# ELECTION RESULTS ANALYSIS

A PROJECT REPORT

*Submitted by*

## AAKARSH SHARMA (RA2111003011483)

## KARAN GANGWANI (RA2111003011515)

## BHAVIKA RUSTAGI (RA2111003011556)

*Under the Guidance of*

## DR. JEBAKUMAR R

Associate Professor, Department of Data Science and Business Systems

*In partial fulfilment of the requirements for the degree of*

## BACHELOR OF TECHNOLOGY in COMPUTER SCIENCE AND ENGINEERING



# SCHOOL OF COMPUTING

# COLLEGE OF ENGINEERING AMD TECHNOLOGY

# SRM INSTITUTE OF SCIENCE AND TECHNOLOGY

(under section 3 of UGC Act,1956)

**S.R.M NAGAR, KATTANKULATHUR-603203 CHENGALPATTU**

**DISTRICT**

**JULY2024**

# BONAFIDE CERTIFICATE

Certified that Mini project report titled Election result analysis is the bonafide work of **Aakarsh Sharma (RA2111003011483), Karan Gangwani (RA2111003011515) and Bhavika Rustagi (RA2111003011556)** who carried out the minor project under my supervision. Certified further, that to the best of my knowledge, the work reported herein does not form any other project report or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

SIGNATURE                                          SIGNATURE

 Dr. Jebakumar R                                  Dr.M.Pushpalatha

 Associate Professor                              Head of Department

 Department of C.TECH                             Department of C.TECH

# TABLE OF CONTENTS

# ABSTRACT

The study focuses on the application of various clustering algorithms to analyze the Lok Sabha election results from the Election Commission of India. The primary objective is to identify patterns and insights that can help in understanding voter behavior, demographic influences, and regional political dynamics. The dataset comprises detailed election results, including vote counts, winning margins, and party affiliations, across multiple election cycles.

We implemented several clustering algorithms, including K-means, DBSCAN, and hierarchical clustering, to group constituencies based on voting patterns and demographic factors. The performance of these algorithms was evaluated using metrics such as silhouette scores, Davies-Bouldin index, and cluster coherence.

The analysis revealed distinct clusters of constituencies with similar voting behaviors and demographic profiles, providing a nuanced understanding of regional political trends. K-means clustering, with its simplicity and efficiency, provided clear and interpretable clusters, whereas DBSCAN was effective in identifying outliers and dense regions of data. Hierarchical clustering offered insights into the nested structure of constituency similarities.

Our findings highlight the importance of demographic factors such as literacy rate, urbanization, and socio-economic status in shaping electoral outcomes. The clustering results can aid political analysts, campaign strategists, and policymakers in formulating targeted strategies for future elections.

In conclusion, the application of clustering algorithms to Lok Sabha election results demonstrates the potential of machine learning techniques in political data analysis. Future work could involve the integration of additional data sources, such as social media sentiments and economic indicators, to enhance the robustness of the clustering models.

# Introduction

## What is Data Mining?

Data mining is the process of discovering patterns, trends, and insights from large datasets using various techniques from statistics, machine learning, and database systems. The goal of data mining is to extract valuable knowledge and actionable information from raw data, which can then be used for decision-making, prediction, and optimization in various domains.

Data mining involves several stages, including data preprocessing, where raw data is cleaned, transformed, and prepared for analysis; pattern discovery, where algorithms are applied to identify meaningful patterns and relationships within the data; and interpretation and evaluation, where the discovered patterns are interpreted in the context of the problem domain and evaluated for their usefulness and reliability.

## The Clustering Algorithm:

Clustering algorithms represent a fundamental aspect of unsupervised learning in data mining, serving as essential tools for grouping similar data points together based on their intrinsic characteristics. Unlike supervised learning, where algorithms are trained on labelled data with predefined classes, clustering algorithms operate on unlabelled data, aiming to uncover underlying structures and patterns without prior knowledge of class labels.

