





```
[88]: 'analyze',
      'model',
      'optimize',
      'solve',
      'visualize',
      'extract',
      'build',
      'deploy',
      'preprocess',
      'train',
      'interpret',
      'wrangle',
      'collect',
      'consolidate',
      'data engineer',
      'design',
      'build',
      'teach',
      'assess',
      'execute',
      'integrate',
      'efficient',
      'advanced',
      'deep',
      'complex',
      'statistical',
      'predictive',
      'analytical',
      'visual',
      'real-time',
      'algorithmic',
      'technical',
      'machine',
      'artificial',
      'neural',
      'quantitative',
      'qualitative',
      'structured',
      'architecture',
      'data science',
      'data',
      'python',
      'powerbi',
      'python programming',
      'big data',
      'machine learning',
      'artificial intelligence',
      'data science',
      'data analytics',
      'data scientist',
      'data engineer',
      'data analyst',
      'data mining',
      'data visualization',
      'predictive modeling',
      'statistical modeling',
      'supervised learning',
      'unsupervised learning',
      'deep learning',
      'neural networks',
      'natural language processing',
      'reinforcement learning',
      'ensemble learning',
      'feature engineering',
      'model training',
      'model evaluation',
      'model deployment',
      'data preprocessing',
      'data cleaning',
      'data wrangling',
      'feature selection',
      'cross-validation',
      'hyperparameter tuning',
      'model interpretability',
      'bias-variance tradeoff',
      'accuracy',
      'precision',
      'recall',
      'f1-score',
      'confusion matrix',
      'ROC curve',
      'AUC',
      'gradient descent',
      'random forest',
      'decision tree',
      'support vector machine',
      'k-nearest neighbors',
      'clustering',
      'classification',
      'regression',
      'evaluation metrics',
      'model interpretation',
      'model explainability',
      'model performance',
      'model selection',
      'model optimization',
      'model deployment',
      'data-driven',
      'analytical',
      'problem-solving',
      'critical thinking',
      'detail-oriented',
      'communication',
      'collaboration',
      'teamwork',
      'leadership',
      'adaptability',
      'flexibility',
      'innovation',
      'strategic thinking',
      'business acumen',
      'presentation skills',
      'project management',
      'time management',
      'self-motivation',
      'multitasking',
      'organizational skills',
      'attention to detail',
      'technical skills',
      'analytical skills',
      'qualitative skills',
      'analytical',
      'problem-solving',
      'detail-oriented',
      'critical thinking',
      'logical',
      'creative',
      'flexible',
      'attention to detail',
      'mathematical',
      'statistical knowledge',
      'technical',
      'business acumen',
      'effective communication',
      'data-driven',
      'data privacy',
      'adaptable',
      'collaborative',
      'project management',
      'results-oriented',
      'self-motivated',
      'continuous learning',
      'data interpretation',
      'data analysis',
      'data visualization',
      'data modeling',
      'data mining',
      'data cleaning',
      'data wrangling',
      'data storytelling',
      'SQL querying',
      'python programming',
      'R programming',
      'machine learning',
      'predictive modeling',
      'statistical analysis',
      'reporting',
      'dashboarding',
      'data governance',
      'data privacy',
      'data security',
      'data ethics',
      'data quality assurance',
      'ETL (Extract, Transform, Load)',
      'data warehouse management',
      'database querying',
      'data integration',
      'data transformation',
      'data validation',
      'data manipulation',
      'data insights',
      'data strategy',
      'business intelligence',
      'visualization tools',
      'statistical techniques',
      'machine learning algorithms',
      'data-driven decision-making',
      'collaboration tools',
      'presentation skills',
      'client interaction',
      'problem identification',
      'solution development',
      'data governance frameworks',
      'quality assurance',
      'data architecture',
      'data documentation',
      'data security measures',
      'data privacy regulations',
