# CS5228 LECTURE 2: CLUSTERING
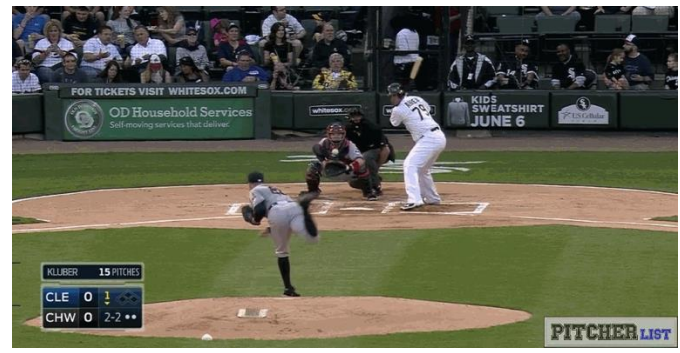
Bryan Hooi

School of Computing

National University of Singapore
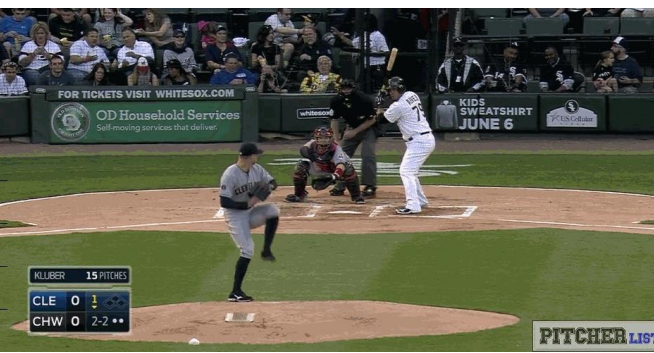
# OUTLINE