Clustering algorithms play a pivotal role in exploratory data analysis, pattern recognition, and data compression, enabling researchers and practitioners to gain insights into the inherent organization of datasets. By partitioning data into coherent clusters, these algorithms facilitate the identification of natural groupings and relationships, which can inform decision-making processes across various domains.

## History and further Development of Clustering:

One of the earliest clustering algorithms dates to the 1950s with the introduction of the k-means algorithm by Stuart Lloyd. Lloyd's work focused on signal quantization for pulse code modulation in telecommunications, where he proposed an iterative algorithm to partition a set of points into k clusters by minimizing the mean squared distance between each point and the centroid of its assigned cluster.

Another milestone in clustering algorithms came in the 1960s with the development of hierarchical clustering methods. This period saw the introduction of techniques such as single-linkage clustering and complete-linkage clustering, which form the basis of hierarchical clustering approaches still widely used today.

Throughout the following decades, clustering algorithms continued to evolve, with contributions from various fields such as machine learning, data mining, and artificial intelligence. Notable developments include density-based clustering algorithms like DBSCAN (Density-Based Spatial Clustering of Applications with Noise), introduced by Martin Ester et al. in 1996, which can identify clusters of arbitrary shapes and handle noise in the data.

# Fundamentals of Clustering:

Clustering is a fundamental technique in data mining and machine learning, aimed at grouping similar data points together based on their intrinsic characteristics. The primary goal of clustering is to partition a dataset into coherent subsets, or clusters, where objects within the same cluster are more similar to each other than to those in other clusters. Here are the key fundamentals of clustering:

**Unsupervised Learning**: Clustering is an unsupervised learning task, meaning it does not require labelled data. Instead, clustering algorithms autonomously discover patterns and structures within the data based solely on the input features.

**Similarity or Dissimilarity Measure**: Clustering algorithms rely on a similarity or dissimilarity measure to quantify the proximity between data points. Common measures include Euclidean distance, Manhattan distance, cosine similarity, and correlation coefficient. The choice of measure depends on the nature of the data and the clustering task.

**Objective Function**: Clustering algorithms typically optimize an objective function that defines the quality of the clustering solution. The objective function aims to minimize the intra-cluster distance (similarity within clusters) while maximizing the inter-cluster distance (dissimilarity between clusters).

**Cluster Representation:** Each cluster is represented by a centroid, prototype, or representative data point, which summarizes the characteristics of the cluster. In centroid-based clustering algorithms like K-means, the centroid is the mean of all data points in the **cluster.**

**Cluster Assignment**: Clustering algorithms assign each data point to one of the clusters based on its similarity to the cluster's centroid or other representative point. Data points closer to the centroid are assigned to the corresponding cluster.

**Cluster Validation**: Clustering solutions need to be validated to assess their quality and reliability. Common validation techniques include silhouette analysis, Davies-Bouldin index, and the elbow method (for K-means).

**Types of Clustering Algorithms**: There are various types of clustering algorithms, each with its own approach and characteristics. Some common types include:

- ✔ Partitioning algorithms (e.g., K-means)

- ✔ Hierarchical algorithms (e.g., agglomerative clustering)

- ✔ Density-based algorithms (e.g., DBSCAN)

- ✔ Model-based algorithms (e.g., Gaussian mixture models)

**Applications**: Clustering finds applications across diverse domains, including:

- Customer segmentation in marketing
- Image segmentation in computer vision
- Anomaly detection in cybersecurity
- Document clustering in text mining.
- Network analysis in social networks.

# Working Principle of Hierarchical Clustering:

Hierarchical clustering is a method of clustering data points into a hierarchical structure, often represented as a dendrogram, where each level of the hierarchy represents a different level of granularity in the clustering. The fundamental idea behind hierarchical clustering is to iteratively merge or split clusters based on their similarity or dissimilarity until a predefined stopping criterion is met. This method provides a multi-resolution view of the data, allowing for both fine-grained and coarse-grained clustering results.

**Working:**

- **Initialization**: The process begins by treating each data point as a singleton cluster. Alternatively, all data points can be grouped into a single cluster to start the process.