      'data visualization techniques',
      'effective reporting',
      'business understanding',
      'domain-specific knowledge',
      'data analysis best practices',
      'continuous improvement',
      'strategic thinking',
      'innovative mindset',
      'stakeholder management',
      'risk assessment',
      'effective communication',
      'data storytelling expertise',
      'data presentation skills',
      'data management',
      'data governance policies',
      'data privacy frameworks',
      'data security protocols',
      'data quality management',
      'data integration methods',
      'data visualization strategies',
      'data validation techniques',
      'data manipulation skills',
      'data warehouse design',
      'SQL optimization',
      'programming proficiency',
      'analytical techniques',
      'algorithm implementation',
      'statistical hypothesis testing',
      'visualization platforms',
      'business intelligence tools',
      'data exploration methods',
      'data interpretation frameworks',
      'data reporting methodologies',
      'data collaboration platforms',
      'project management methodologies',
      'python',
      'R',
      'SQL',
      'Java',
      'Scala',
      'MATLAB',
      'Tableau',
      'Power BI',
      'C++',
      'Pandas',
      'NumPy',
      'SciPy',
      'scikit-learn',
      'TensorFlow',
      'Keras',
      'PyTorch',
      'NLTK',
      'Natural Language Processing',
      'OpenCV',
      'Apache Spark',
      'Hadoop',
      'AWS',
      'Azure',
      'Google Cloud',
      'Docker',
      'Git',
      'Jenkins',
      'Visual Studio Code',
      'Sublime Text',
      'Dataiku',
      'Figma',
      'Altair',
      'Plotly',
      'Redshift',
      'MySQL',
      'PostgreSQL',
      'MongoDB',
      'Cassandra',
      'Hive',
      'Apache Kafka',
      'Apache Airflow',
      'Databricks',
      'Dask',
      'Streamlit',
      'Flask',
      'Django',
      'Pylot',
      'Seaborn',
      'Matplotlib',
      'Bokeh',
      'Ggplot2',
      'Statsmodels',
      'Xgboost',
      'LightGBM',
      'Catboost',
      'SVM',
      'Random Forest',
      'Linear Regression',
      'Logistic Regression',
      'Reinforcement Learning',
      'Time Series Analysis',
      'Deep Learning',
      'Neural Networks',
      'Data Mining',
      'Data Wrangling',
      'Feature Engineering',
      'Model Evaluation',
      'Model Deployment',
      'Data Cleaning',
      'Data Visualization',
      'Exploratory Data Analysis',
      'Machine Learning',
      'Artificial Intelligence',
      'Predictive Analytics',
      'Statistical Analysis',
      'Text Mining',
      'Image Processing',
      'Web Scraping',
      'Version Control',
      'Collaborative Filtering',
      'Recommendation Systems',
      'Knowledge Discovery',
      'Dimensionality Reduction',
      'Ensemble Learning',
      'Model Selection',
      'Cross-Validation',
      'Hyperparameter Tuning',
      'Model Interpretability',
      'Model Metrics',
      'Feature Selection',
      'Data Engineering',
      'Model Deployment',
      'Cloud Computing',
      'Containerization',
      'API Development',
      'Web Development',
      'Time Series Forecasting',
      'Unsupervised Learning',
      'Supervised Learning',
      'Natural Language Understanding',
      'Computer Vision',
      'Parallel Computing',
      'Data Privacy',
      'Data Governance',
      'Data Quality',
      'Data Security',
      'Data Ethics',
      'Data Usability',
      'Data Integration',
      'Data Architecture',
      'ETL',
      'Big Data',
      'Internet of Things',
      'Blockchain',
      'Cybersecurity']

In [89]: # Remove punctuation and convert text to lowercase
        Indf['job description'] = Indf['job description'].str.replace(r'[^\w\s\n]', ''), str.lower()

C:\Users\juydydutta\AppData\Local\Temp\ipykernel_1156\3138801493.py:2: FutureWarning: The default value of regex
x will change from True to False in a future version.
  Indf['job description'] = Indf['job description'].str.replace(r'[^\w\s\n]', ''), str.lower()

In [90]: # Remove symbols from the text
        Indf['job description'] = Indf['job description'].str.replace(r'[@%&*()\[\]{}', '')

# Print the updated job descriptions
Indf['job description']

C:\Users\juydydutta\AppData\Local\Temp\ipykernel_1156\3202076043.py:2: FutureWarning: The default value of regex
x will change from True to False in a future version.
