# Clustering Analysis

By Karis Kim

# Clustering:
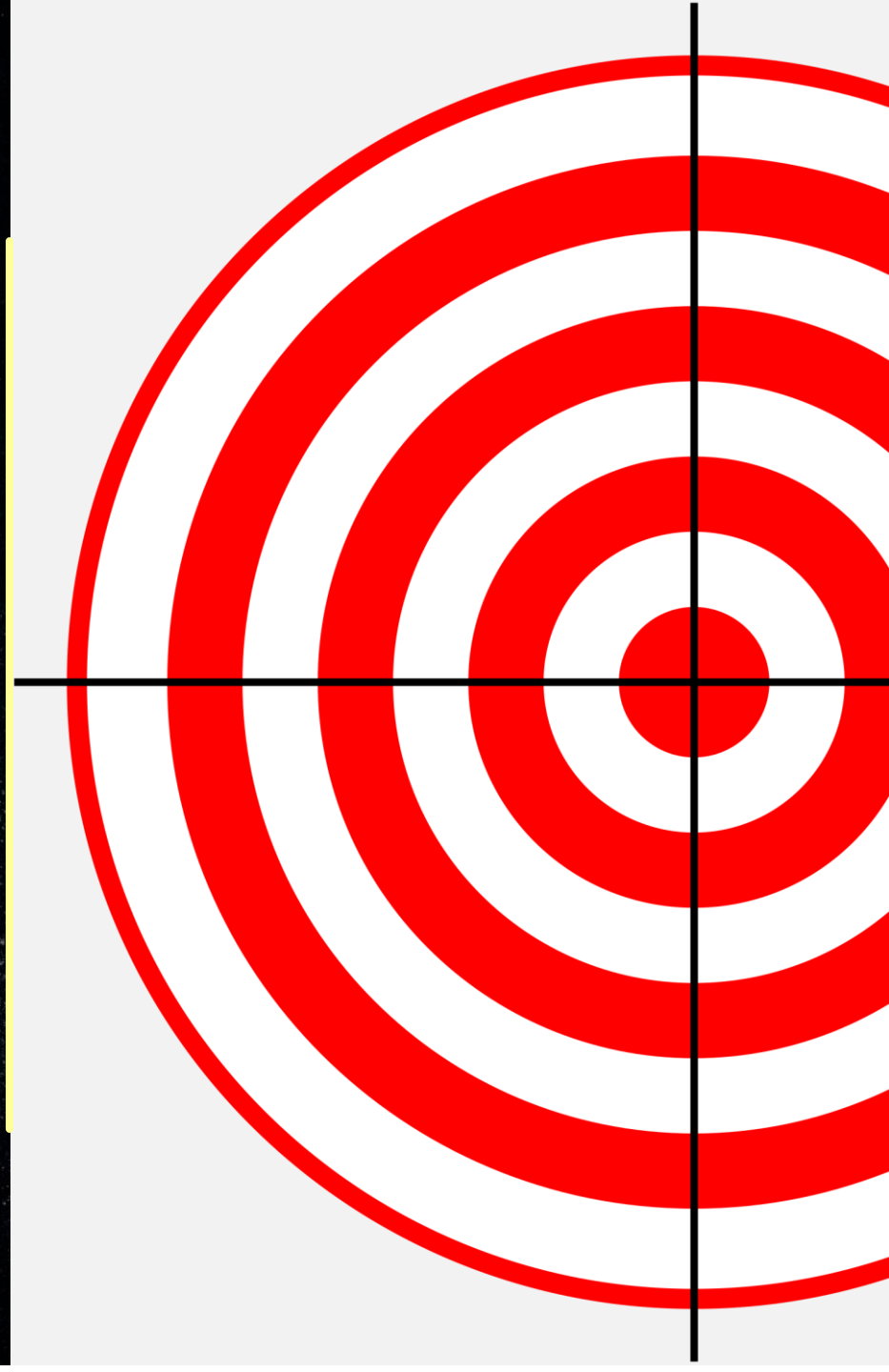
/ˈkləstərNG/

Process of finding meaningful groups in data

# Goal:

Exploring,
not predicting

# Applications of Clustering

## Describing

1. Marketing
2. Document Clustering