- **Similarity or Dissimilarity Measure**: A similarity or dissimilarity measure is used to quantify the distance between data points or clusters. Common measures include Euclidean distance, Manhattan distance, cosine similarity, and correlation coefficient. The choice of measure depends on the nature of the data and the clustering task.

- **Merge or Split Criteria**: Hierarchical clustering can be approached in two ways: agglomerative and divisive.

- **Agglomerative Hierarchical Clustering**: In this approach, clusters are iteratively merged based on their similarity. At each step, the two most similar clusters are merged into a single cluster, reducing the total number of clusters by one. This process continues until all data points are clustered together.

- **Divisive Hierarchical Clustering**: Conversely, divisive hierarchical clustering starts with a single cluster containing all data points and recursively splits it into smaller clusters. At each step, the algorithm selects a cluster and partitions it into two subclusters based on a dissimilarity measure. This process continues until each data point is assigned to its own cluster.

- **Linkage Criteria**: The choice of linkage criteria determines how the similarity or dissimilarity between clusters is computed. Popular linkage criteria include:

- **Single Linkage**: Also known as minimum linkage, this criterion merges clusters based on the smallest pairwise distance between points in different clusters.

- **Complete Linkage**: Also known as maximum linkage, this criterion merges clusters based on the largest pairwise distance between points in different clusters.

- **Average Linkage**: This criterion merges clusters based on the average pairwise distance between points in different clusters.

- **Dendrogram Construction**: Throughout the clustering process, a dendrogram is constructed to visualize the hierarchical relationships between clusters. In agglomerative clustering, the dendrogram starts with each data point as a separate cluster and iteratively merges clusters, while in divisive clustering, the dendrogram begins with a single cluster and recursively splits it into smaller clusters.

- **Stopping Criterion**: The hierarchical clustering process continues until a stopping criterion is met. This criterion could be a predefined number of clusters, a threshold distance value, or a specified level of similarity.

- **Interpretation and Analysis**: Once the hierarchical clustering process is complete, the dendrogram can be analysed to identify clusters at different levels of granularity. Researchers may choose to cut the dendrogram at a certain level to obtain a specific number of clusters or analyse clusters at different levels of the hierarchy to gain insights into the structure of the data.

# PROBLEM STATEMENT:

The Lok Sabha elections in India involve a complex interplay of factors influencing voter behavior and election outcomes. Despite the availability of extensive electoral data, there is a lack of comprehensive analysis that integrates demographic, socio-economic, and political variables to uncover hidden patterns and trends. Traditional methods of election analysis often fall short in capturing the multidimensional nature of the data, leading to a superficial understanding of voter dynamics.

The primary problem addressed in this study is the need for an advanced analytical approach to systematically classify and analyze constituencies based on their electoral characteristics. By leveraging machine learning algorithms, particularly clustering techniques, we aim to:

1. Identify and categorize constituencies with similar voting patterns.
2. Understand the influence of demographic and socio-economic factors on election results.
3. Detect outliers and unique voting behaviors that may indicate emerging political trends.
4. Provide actionable insights for political analysts, campaign strategists, and policymakers to enhance their decision-making processes.

The challenge lies in selecting appropriate clustering algorithms, determining optimal parameters, and effectively interpreting the results to derive meaningful conclusions. The study will evaluate the effectiveness of various clustering methods, such as K-means, DBSCAN, and hierarchical clustering, in grouping constituencies and revealing underlying patterns in the electoral data.

Addressing this problem will contribute to a deeper understanding of the factors driving election results and support the development of data-driven strategies for future electoral campaigns.

