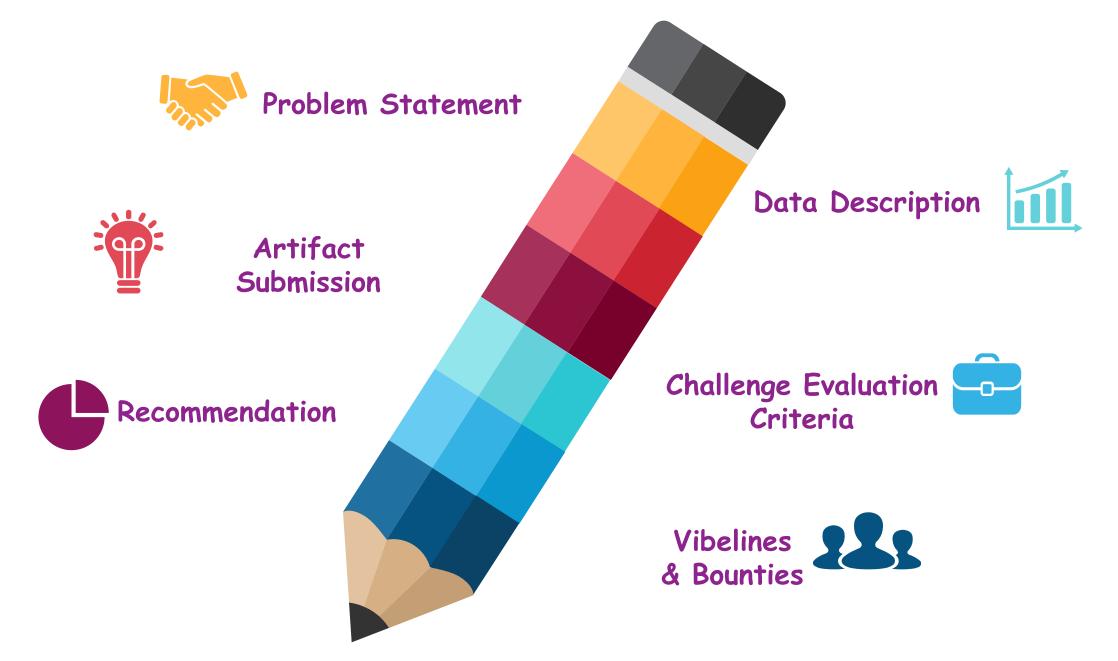
Project on Clustering









Problem Statement - Cereal Product Segmentation for Nutritional Strategy

Goal

The primary goal of this project is to apply three distinct clustering methodologies (K-Means, Hierarchical, and DBSCAN) to a dataset of popular cereals. This analysis will identify natural segments of cereals based on their nutritional profiles, allowing a manufacturer or retailer to understand the existing competitive landscape and inform future product development or marketing efforts.

Dataset

Cereal Nutritional Data: The dataset includes various key features for a range of cereals:

- Identification: Cereal Name, Manufacturer
- Nutritional Constituents (Key Features for Clustering): Calories, Protein (g), Fat (g), Sugars (g), and Vitamin and Minerals.

Objectives

- 1. Preparation & Scaling: Analyze the distribution of the nutritional features (especially Calories, Protein, Fat, and Sugars) and apply appropriate scaling (e.g., standardizing) to ensure all features contribute equally to the distance calculation.
- 2. Partitional Clustering (K-Means):
 - Use the Elbow Method and Silhouette Score to determine the optimal number of clusters (\$k\$).
 - Apply K-Means with the chosen \$k\$ to segment the cereals.
- 3. Hierarchical Clustering:
 - Create a dendrogram to visually explore the natural groupings in the data.
 - Apply Hierarchical Clustering (using different linkages like Ward or Complete) to validate the segments found by K-Means.
- 4. Density-Based Clustering (DBSCAN):
 - Apply DBSCAN to identify dense regions of cereals, treating any isolated products as potential outliers or unique, unclustered items.
- 5. Cluster Interpretation: Analyze the mean nutritional values for the final, chosen clusters (e.g., "Cluster 1: High Sugar/Low Protein," "Cluster 2: Healthy/Low Calorie").

