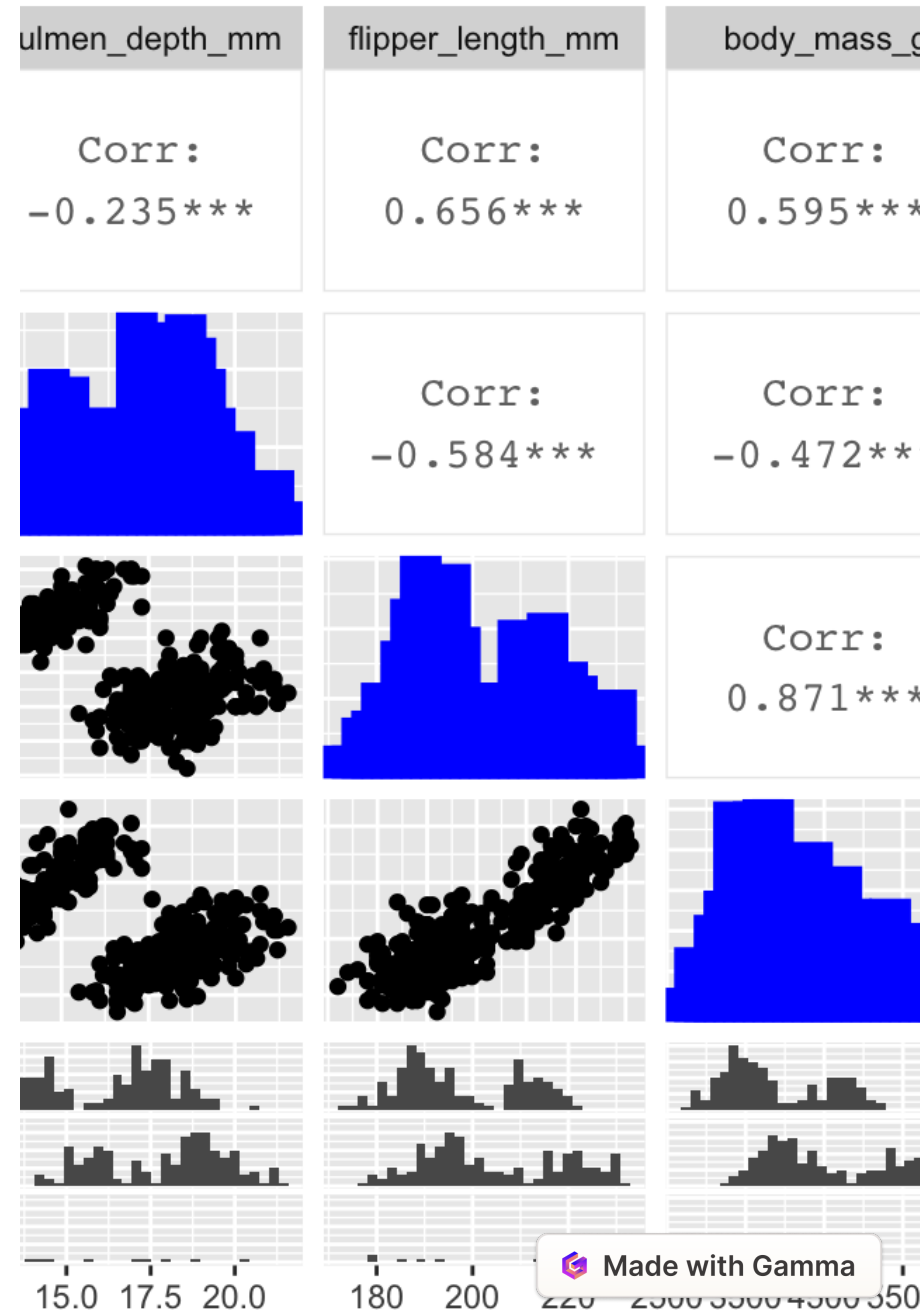# KMeans: A Guide to Unsupervised Machine Learning

KMeans is a popular unsupervised machine learning algorithm used for clustering tasks. This guide provides an explanation of its workings, applications, and the simplified mathematics behind it.

# How It Works

- Clustering Objective: KMeans partitions data points into distinct non-overlapping clusters by assigning each data point to the cluster with the nearest mean.

- Algorithm Steps:

1. **Initialization:** Choose initial centroids randomly from the data points.

2. **Assignment Step:** Assign each data point to the nearest centroid, forming clusters.

3. **Update Step:** Recalculate the centroids as the mean of all points in the cluster.

4. **Iterate:** Repeat the assignment and update steps until the centroids converge.

# When to Use It

**Market Segmentation**

Identify different customer groups based on purchasing behavior.

**Document Clustering**

Group similar documents together in text mining.

**Image Segmentation**

Partition an image into regions with similar pixel intensities.

**General Data Analysis**

Discover underlying patterns or groupings in data.

# Simplified Maths Behind the Model

- Choosing Centroids: Start by randomly selecting K points from the dataset as initial centroids.

- Assigning Points to Centroids: For each data point, find the nearest centroid using distance measures like Euclidean distance.