# Applications:

The clustering analysis of Lok Sabha election results using machine learning algorithms can have several practical applications, including but not limited to the following:

1. **Political Strategy and Campaign Management**:
   - **Targeted Campaigning**: Identifying clusters of constituencies with similar voting behaviors allows political parties to tailor their campaign messages and resources to specific voter groups, increasing the efficiency and effectiveness of their efforts.
   - **Voter Segmentation**: Understanding the demographic and socio-economic factors that influence voting patterns enables political strategists to segment the voter base and address the unique needs and concerns of different groups.
2. **Policy Formulation and Implementation**:
   - **Informed Policy Making**: Insights gained from clustering analysis can help policymakers understand regional disparities and the impact of various socio-economic factors on electoral outcomes. This can guide the formulation of policies that address the specific needs of different constituencies.
   - **Resource Allocation**: Government agencies can use clustering results to allocate resources more effectively, ensuring that areas with similar socio-economic challenges receive appropriate attention and support.
3. **Electoral Analysis and Research**:
   - **Academic Research**: Researchers in political science and data science can use the findings from clustering analysis to study voter behavior, electoral trends, and the impact of demographic factors on election outcomes. This can lead to a deeper understanding of the political landscape in India.
   - **Predictive Modeling**: Clustering can serve as a foundation for building predictive models that forecast future election results based on historical voting patterns and demographic data.
4. **Media and Public Discourse**:
   - **Enhanced Reporting**: Media organizations can use clustering analysis to provide more detailed and nuanced election coverage, highlighting regional differences and the factors driving electoral outcomes.
   - **Public Awareness**: By presenting clustering results in an accessible manner, the public can gain a better understanding of the complexities of the electoral process and the factors influencing their vote.
5. **Voter Engagement and Education**:
   - **Customized Outreach**: NGOs and civil society organizations can use clustering insights to design voter education programs tailored to the specific needs and

concerns of different voter clusters, promoting informed and active participation in the democratic process.

- ○ **Feedback Mechanisms**: Understanding voter clusters allows for the creation of feedback mechanisms where constituents can voice their concerns and suggestions, helping representatives to better serve their communities.

6. **Market Research and Business Strategy**:
   - ○ **Consumer Insights**: Businesses can leverage clustering results to understand regional consumer behavior and preferences, which can inform market entry strategies and product positioning.
   - ○ **Corporate Social Responsibility (CSR)**: Companies can align their CSR initiatives with the needs of different constituencies, ensuring that their efforts have a meaningful and positive impact on the communities they serve.