Introduction → Clustering Concepts → K-Means Algorithm → Proof of Convergence → Drawbacks → Variants
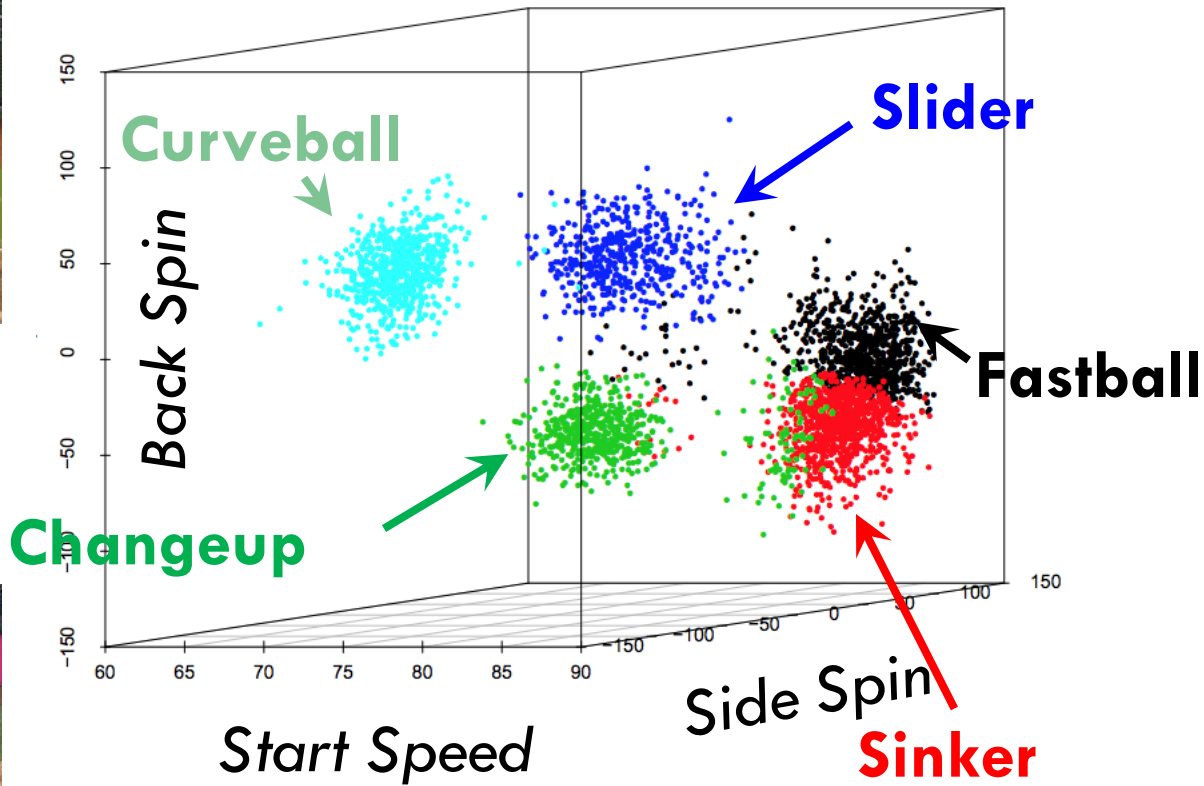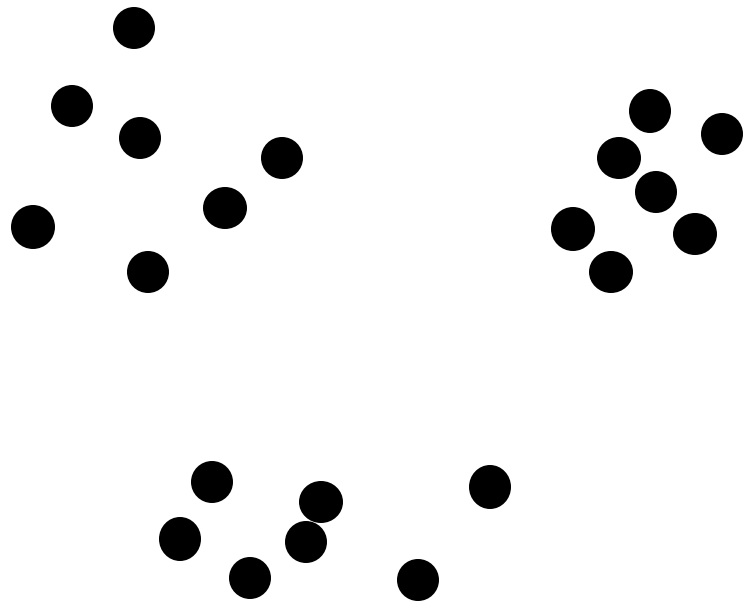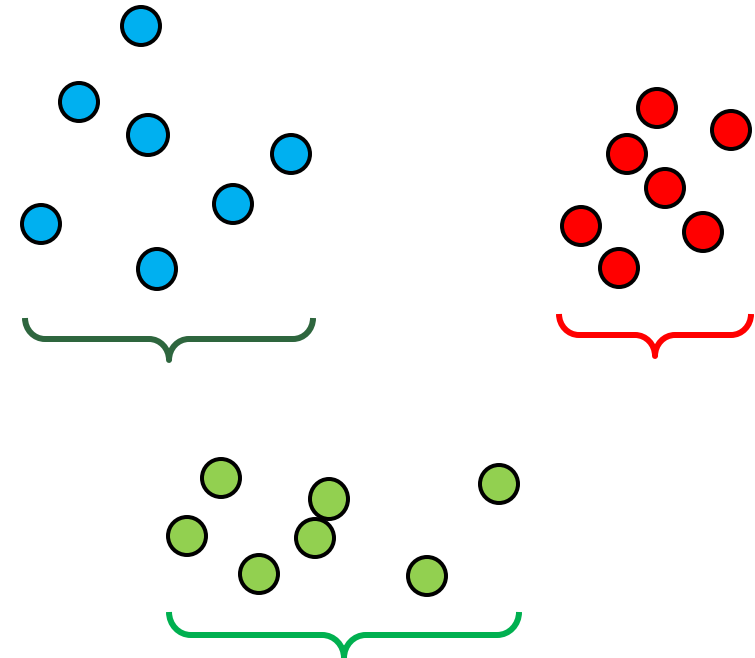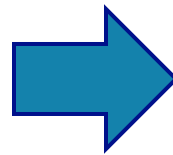
# GOAL OF CLUSTERING

Clustering separates **unlabelled** data into **groups** of similar points.

Clusters should have high **intra-cluster similarity,** and low **inter-cluster similarity.**



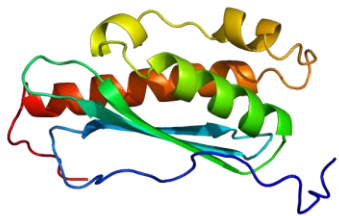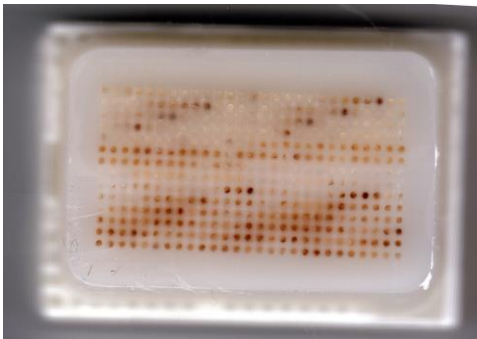**Unlabelled data**                    **Groups**

