

Code Description

The Python code performs K-Means clustering on an EEG dataset. Here's a step-by-step breakdown:

1. **Import Libraries:** It begins by importing essential libraries: pandas for data manipulation, NumPy for numerical operations, matplotlib and seaborn for plotting, StandardScaler for data normalization, KMeans for clustering, and PCA for dimensionality reduction.
2. **Load Data:** The code loads the EEG dataset from a CSV file ("s00.csv") into a pandas DataFrame.
3. **Rename Columns:** It renames the columns of the DataFrame to generic names ("Feature_1", "Feature_2", etc.) to avoid potential issues with column names that might be numerical or problematic for processing.
4. **Data Preprocessing:**
 - It removes the last column, assuming it represents the labels (for unsupervised learning, we don't use the labels).
 - It scales the data using StandardScaler to normalize the features. This is crucial for K-Means as it's distance-based, and features with larger scales can disproportionately influence the clustering.
5. **Elbow Method:**
 - It calculates the inertia (within-cluster sum of squares) for different numbers of clusters (k) ranging from 1 to 10.
 - It plots the elbow curve (inertia vs. number of clusters) to help determine the optimal number of clusters.
6. **K-Means Clustering:**
 - Based on the elbow plot (in this case, the code assumes 3 clusters), it trains a K-Means model with the selected optimal 'k'.
 - It assigns cluster labels to the original DataFrame in a new 'Cluster' column.
7. **Dimensionality Reduction with PCA:**
 - It applies Principal Component Analysis (PCA) to reduce the dimensionality of the data to 2 principal components. This is done for visualization purposes, as it's easier to plot clusters in a 2D space.
 - It adds the two principal components ('PCA1' and 'PCA2') as new columns to the DataFrame.
8. **Cluster Plot:**
 - It generates a scatter plot to visualize the clusters in the 2D PCA space. The x-axis represents the first principal component (PCA1), the y-axis represents the second principal component (PCA2), and the color of each point indicates its assigned cluster.
9. **Cluster Distribution:**
 - Finally, it prints the distribution of data points across the different clusters.

Analysis of the Plots and Output

- **Elbow Method Plot:** The elbow plot shows a rapid decrease in inertia as 'k' increases from 1 to 3. After k=3, the rate of decrease slows down significantly, forming an "elbow" in the curve. This suggests that 3 is a good choice for the number of clusters, as adding more clusters beyond 3 provides diminishing returns in terms of reducing inertia.
- **Cluster Plot:** The scatter plot visually represents the K-Means clustering in the 2D PCA space.
 - The data points are clearly separated into three distinct clusters, indicated by different colors (purple, yellow, and teal).
 - The clusters appear reasonably well-defined with minimal overlap, suggesting that K-Means effectively partitioned the data.
 - PCA has successfully captured the variance in the data to allow for good visual separation of the clusters.
- **Cluster Distribution Output:** The cluster distribution shows the number of data points assigned to each cluster:
 - Cluster
 - 2 14971
 - 1 8059
 - 0 7969
 - Name: count, dtype: int64

Cluster 2 has a significantly larger number of data points compared to clusters 0 and 1. This indicates a potential imbalance in the cluster sizes.

Optimal Result and Comprehensive Report

Based on the elbow method plot and the cluster visualization, the optimal number of clusters chosen (k=3) appears to be a reasonable choice. The elbow plot suggests that 3 clusters capture a significant amount of variance in the data, and the scatter plot visually confirms that the data is well-separated into three groups.

Comprehensive Report on EEG Signal Clustering

This analysis aimed to identify distinct patterns within EEG data using K-Means clustering. The process involved loading, preprocessing, and clustering the EEG dataset, followed by dimensionality reduction for visualization and analysis of cluster distribution.

The elbow method was employed to determine the optimal number of clusters. The resulting plot indicated that three clusters provided a good balance between minimizing inertia and avoiding over-segmentation. K-Means clustering was then performed with k=3, and the data was projected onto a 2D space using PCA for visualization.

The cluster plot effectively illustrates the separation of EEG data into three distinct groups. This visualization confirms that K-Means successfully partitioned the data based on underlying

patterns. However, the cluster distribution reveals an imbalance, with one cluster containing a disproportionately larger number of data points.

Potential Implications and Further Considerations:

- **EEG Pattern Analysis:** The identified clusters may represent different brain states or activities captured by the EEG signals. Further analysis could involve correlating these clusters with known cognitive processes or experimental conditions.
- **Data Imbalance:** The observed cluster imbalance might warrant further investigation. Techniques like oversampling or undersampling could be explored to mitigate the imbalance and potentially improve the clustering results.
- **Feature Importance:** Analyzing the contribution of original features to the principal components could provide insights into which EEG features are most important for distinguishing the clusters.
- **Alternative Clustering Methods:** Exploring other clustering algorithms, such as hierarchical clustering or DBSCAN, could provide alternative perspectives on the data structure.

In conclusion, the K-Means clustering analysis, guided by the elbow method and visualized through PCA, effectively revealed three distinct clusters within the EEG data. The cluster distribution highlighted a potential imbalance that could be addressed in future analyses. Overall, this analysis provides a foundation for further exploration of EEG patterns and their potential relevance to brain function.