Deliverable

A comprehensive report comparing the results of all three clustering methods (K-Means, Hierarchical, DBSCAN). The final analysis should define the cereal segments and provide insights, such as identifying a gap in the market for a specific nutritional profile.



Project Outcomes

The outcome of this project will be a clear segmentation of cereals based on their nutritional profiles. This insight can be invaluable for:

- Cereal Manufacturers: Informing product development (e.g., identifying gaps in the market, creating new cereals for specific health segments), and tailoring marketing messages to target consumers interested in specific nutritional benefits.
- Marketers: Developing targeted advertising campaigns that highlight the nutritional aspects appealing to different consumer preferences.
- Health Professionals/Consumers: Providing a simplified way to understand and categorize cereals for dietary planning or healthy eating choices.



Data Description

Dataset Details:

· Dataset Name: Cereal dataset with Nutritional constituent

Column description (Key Features for Clustering):

- 1. Cereal Name: name of the cereal
- 2. Manufacturer: manufacturer of the cereal
- 3. Calories: calories consumed per 100g
- 4. Protein (g): protein in grams per 100g
- 5. Fat: fat per 100g
- 6. Sugars: sugar per 100g
- 7. Vitamin and Minerals: vitamin and minerals per 100g



Artifact Submission

Your submission must include the following five artifacts, all packaged within a single GitHub repository.

1. Jupyter Notebook (.ipynb) (Add separate ipynb files for Kmeans , Hierarchical and DBSCAN Clustering)

This is the core of your submission. Your Jupyter Notebook should be a complete, well-documented narrative of your data analysis journey. It must include:

- Detailed Explanations: Use Markdown cells to explain your thought process, the questions you are trying to answer, and the insights you've uncovered.
- Clean Code: The code should be well-structured, easy to read, and free of unnecessary clutter.
- Comprehensive Comments: Use comments to explain complex logic and the purpose of different code blocks.
- · Key Visualizations: All visualizations should be clear, properly labeled, and directly support your findings.

2. Presentation (.pptx or .pdf)

Create a compelling presentation that summarizes your team's analysis and key findings. This presentation should serve as your final pitch. It must include:

- Executive Summary: A concise overview of your findings.
- · Key Insights: The most important takeaways from your analysis.
- Data-Driven Recommendations: Actionable steps that can be taken based on your insights.
- Supporting Visualizations: A selection of your best visualizations to illustrate your points.

3. README File (.md)

The README file is the first thing we'll look at. It should serve as a quick guide to your project and provide essential details. It must include:

- Project Title :
- Brief Problem Statement: A summary of the project and your approach.
- · Summary of Findings: A bullet-point summary of your most significant insights.

4. Attached Dataset

Please include the original dataset (.csv or other format) within your repository. This ensures the judges can reproduce your analysis without any issues.

5. GitHub Repository

Your final submission will be your GitHub repository. The repository name must follow this exact format: Clustering_ProjectName_TMP



Challenge Evaluation Criteria

Criteria Name	Criteria weight
Data Understanding and Exploratory Data Analysis	20%
Data preprocessing and feature engineering	25%
Model building and evaluation	30%
Business Recommendation	15%
Coding guidelines and standards	10%



Recommendation for Kmeans Clustering Implementation

1. Data Preprocessing:

- Selecting only the numerical columns representing nutritional constituents (Calories, Protein (g), Fat, Sugars, Vitamin and Minerals).
- Crucially, performing feature scaling on these columns (e.g., using StandardScaler). This is essential because K-Means is a distance-based algorithm, and features like Calories or Sugars might have much larger numerical ranges than Fat or Protein, disproportionately influencing the distance calculations if not scaled.

2. Determining the Optimal 'K':

• Applying the Elbow Method (plotting inertia for various K values) and/or Silhouette Score to determine the most appropriate number of clusters (K) for the cereal dataset. This will help identify natural groupings of cereals based on their nutritional makeup.

3. K-Means Implementation:

- Applying the K-Means algorithm with the chosen K to the scaled nutritional data.
- The algorithm will assign a cluster label to each cereal, grouping those with similar nutritional profiles.