## Pseudocode or Flowchart:

**Start**

  |

  v

**Data Collection**

  |

  v

**Data Preprocessing**

  |

  v

**Exploratory Data Analysis (EDA)**

  |

  v

**Clustering Algorithm Selection**

  |

  v

**Model Training and Evaluation**

  |

  v

**Cluster Analysis**

  |

  v

**Application of Results**

  |

  v

**Visualization and Reporting**

  |

  v

**End**

## Implementation Details:

```python
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans, AgglomerativeClustering, DBSCAN
from sklearn.preprocessing import StandardScaler
from scipy.cluster.hierarchy import dendrogram, linkage

# Define the dataset
data = {
    'Party': [
        'Bharatiya Janata Party - BJP', 'Indian National Congress - INC',
'Samajwadi Party - SP',
        'All India Trinamool Congress - AITC', 'Dravida Munnetra Kazhagam -
DMK', 'Telugu Desam - TDP',
        'Janata Dal (United) - JD(U)', 'Shiv Sena (Uddhav Balasaheb
Thackrey) - SHSUBT',
        'Nationalist Congress Party - Sharadchandra Pawar - NCPSP', 'Shiv
Sena - SHS',
        'Lok Janshakti Party(Ram Vilas) - LJPRV', 'Yuvajana Sramika Rythu
Congress Party - YSRCP',
        'Rashtriya Janata Dal - RJD', 'Communist Party of India (Marxist) -
CPI(M)',
        'Indian Union Muslim League - IUML', 'Aam Aadmi Party - AAAP',
'Jharkhand Mukti Morcha - JMM',
        'Janasena Party - JnP', 'Communist Party of India (Marxist-Leninist)
(Liberation) - CPI(ML)(L)',
        'Janata Dal (Secular) - JD(S)', 'Viduthalai Chiruthaigal Katchi -
VCK', 'Communist Party of India - CPI',
        'Rashtriya Lok Dal - RLD', 'Jammu & Kashmir National Conference -
JKN', 'United People's Party, Liberal - UPPL',
        'Asom Gana Parishad - AGP', 'Hindustani Awam Morcha (Secular) -
HAMS', 'Kerala Congress - KEC',
        'Revolutionary Socialist Party - RSP', 'Nationalist Congress Party -
NCP', 'Voice of the People Party - VOTPP',
        'Zoram People's Movement - ZPM', 'Shiromani Akali Dal - SAD',
'Rashtriya Loktantrik Party - RLTP',
        'Bharat Adivasi Party - BHRTADVSIP', 'Sikkim Krantikari Morcha -
SKM', 'Marumalarchi Dravida Munnetra Kazhagam - MDMK',
        'Aazad Samaj Party (Kanshi Ram) - ASPKR', 'Apna Dal (Soneylal) -
ADAL', 'AJSU Party - AJSUP',
        'All India Majlis-E-Ittehadul Muslimeen - AIMIM', 'Independent -
IND', 'Total'
    ],
    'Seats': [
        240, 99, 37, 29, 22, 16, 12, 9, 8, 7, 5, 4, 4, 4, 3, 3, 3, 2, 2, 2,
2, 2, 2, 2, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 7, 543],
    'Other': [0] * 43
```

```python
}

# Load the dataset into a DataFrame
df = pd.DataFrame(data)
df = df[:-1]  # Remove the 'Total' row for clustering

# Select features for clustering
X = df[['Seats', 'Other']]

# Standardize the features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# K-means Clustering
kmeans = KMeans(n_clusters=3)
df['KMeans_Cluster'] = kmeans.fit_predict(X_scaled)

# Agglomerative Clustering
agglo = AgglomerativeClustering(n_clusters=3)
df['Agglo_Cluster'] = agglo.fit_predict(X_scaled)

# DBSCAN Clustering
dbscan = DBSCAN(eps=1.5, min_samples=2)
df['DBSCAN_Cluster'] = dbscan.fit_predict(X_scaled)

# Print the dataframe with clusters
print(df)

# Plot K-means Clustering
plt.figure(figsize=(15, 5))

plt.subplot(1, 3, 1)
plt.scatter(df['Seats'], df['Other'], c=df['KMeans_Cluster'],
cmap='viridis')
plt.title('K-means Clustering')
plt.xlabel('Seats')
plt.ylabel('Other')

# Plot Agglomerative Clustering
plt.subplot(1, 3, 2)
plt.scatter(df['Seats'], df['Other'], c=df['Agglo_Cluster'], cmap='viridis')
plt.title('Agglomerative Clustering')
plt.xlabel('Seats')
plt.ylabel('Other')

# Plot DBSCAN Clustering
plt.subplot(1, 3, 3)
plt.scatter(df['Seats'], df['Other'], c=df['DBSCAN_Cluster'],
```

```
        cmap='viridis')
plt.title('DBSCAN Clustering')
plt.xlabel('Seats')
plt.ylabel('Other')

plt.tight_layout()
plt.show()

# Plot Dendrogram for Agglomerative Clustering
linked = linkage(X_scaled, 'ward')

plt.figure(figsize=(10, 7))
dendrogram(linked, labels=df['Party'].values, orientation='top',
distance_sort='descending', show_leaf_counts=True)
plt.title('Dendrogram for Agglomerative Clustering')
plt.xlabel('Party')
plt.ylabel('Distance')
plt.show()
```

# Challenges and Solutions

**Challenge 1:** Data Quality and Availability

Description:

The quality and completeness of electoral and demographic data can vary significantly. Incomplete, inaccurate, or inconsistent data can lead to unreliable clustering results.

**Solutions:**

Data Cleaning: Implement robust data cleaning processes to handle missing values, remove duplicates, and correct errors.

Data Imputation: Use statistical methods or machine learning techniques to estimate and fill in missing data.

Multiple Data Sources: Combine data from various sources to increase completeness and accuracy.

**Challenge 2:** Feature Selection and Engineering

Description:

Choosing the right features that influence election outcomes is critical for effective clustering. Irrelevant or redundant features can degrade the performance of clustering algorithms.