3. Session Grouping

## Preprocessing

1. Reduce Dimensionality
2. Object Reduction

# Clustering Types

**1** Exclusive/Strict Partitioning

**2** Overlapping/ Multiview

**3** Hierarchical

**4** Fuzzy/ Probabilistic

# Clustering Algorithm Types

**1** Prototype/ Centroid/ Center Based
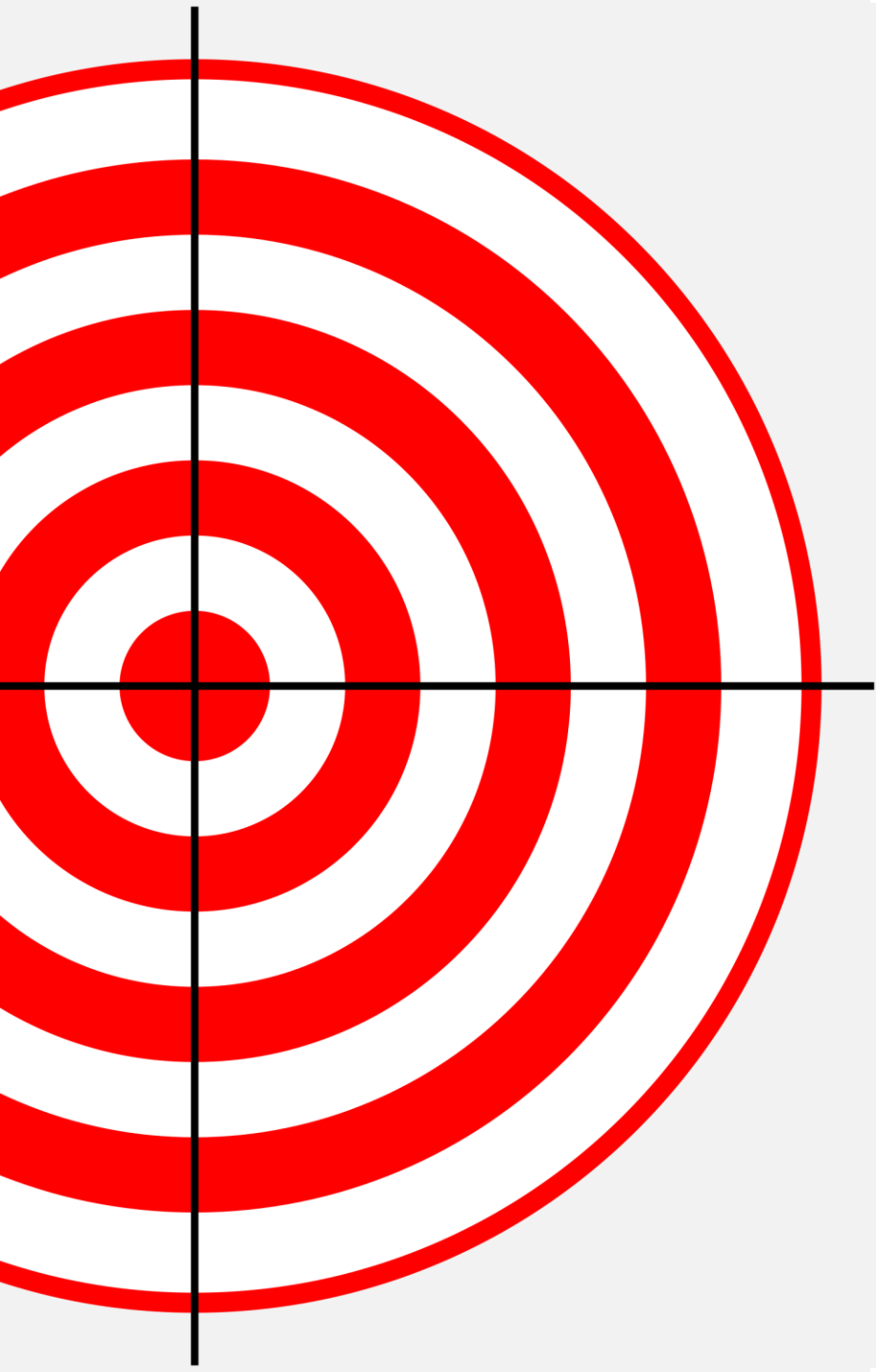
**2** Density

**3** Hierarchical

**4** Model/ Distribution Based

# k-Means:
/k-meens/

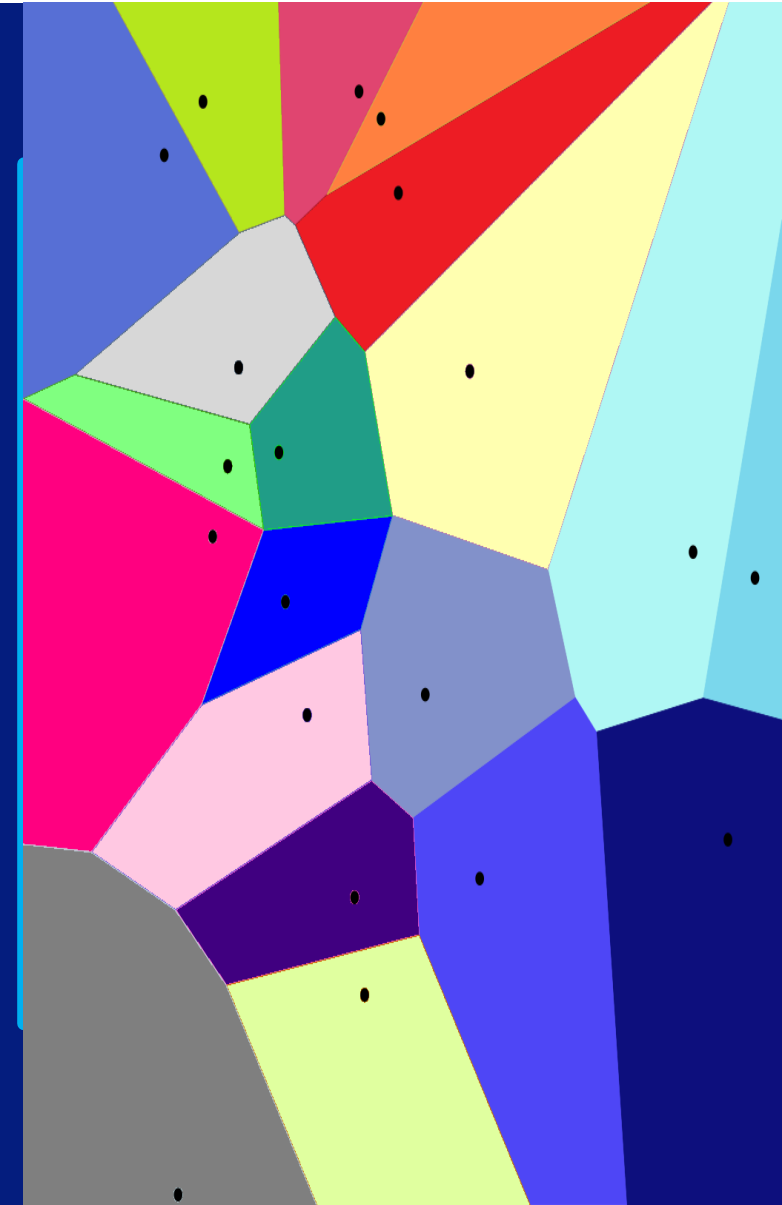Data set is divided into user-specified _k_ clusters around the nearest centroid