  Indf['job description'] = Indf['job description'].str.replace(r'[@%&*()\[\]{}', '')

Out[90]: 0      about the job\this role has been designated a...
1      about the job\about us\braintrust is a use...
2      about the job\about us\braintrust is a use...
3      about the job\american express invites you to...
4      about the job\job descriptions promanage it s...
...
95     about the job\jponentester enterprise is a co...
96     about the job\greetings\contractual hiring, ernst, y...
97     about the job\ance you passionate about music...
98     about the job\job description\american express...
99     about the job\quadeye is one of the largest i...
Name: job description, Length: 100, dtype: object

In [91]: # Tokenize the text
        Indf['job description'] = Indf['job description'].apply(word_tokenize)
        Indf['job description']

Out[91]: 0      [about, the, job, this, role, has, been, desig...
1      [about, the, job, about, us, braintrust, is, a...
2      [about, the, job, at, snapshot, we, believe, t...
3      [about, the, job, american, express, invites...
4      [about, the, job, job, descriptions, promanage...
...
95     [about, the, job, jponentester, enterprise, in...
96     [about, the, job, greetings, contractual, hire...
97     [about, the, job, ance, you, passionate, abou...
98     [about, the, job, job, description, excellen...
99     [about, the, job, quadeye, is, one, of, the, l...
Name: job description, Length: 100, dtype: object

In [92]: # Removing stopwords
        stopwords_list = stopwords.words('english')
        Indf['job description'] = Indf['job description'].apply(lambda tokens: [word for word in tokens if word not in
        Indf['job description']

Out[92]: 0      [job, role, designated, sedgea, means, primari...
1      [job, us, braintrust, userowned, talent, netwo...
2      [job, snapshot, believe, future, work, remote...
3      [job, american, express, invites, share, resume...
4      [job, job, descriptions, promanage, solutions...
...
95     [job, jponentester, enterprise, consultation...
96     [job, greetings, contractual, hiring, ernst, y...
97     [job, passionate, music, experience, working...
98     [job, job, description, excellent, spoken, wri...
99     [job, quadeye, one, largest, indiabased, prop...
Name: job description, Length: 100, dtype: object

In [93]: # Lemmatizing the text
        lemmatizer = WordNetLemmatizer()
        Indf['job description'] = Indf['job description'].apply(lambda tokens: [lemmatizer.lemmatize(word) for word in
        Indf['job description']

Out[93]: 0      [job, role, designated, sedgea, mean, primari...
1      [job, us, braintrust, userowned, talent, netwo...
2      [job, snapshot, believe, future, work, remote...
3      [job, american, express, invite, share, resume...
4      [job, job, description, promanage, solution...
...
95     [job, jponentester, enterprise, consultation...
96     [job, greeting, contractual, hiring, ernst, y...
97     [job, passionate, music, experience, working...
98     [job, job, description, excellent, spoken, wri...
99     [job, quadeye, one, largest, indiabased, prop...
Name: job description, Length: 100, dtype: object

In [94]: # Filter relevant words related to data science
        Indf['job description'] = Indf['job description'].apply(lambda tokens: [word for word in tokens if word in dat
        Indf['job description']

Out[94]: 0      [analyze, data, complex, extract, data, data, ...
1      [algorithmic, model, model, data, data, statis...
2      [predictive, model, data, statistical, data, m...
3      [data, leadership, design, collaboration, anal...
4      [technical, technical, optimize, optimize, opt...
...
95     [predictive, model, data, statistical, data, m...
96     [communication, data, data, data, accuracy, an...
97     [optimize, analytical, data, communication, bu...
98     [communication, data, data, data, communication]
99     [track, quantitative, explore, quantitative, s...
Name: job description, Length: 100, dtype: object

In [95]: # Checking the final dataset
        Indf['job description']

Out[95]: 0      [analyze, data, complex, extract, data, data, ...
1      [algorithmic, model, model, data, data, statis...
2      [predictive, model, data, statistical, data, m...
3      [data, leadership, design, collaboration, anal...
4      [technical, technical, optimize, optimize, opt...
...
95     [predictive, model, data, statistical, data, m...