# APPLICATIONS OF CLUSTERING

Many applications:



**Microbiology:** find groups of related genes (or proteins etc.)



**Recommendation & Social Networks:** find groups of similar users



**Search & Information Retrieval:** grouping similar search (or news etc.) results

# CLUSTERING VS. CLASSIFICATION

**Clustering**

**Classification**

# OUTLINE

# WHAT DOES SIMILARITY MEAN?



(These are quite similar at the
**pixel level,** but not **semantically**)

*(In terms of their
"meaning")*

# DEFINITION OF A DISTANCE METRIC

Given a set $S$, a **distance metric** is a **nonnegative** function $d : S \times S \to \mathbb{R}^{\geq 0}$ satisfying the properties:

*Equivalent to*

*Nonnegative real numbers*

- Uniqueness: $d(a, b) = 0 \Leftrightarrow a = b$

    *(We don't want there to be objects that we cannot tell apart)*

- Symmetry: $d(a, b) = d(b, a)$

    *(If Alice is like Bob, then Bob is like Alice)*

- Triangle Inequality: $d(a, b) \leq d(a, c) + d(c, b)$

    *(Otherwise, Alice could be very like Carol, and Carol very like Bob, but Alice very unlike Bob)*

# COMMON DISTANCE / SIMILARITY METRICS

▪ **Euclidean distance**

$$d(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|_2 = \sqrt{\sum_{i=1}^{p}(a_i - b_i)^2)}$$

# COMMON DISTANCE / SIMILARITY METRICS

- **Manhattan distance**

$$d(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|_1 = \sum_{i=1}^{p} |a_i - b_i|$$

# COMMON DISTANCE / SIMILARITY METRICS

- **Cosine distance**

$$d(\mathbf{a}, \mathbf{b}) = \cos\theta = \frac{\mathbf{a} \cdot \mathbf{b}}{\|a\| \cdot \|b\|}$$

# COMMON DISTANCE / SIMILARITY METRICS

- **Jaccard Similarity**
  (between **sets** A and B)

$$s_{\mathrm{Jaccard}}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

A = { 🍞 , 🥛 }    B = { 🧀 , 🥛 }

$$s_{\mathrm{Jaccard}} = \frac{🥛}{🧀 , 🍞 , 🥛} = 1/3$$

- **Jaccard Distance**

$$d_{\mathrm{Jaccard}}(A, B) = 1 - s_{\mathrm{Jaccard}}(A, B)$$

# OVERVIEW OF CLUSTERING APPROACHES

- **Center-based:** each cluster is characterized by its center



**Centers**

# OVERVIEW OF CLUSTERING APPROACHES

- **Hierarchical:** points are organized according to a hierarchy (or tree structure)

# OVERVIEW OF CLUSTERING APPROACHES

- **Density-based:** clusters are high-density regions surrounded by low-density regions