4. Cluster Profiling:

 Analyzing the characteristics of each identified cluster. For example, a cluster might be defined by high Sugars and Calories ("Sweet & High-Calorie Cereals"), while another might show high Protein (g) and Vitamin and Minerals ("Nutrient-Dense Cereals"). This involves calculating the average nutritional values for each cluster.

5. Visualization:

Since there are multiple numerical features, dimensionality reduction techniques like PCA or t-SNE can be applied before or after clustering to visualize the clusters in 2D or 3D, making the groupings visually apparent.

Recommendation for Hierarchical Clustering Implementation

1. Data Preprocessing:

- Selecting only the numerical columns representing nutritional constituents (Calories, Protein (g), Fat, Sugars, Vitamin and Minerals).
- Crucially, performing feature scaling on these columns (e.g., using StandardScaler). This is essential because hierarchical clustering is a distance-based algorithm, and features like Calories or Sugars might have much larger numerical ranges than Fat or Protein, disproportionately influencing the distance calculations if not scaled.

2. Distance Matrix Calculation:

Calculating the pairwise distances between all cereal data points using a chosen distance metric (e.g., Euclidean distance).

3. Linkage Method Application:

Applying a chosen linkage criterion (e.g., Ward's method, Average linkage) to define the distance between clusters and build the
hierarchy. Ward's method is often a good starting point for numerical data as it tends to produce compact clusters.

4. Dendrogram Visualization:

• Generating and visualizing the dendrogram. This tree-like diagram will show how individual cereals are successively merged into larger clusters based on their nutritional similarity.

5. Determining the Number of Clusters:

 By visually inspecting the dendrogram, identifying a suitable "cut-off" point (a horizontal line) that yields a meaningful number of clusters. This allows for flexibility in choosing the granularity of segmentation, for example, distinguishing between "Breakfast Cereals," "Snack Cereals," or "Dietary Specific Cereals."

6. Cluster Assignment & Profiling:

- Assigning each cereal to a specific cluster based on the chosen cut-off point on the dendrogram.
- Analyzing the characteristics (e.g., average nutritional values) of cereals within each identified cluster. For example, one cluster might be "High-Sugar, Low-Protein Kids' Cereals," while another could be "High-Fiber, Low-Fat Adult Cereals."



Recommendation for DBSCAN Clustering Implementation

1. Data Preprocessing:

- Selecting only the numerical columns representing nutritional constituents (Calories, Protein (g), Fat, Sugars, Vitamin and Minerals).
- Crucially, performing feature scaling on these columns (e.g., using StandardScaler). This is essential because DBSCAN is a distance-based algorithm, and features like Calories or Sugars might have much larger numerical ranges than Fat or Protein, disproportionately influencing the distance calculations if not scaled.

2. DBSCAN Implementation:

- Applying the DBSCAN algorithm to the scaled nutritional data.
- Careful tuning of the eps and min_samples hyperparameters will be critical. The choice of these parameters will directly influence what is considered a "dense region" of normal cereal nutritional profiles and, thus, how clusters are formed and which cereals are labeled as "noise" (outliers).

3. Cluster Assignment & Anomaly Identification:

• DBSCAN will assign a cluster label to each cereal. Cereals assigned to a numerical cluster belong to a segment, while those labeled -1 are identified as noise points (anomalies).

4. Cluster Profiling & Anomaly Analysis:

- Analyzing the characteristics of each identified cluster. For example, one cluster might represent cereals high in Sugars and Calories ("Sweet & Indulgent Cereals"). Another could be "High-Protein, Low-Fat Cereals," and a third "Vitamin-Fortified Options."
- Investigating the cereals flagged as noise points. These would be cereals whose nutritional compositions are so unique or extreme that they
 don't fit into any dense cluster. This could include cereals with unusually high or low values for certain nutrients compared to the rest of the
 market.

5. Visualization:

• Since there are multiple numerical features, dimensionality reduction techniques like PCA or t-SNE can be applied before or after clustering to visualize the clusters and noise points in 2D or 3D, making the groupings and outliers visually apparent.



Vibelines

VIBELINES

```
Kickoff Vibes - T
Wrapoff Vibes - T + 7
Spotlight Vibes - T + 14
```



Lets Go