96     [communication, data, data, data, accuracy, an...
97     [optimize, analytical, data, communication, bu...
98     [communication, data, data, data, communication]
99     [track, quantitative, explore, quantitative, s...
Name: job description, Length: 100, dtype: object

In [96]: # Convert the column to string type
        Indf['job description'] = Indf['job description'].astype(str)

# Count distinct comma-separated words
word_counts = Indf['job description'].str.split(',').explode().str.strip().value_counts()

# Display the word counts
print(word_counts)

'data' 672
'communication' 69
'model' 67
'machine' 57
'analytical' 52
'collect' 1
'identify' 1
'algorithmic' 1
'design' 1
Name: job description, Length: 99, dtype: int64

In [97]: # Convert the column to string type
        Indf['job description'] = Indf['job description'].astype(str)

# Count distinct comma-separated words
word_counts = Indf['job description'].str.split(',').explode().str.strip().value_counts()

# Separate the words and counts into variables
words = word_counts.index.tolist()
counts = word_counts.values.tolist()

# Display the words and their counts
print(words, counts)

'data' 672
'communication' 69
'model' 67
'machine' 57
'analytical' 52
'collect' 1
'identify' 1
'algorithmic' 1
'design' 1
Name: job description, dtype: int64

In [98]: # Using the top keywords
        word_counts.head()

Out[98]: 'data' 672
        'communication' 69
        'model' 67
        'machine' 57
        'analytical' 52
        -- -- --
94     'collect' 1
95     'identify' 1
96     'collaboration' 1
97     'algorithmic' 1
98     'design' 1

99 rows x 2 columns

In [99]: # Create a DataFrame with words and counts
        word_counts_df = pd.DataFrame({'Word': words, 'Count': counts})

# Display the DataFrame
print(word_counts_df)

          Word  Count
0         'data'     672
1  'communication'     69
2         'model'     67
3         'machine'     57
4  'analytical'     52
...
94     'collect'         1
95     'identify'         1
96  'collaboration'         1
97  'algorithmic'         1
98     'design'         1

99 rows x 2 columns

In [100]: # Remove 'r' and 'n' from the dataset
        word_counts_df['Word'] = word_counts_df['Word'].str.replace(r'[\r\n]', '', regex=True)

# Print the updated dataset
print(word_counts_df['Word'])

0         'data'
1  'communication'
2         'model'
3         'machine'
4  'analytical'
...
94     'collect'
95     'identify'
96  'collaboration'
97  'algorithmic'
98     'design'
Name: Word, Length: 99, dtype: object

In [101]: # Dropping the outlier word data
        word_counts = word_counts_df[word_counts_df['Word'] != 'data']

Out[101]:          Word  Count
1  'communication'     69
2         'model'     67
3         'machine'     57
4  'analytical'     52
-- -- --
94     'collect'         1
95     'identify'         1
96  'collaboration'         1
97  'algorithmic'         1
98     'design'         1

95 rows x 1 columns

In [102]: # Set the range of clusters to try
        max_clusters = 3

# Initialize an empty list to store the within-cluster sum of squares (WCSS)
wcss = []

# Fit K-Means clustering for different number of clusters
for n_clusters in range(1, max_clusters+1):
    kmeans = KMeans(n_clusters=n_clusters)
    kmeans.fit(X)
    wcss.append(kmeans.inertia_)

C:\Users\juydydutta\Anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to
o have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it b
y setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(
C:\Users\juydydutta\Anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to
o have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it b
y setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(
C:\Users\juydydutta\Anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to
o have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it b
y setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(

In [104]: # Plot the Elbow Curve
        plt.plot(range(1, max_clusters+1), wcss)
        plt.xlabel('Number of Clusters')
        plt.ylabel('WCSS')
        plt.title('Elbow Curve')

# Find the optimal number of clusters
diff = [wcss[i] - wcss[i-1] for i in range(1, len(wcss))]
optimal_clusters = diff.index(max(diff)) + 1

# Add a vertical line at the optimal number of clusters
plt.axvline(x=optimal_clusters, color='r', linestyle='--', label='Optimal Clusters')
plt.legend()
plt.show()

print(f'Optimal number of clusters:', optimal_clusters)

Elbow Curve