# Goal:

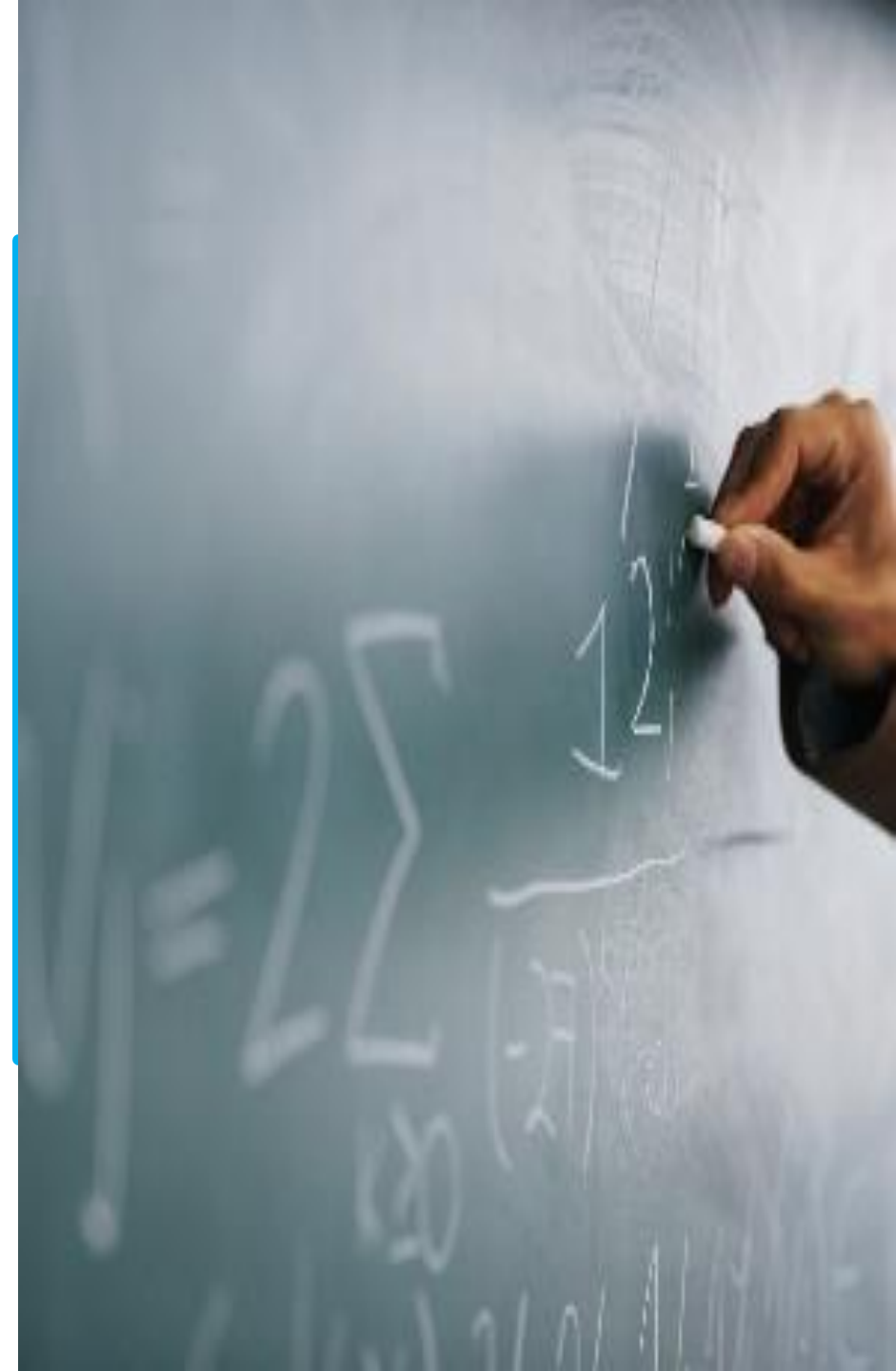To find the prototype data point for each cluster

# k-Means Process

1. Set *k* # of clusters,
   initiate *k* random centroids

2. Assign data points to nearest centroid

3. Find new centroids (most representative)
   as centroid w/minimal SSE
   = new mean of cluster

4. Repeat step 2 and 3

5. When no significant change,
   declare final centroid as prototype

# HOW TO PARTITION?

Define a proximity measure using Euclidean distance

# RapidMiner vs. Orange

**RapidMiner**

- **Performance operator** measures avr within centroid distance & Davies Bouldin index

- Lower the index the better*

- No limit on sample size

- Cluster Model Visualizer operator to see clusters

**Orange**

- **Silhouette scores** measure how well on avr each data pt fits into designated cluster

- Higher the score the better

- Not computed >5000 samples*

- Cluster visualized automatically; syncs interactively to Silhouette

# Comparison Summary

|  | RapidMiner | Orange |
|---|:---:|:---:|
| Measures cluster quality | Y | Y |
| Manually select centroids* | N | Y |
| Sample size limit for cluster quality* | N | Y |
| Normalize data* | Y | N |
| Cluster results visualization | Y | Y |

# Sample Data Information

## KSU Marietta Campus Atrium Bldg 3 FL Traffic Pattern Analysis

- 11790 examples with 2 attributes (x & y coordinates)

- Data captured by beacons every 10 sec for 1 hour

- Goal: Monitor traffic pattern clusters to inform virtual/mixed reality application for indoor navigation assistance for the blind

- Idea was to use beacons to detect heavy/light traffic areas indoors and warn blind users, as change from heavy to light/no traffic could indicate construction or path blockage

# Thank You