# OUTLINE

# K-MEANS ALGORITHM: STEPS

**1. Initialization**: Pick K random
points as centers

2. Repeat:

   a) **Assignment:** assign each point to
   nearest cluster

   b) **Update:** move each cluster center
   to average of its assigned points

   **Stop** if no assignments change



**Centers**

# K-MEANS ALGORITHM: STEPS

**1. Initialization**: Pick K random points as centers

**2. Repeat:**

**a) Assignment:** assign each point to nearest cluster

**b) Update:** move each cluster center to average of its assigned points

**Stop** if no assignments change

# K-MEANS ALGORITHM: STEPS

**1. Initialization**: Pick K random points as centers

**2. Repeat:**

**a) Assignment:** assign each point to nearest cluster

**b) Update:** move each cluster center to average of its assigned points

**Stop** if no assignments change
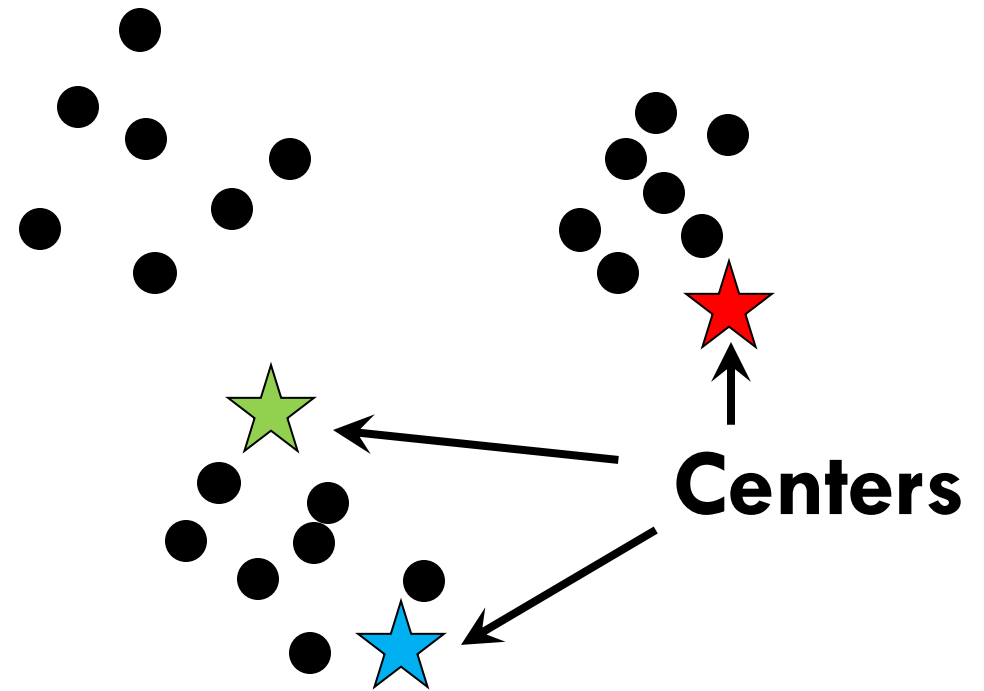
# K-MEANS ALGORITHM: STEPS

**1. Initialization**: Pick K random points as centers

**2. Repeat:**

**a) Assignment:** assign each point to nearest cluster

**b) Update:** move each cluster center to average of its assigned points

**Stop** if no assignments change

# K-MEANS ALGORITHM: STEPS

**1. Initialization**: Pick K random points as centers

**2. Repeat:**

**a) Assignment:** assign each point to nearest cluster

**b) Update:** move each cluster center to average of its assigned points

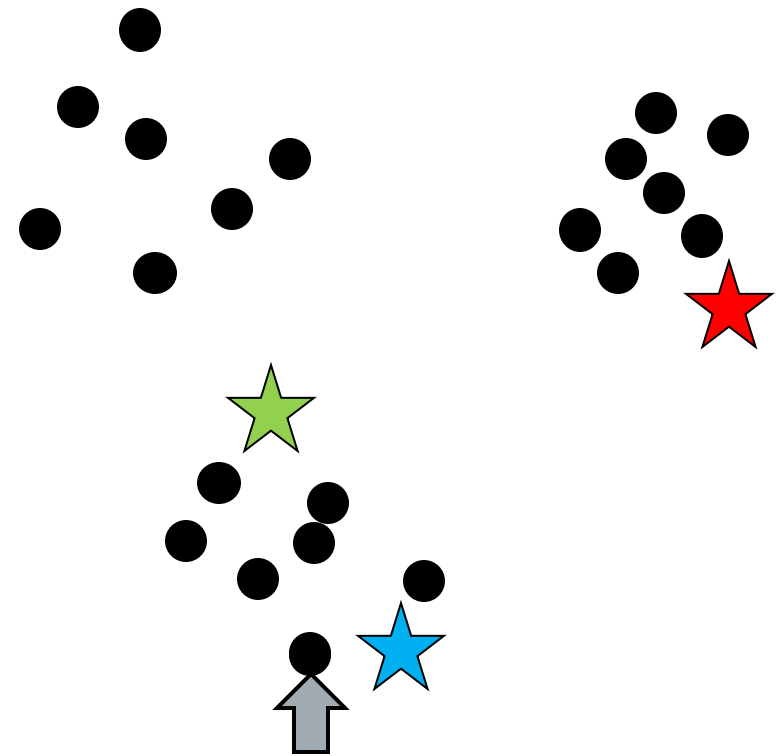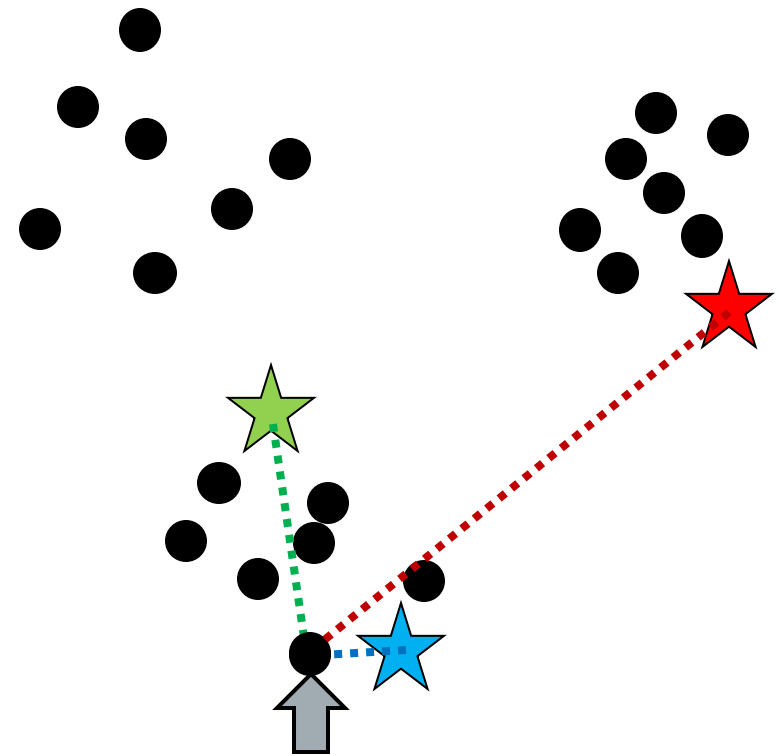**Stop** if no assignments change