Optimal number of clusters: 2



In [105]: # Assuming you have already determined the optimal number of clusters
        n_clusters = 2

# Fit K-Means clustering with the optimal number of clusters
kmeans = KMeans(n_clusters=n_clusters)
kmeans.fit(X)

# Calculate the Silhouette Coefficient
silhouette_avg = silhouette_score(X, kmeans.labels_)
print('Silhouette Coefficient:', silhouette_avg)

# Calculate the Calinski-Harabasz Index
ch_score = calinski_harabasz_score(X, kmeans.labels_)
print('Calinski-Harabasz Index:', ch_score)

C:\Users\juydydutta\Anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to
o have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it b
y setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(
C:\Users\juydydutta\Anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to
o have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it b
y setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(
C:\Users\juydydutta\Anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to
o have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it b
y setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(

In [106]: # Perform clustering on the dataset using KMeans
        kmeans = KMeans(n_clusters=2, random_state=42)
        kmeans.fit(X)
        cluster_labels = kmeans.fit_predict(X)

C:\Users\juydydutta\Anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to
o have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it b
y setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(

In [107]: # Add the cluster labels to the dataset
        word_counts_df['cluster_label'] = cluster_labels

In [108]: # Seeing the new dataset with labels
        word_counts_df.head()

Out[108]:          Word  Count  cluster_label
1  'communication'     69                1
2         'model'     67                1
3         'machine'     57                1
4  'analytical'     52                1
5  'technical'     48                1
-- -- --
94     'collect'         1
95     'identify'         1
96  'collaboration'         1
97  'algorithmic'         1
98     'design'         1

95 rows x 3 columns

In [113]: # Create a word cloud for each cluster
        fig = plt.figure(figsize=(8, 6))
        for cluster in range(2):
            # Check if there are words available for the cluster
            if len(cluster_words) > 0:
                cluster_wordcloud = WordCloud().generate(' '.join(cluster_words))

                # Display the word cloud
                plt.figure(figsize=(8, 6))
                plt.imshow(cluster_wordcloud, interpolation='bilinear')
                plt.title(f'Cluster: {cluster}')
                plt.axis('off')
                plt.show()

            print(f'No words available for Cluster {cluster}')

Cluster 0



Cluster 1


```



```
In [114]: #seeing words in first cluster label
filtered_df_1 = word_counts_df[word_counts_df['cluster_label'] == 1]
filtered_words = list(filtered_df_1['word'])
print(filtered_words)

["communication", "model", "machine", "analytical", "technical", "build", "identify", "complex",
 "advanced", "statistical", "design"]

In [115]: #seeing words in second cluster
filtered_df_2 = word_counts_df[word_counts_df['cluster_label'] == 0]
print(filtered_words)

filtered_words = list(filtered_df_2['word'])
print(filtered_words)

["leadership", "analyze", "tableau", "deep", "quantitative", "collaboration", "predictive", "innovation",
 "communication", "accuracy", "optimize", "solve", "execute", "track", "extract", "evaluate",
 "efficient", "build", "powerbi", "communication", "technical", "explore", "model", "teamwork",
 "structured", "advanced", "innovation", "creativity", "artificial", "innovation", "clustering",
 "model", "deploy", "technical", "classification", "regression", "track", "design", "analytical",
 "machine", "advanced", "visual", "complex", "algorithmic", "identify", "predictive",
 "leadership", "interpret", "optimize", "deep", "build", "unstructured", "complex", "track",
 "qualitative", "machine", "artificial", "efficient", "predictive", "neutral", "accuracy", "tableau",
 "adaptability", "solve", "quantitative", "deep", "artificial", "analyze", "evaluate",
 "communication", "accuracy", "extract", "predict", "solve", "explore", "analyze", "analytical",
 "consolidate", "collect", "identify", "collaboration", "algorithmic", "design"]

In [112]: # Save the clustering model
with open('clustering_model.pkl', 'wb') as f:
    pickle.dump(kmeans, f)

In [ ]:
```