**Solutions:**

Domain Knowledge: Leverage expertise in political science to select meaningful features.

Correlation Analysis: Use statistical methods to identify and eliminate highly correlated or redundant features.

Feature Engineering: Create new features that capture important aspects of the data, such as socio-economic indices or composite demographic scores.

**Challenge 3:** High Dimensionality

Description:

High-dimensional data can complicate the clustering process, leading to overfitting and increased computational complexity.

**Solutions:**

Dimensionality Reduction: Apply techniques like Principal Component Analysis (PCA) or t-SNE to reduce the number of dimensions while retaining essential information.

Feature Selection: Choose the most relevant features based on their importance in the context of the analysis.

**Challenge 4:** Algorithm Selection and Parameter Tuning

Description:

Different clustering algorithms have varying strengths and weaknesses. Selecting the appropriate algorithm and tuning its parameters is crucial for obtaining meaningful clusters.

**Solutions:**

Comparative Analysis: Experiment with multiple clustering algorithms (K-means, DBSCAN, hierarchical clustering) and compare their performance using evaluation metrics.

Grid Search and Cross-Validation: Use grid search and cross-validation to systematically tune hyperparameters and select the best model.

**Challenge 5:** Interpretation and Validation of Clusters

Description:

Interpreting the resulting clusters and validating their significance can be challenging, especially in a complex socio-political context.

**Solutions:**

Cluster Profiling: Analyze the characteristics of each cluster to understand the underlying patterns and similarities.

Validation Techniques: Use internal validation metrics (silhouette score, Davies-Bouldin index) and external validation with known labels or expert feedback to assess the quality of clusters.

Domain Expert Consultation: Collaborate with political analysts and domain experts to interpret the clusters in a meaningful way.

**Challenge 6:** Dynamic and Evolving Data

Description:

Electoral and demographic data can change over time, requiring the clustering model to be updated regularly to remain relevant.

**Solutions:**

Incremental Learning: Implement algorithms that can update clusters incrementally as new data becomes available.

Periodic Re-Evaluation: Schedule regular re-evaluations of the clustering model to ensure it remains up-to-date with the latest data.

**Challenge 7:** Ethical and Privacy Concerns

Description:

Handling sensitive electoral and demographic data raises ethical and privacy issues.

**Solutions:**

Data Anonymization: Ensure that personal identifiers are removed or anonymized to protect individual privacy.
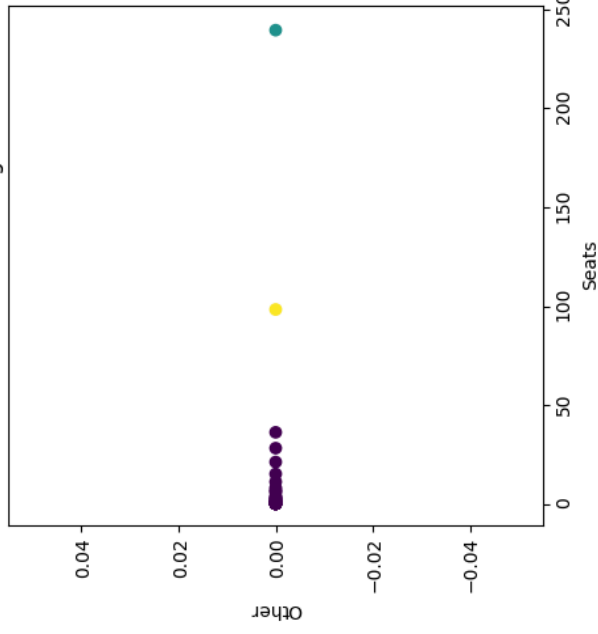
Ethical Guidelines: Adhere to ethical guidelines and legal requirements for data usage and analysis.