# K-MEANS ALGORITHM: STEPS

**1. Initialization**: Pick K random points as centers

**2. Repeat:**

   **a) Assignment:** assign each point to nearest cluster

   **b) Update:** move each cluster center to average of its assigned points
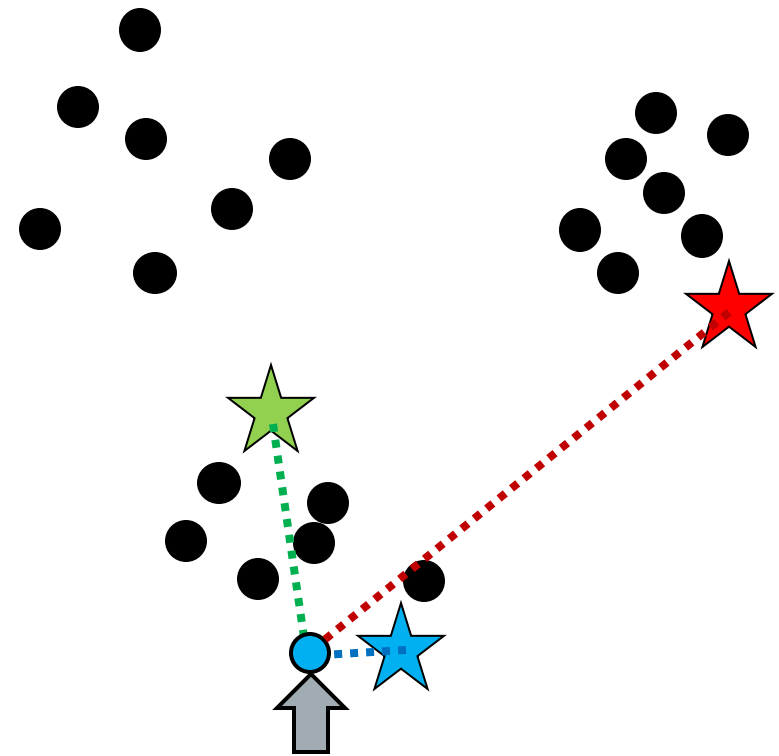
**Stop** if no assignments change