- Calculating New Centroids: Recalculate the centroid of each cluster as the average of all points in that cluster.

- Repeating the Process: Repeat the process until the centroids do not change much between iterations, indicating a stable grouping.
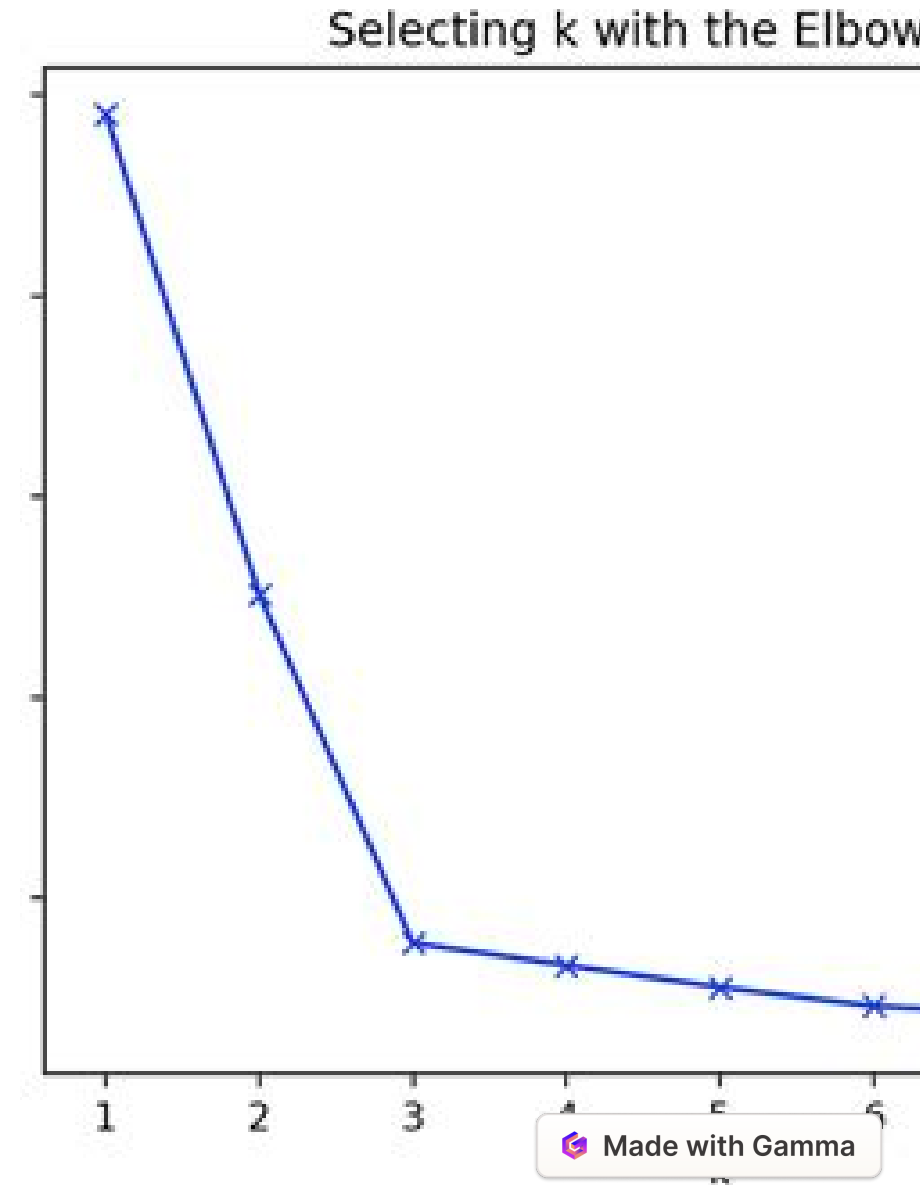
# Choosing the Number of Clusters (K)

## Elbow Method

This method involves running KMeans clustering on the dataset for a range of values of K and calculating the sum of squared distances from each point to its assigned center. Plotting these distances against the number of clusters helps determine an optimal value for K.

The "elbow" of the curve represents a good trade-off between the number of clusters and the within-cluster sum of squares.

By organizing a set of different colored balls into groups based on similarity, KMeans provides a way to discover patterns in data.

Figure 1

Selecting k with the Elbow

# KMeans: Transforming Unlabeled Data

KMeans can be used to transform unlabeled data into a labeled dataset by assigning each data point to a cluster label.

This enables further analysis and creates a structure in an otherwise unstructured dataset.

# Optimizing KMeans Performance

- **Choosing the Right K:** Selecting an appropriate number of clusters increases the accuracy of the model.

- **Scaling the Data:** Preprocessing data by scaling ensures that features with larger scales do not dominate the clustering process.

- **Handling Outliers:** Outliers can significantly affect cluster formation, so removing or treating them carefully is important.

# Evaluating KMeans Results

## 1. Silhouette Score

The silhouette score measures how well each data point fits into its assigned cluster, ranging from -1 to 1. A higher score indicates better separation between clusters.

Here's the step-by-step calculation of the Silhouette Score for a single data point:

1. For each data point $i$, calculate the average distance ($a(i)$) from the data point to all other data points within the same cluster. This measures how well data point $i$ is clustered with its own cluster members.