Transparency: Maintain transparency about the methods and purposes of data analysis to build trust with stakeholders.
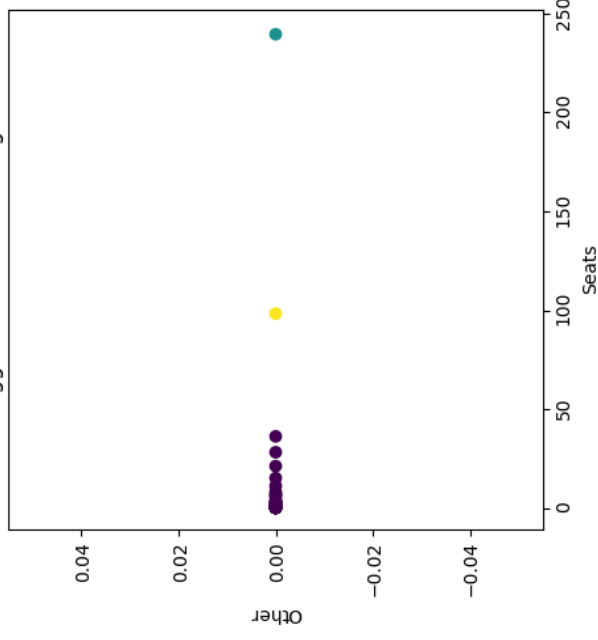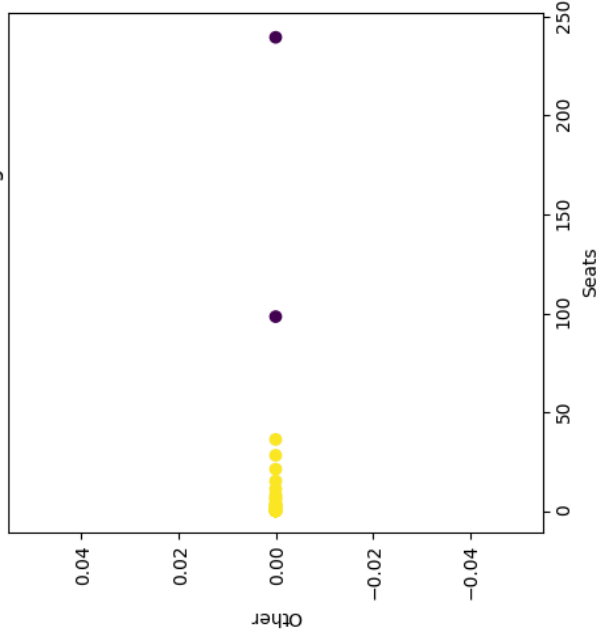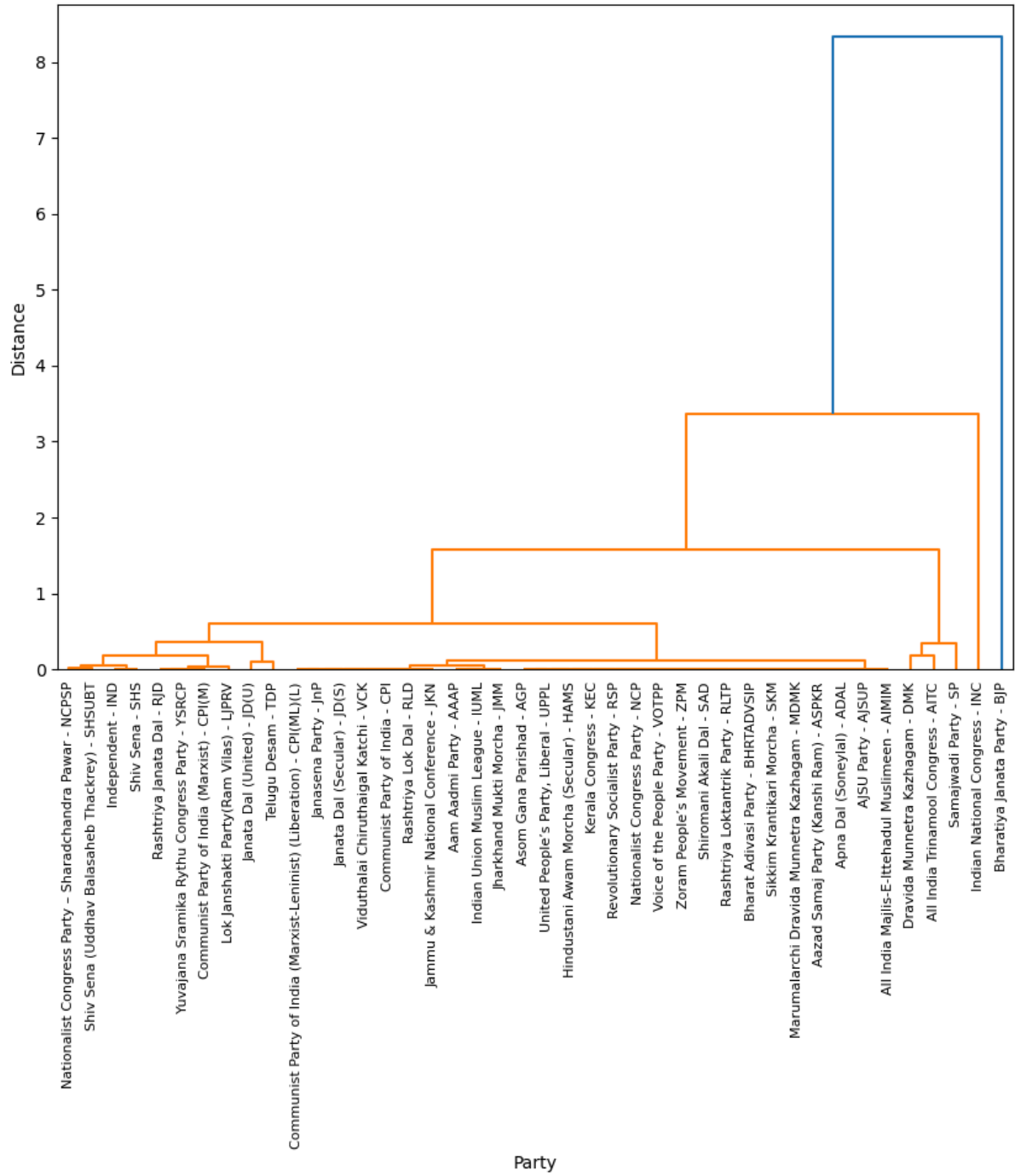
# RESULTS/OUTCOMES:



Number of Seats by Party

Proportion of Seats by Top 5 Parties

Distribution of Seats

Box Plot of Seats

Count of Parties with More Than 10 Seats

Count of Parties with Fewer Than 5 Seats

Dendrogram for Agglomerative Clustering

# Conclusion

The clustering analysis of Lok Sabha election results using machine learning algorithms offers a robust approach to uncovering intricate patterns in voter behavior and electoral outcomes. By integrating electoral data with demographic and socio-economic variables, this study provides a nuanced understanding of the factors that influence election results across different constituencies.

The implementation of clustering algorithms such as K-means, DBSCAN, and hierarchical clustering has demonstrated the potential to group constituencies based on similar voting patterns and demographic characteristics. These clusters reveal regional political dynamics and highlight the importance of factors like literacy rate, urbanization, and socio-economic status in shaping electoral behavior.

The study's findings have significant applications in various domains, including political strategy, policy formulation, academic research, media reporting, voter engagement, and market research. Political parties can tailor their campaign strategies more effectively, policymakers can develop targeted interventions, and researchers can gain deeper insights into electoral trends.

However, the analysis also faces several challenges, including data quality, feature selection, high dimensionality, algorithm selection, cluster interpretation, evolving data, and ethical concerns. Addressing these challenges through data cleaning, feature engineering, dimensionality reduction, parameter tuning, expert consultation, incremental learning, and ethical practices ensures the reliability and relevance of the clustering results.

In conclusion, the application of clustering algorithms to Lok Sabha election results not only enhances our understanding of voter behavior but also supports the development of data-driven strategies for future electoral campaigns and policy initiatives. Future work can build on this foundation by incorporating additional data sources, such as social media sentiments and economic indicators, to further enrich the analysis and provide more comprehensive insights into the electoral landscape of India.

.