# K-MEANS ALGORITHM: STEPS

**1. Initialization**: Pick K random points as centers

**2. Repeat:**

**a) Assignment:** assign each point to nearest cluster

**b) Update:** move each cluster center to average of its assigned points
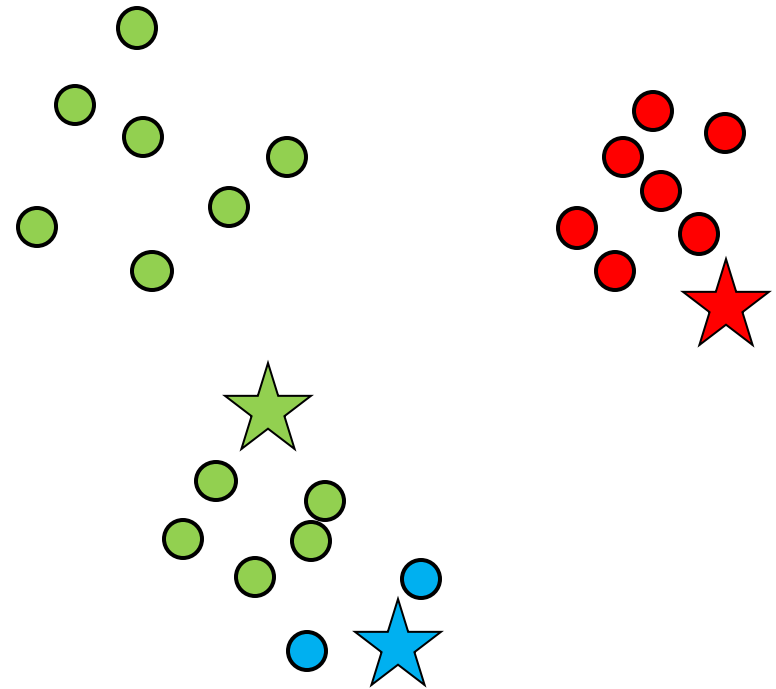
**Stop** if no assignments change
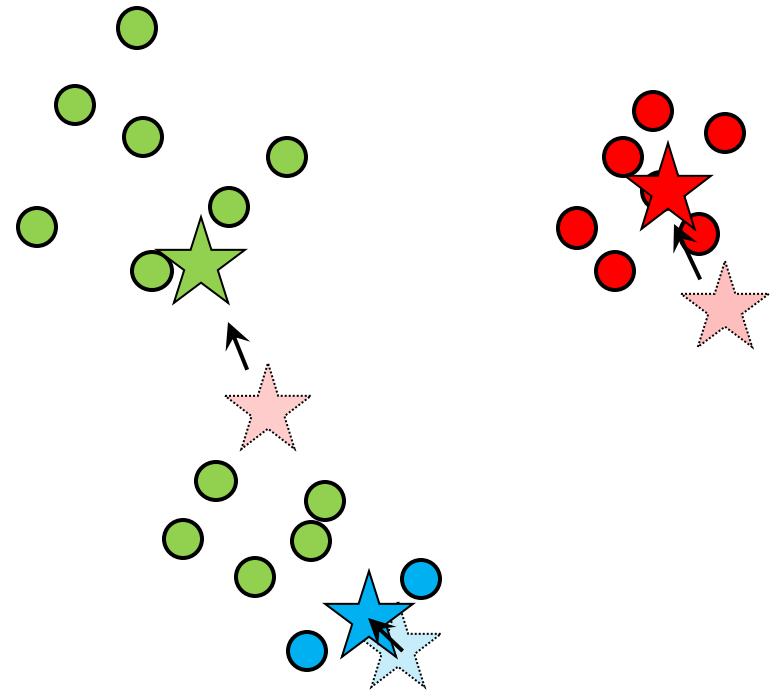
# K-MEANS ALGORITHM: STEPS

