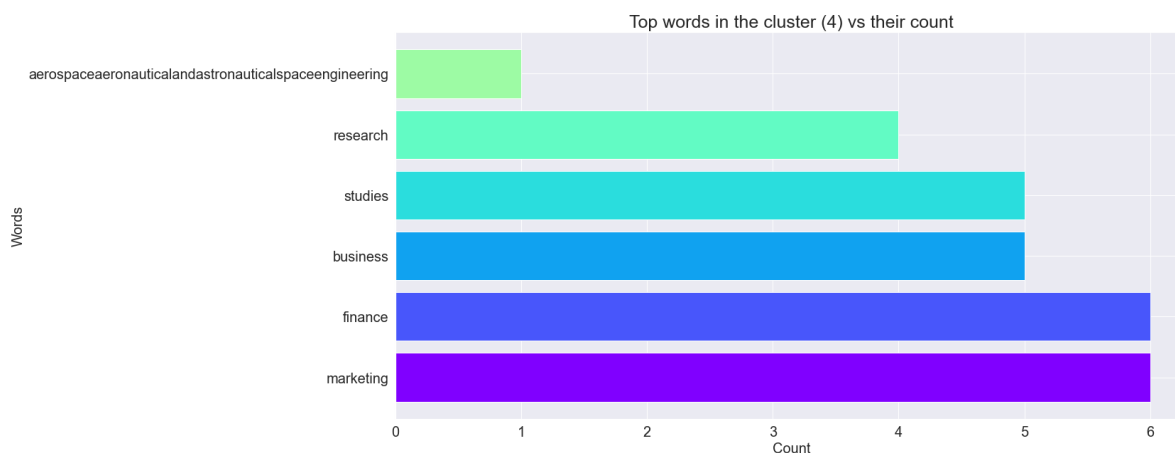
```
1 clus6_df = word_cloud_df[word_cloud_df['label'] == 4]  
2 plot_wordCloud(clus6_df)
```



In [48]:

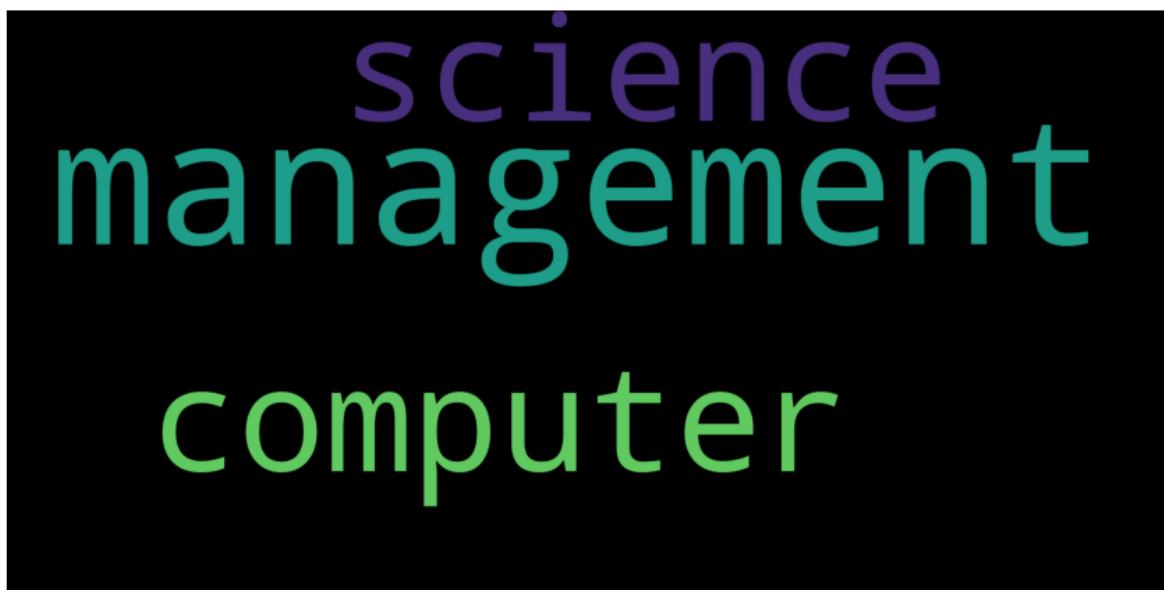
```
1 analyse_word_cloud(clus6_df, 4)
```

label



In [49]:

```
1 clus7_df = word_cloud_df[word_cloud_df['label'] == 5]  
2 plot_wordCloud(clus7_df)
```



In [50]:

```
1 analyse_word_cloud(clus7_df, 5)
```

label