2. For the same data point $i$, calculate the smallest average distance ($b(i)$) from the data point to all data points in a different cluster, where the data point is not a member. This measures how well data point $i$ could belong to another nearby cluster.

3. Calculate the silhouette score $s(i)$ for data point $i$ using the formula:
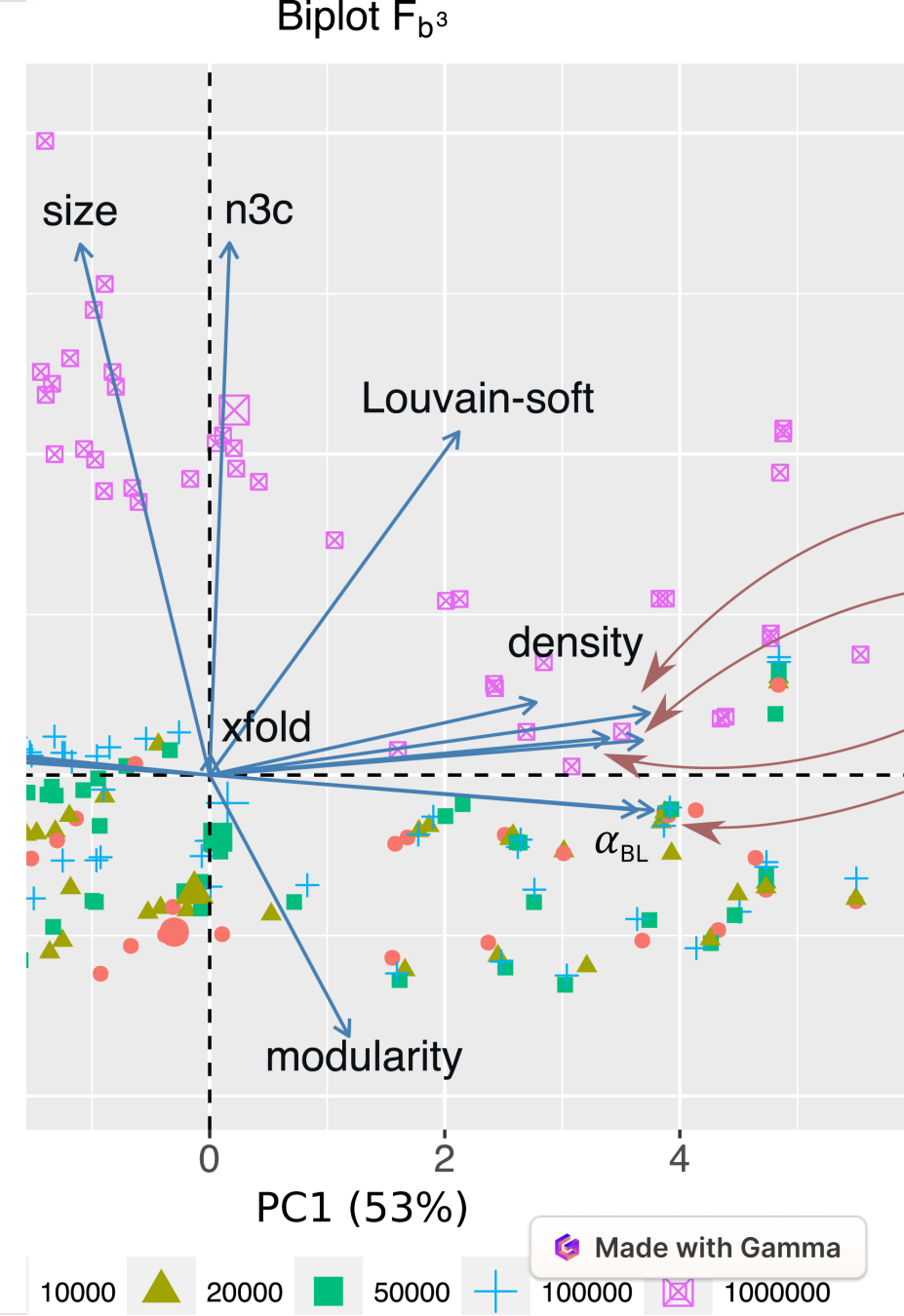
$$s(i) = (b(i) - a(i)) / \max\{a(i), b(i)\}$$

# KMeans: Beyond Hard Clustering

KMeans can also be extended to perform soft clustering, where each data point is assigned a probability distribution indicating its likelihood of belonging to each cluster.

This allows for more nuanced interpretations and a probabilistic perspective on cluster assignments.



Biplot $F_{b^3}$

# Conclusion

## Flexible Algorithm

KMeans is a versatile algorithm that can be applied to a wide range of unsupervised learning problems.

## Discover Hidden Patterns

By grouping data points into clusters, KMeans helps reveal underlying patterns and structures in your data.

## Data Preprocessing

Scaling, outlier handling, and feature engineering play important roles in optimizing KMeans performance.

## Continued Advancements

Ongoing research in KMeans and its variants continues to enhance its performance and applicability in various domains.