**1. Initialization**: Pick K random points as centers

**2. Repeat:**

   **a) Assignment:** assign each point to nearest cluster

   **b) Update:** move each cluster center to average of its assigned points

**Stop** if no assignments change
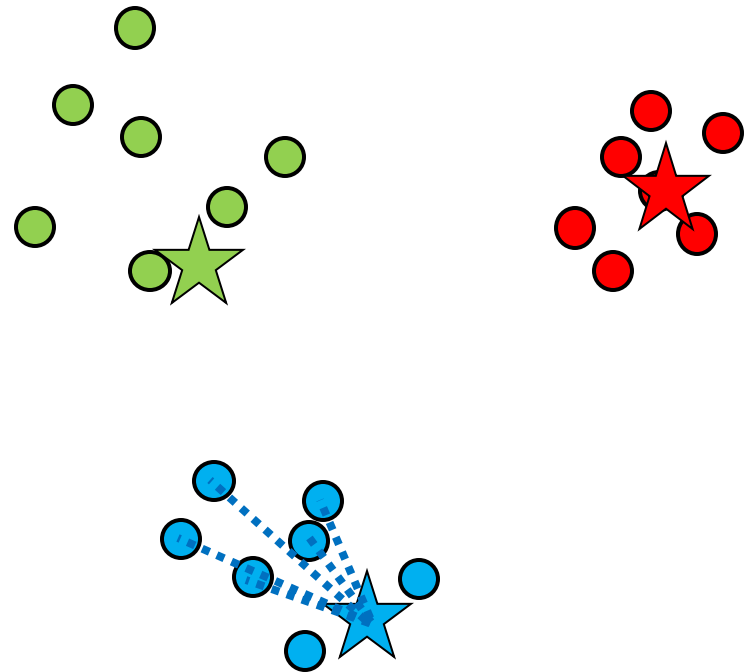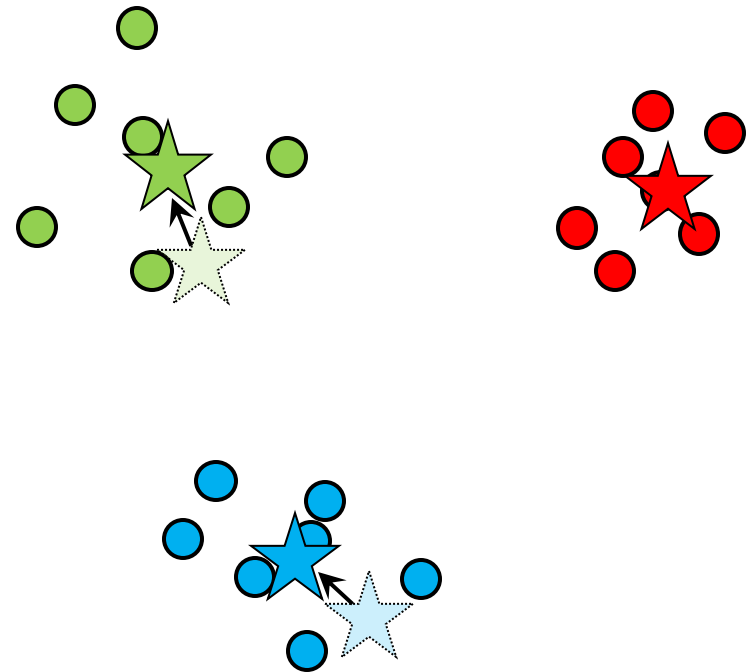
# K-MEANS ALGORITHM: STEPS

**1. Initialization**: Pick K random points as centers

**2. Repeat:**

**a) Assignment:** assign each point to nearest cluster

**b) Update:** move each cluster center to average of its assigned points
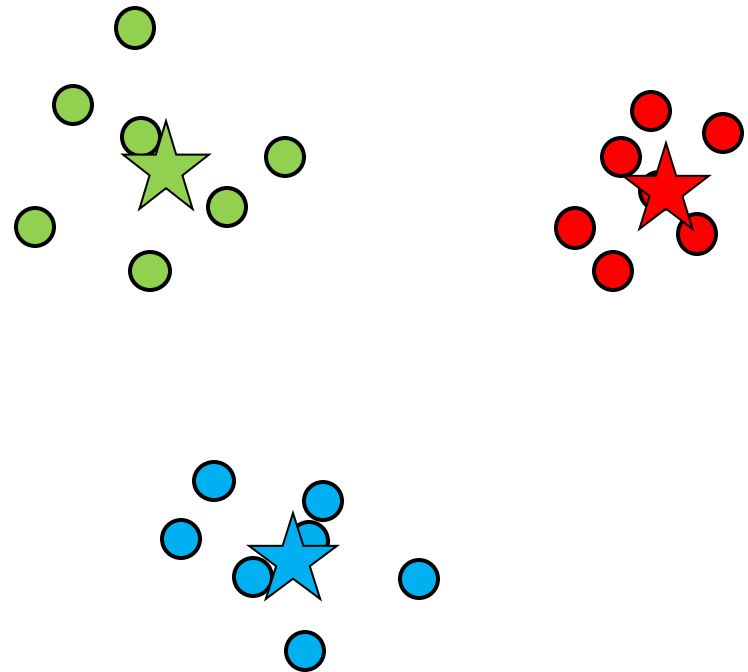
**Stop** if no assignments change
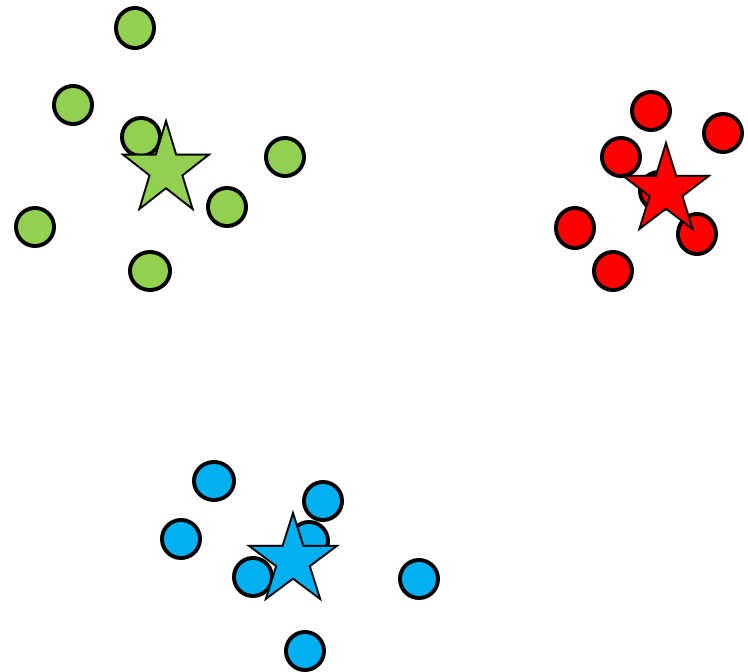
# K-MEANS ALGORITHM: STEPS

**1. Initialization**: Pick K random points as centers
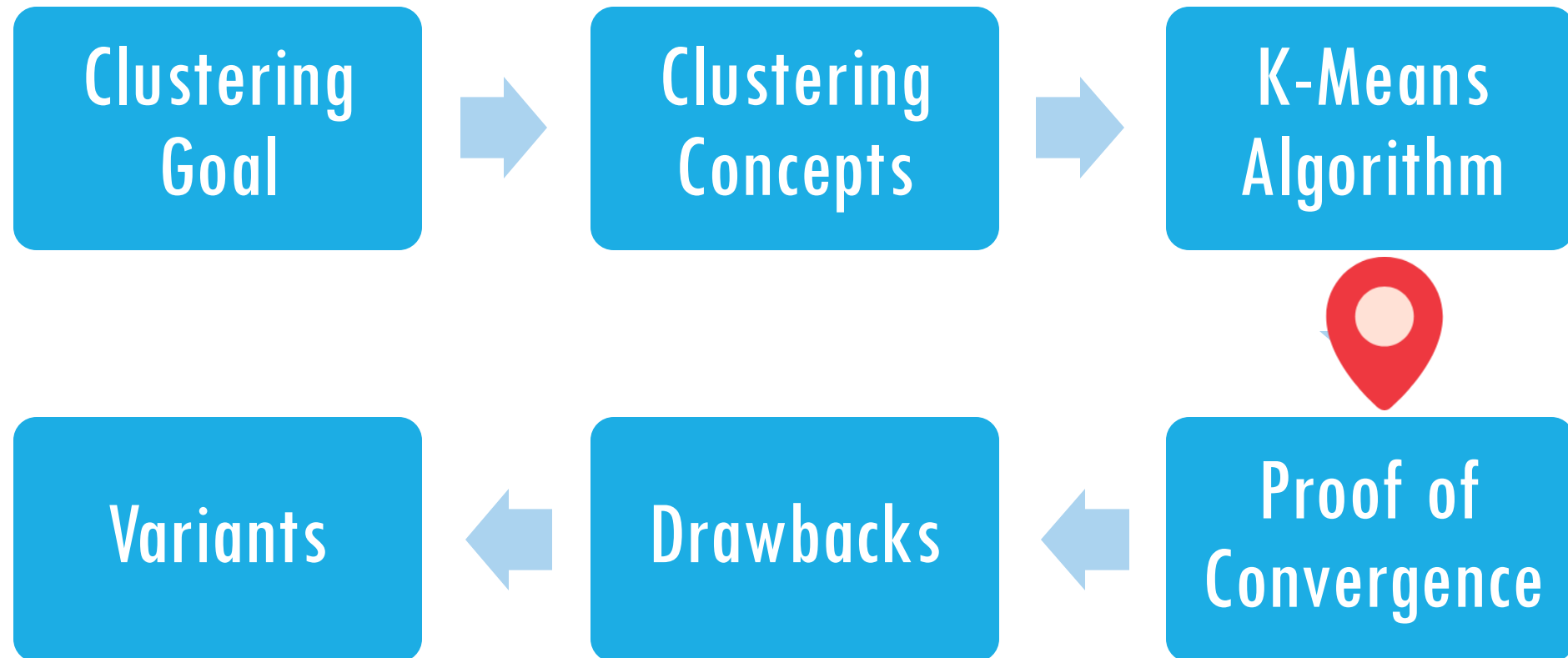
**2. Repeat:**

    **a) Assignment:** assign each point to nearest cluster

    **b) Update:** move each cluster center to average of its assigned points

**Stop** if no assignments change

# OUTLINE

Clustering Goal → Clustering Concepts → K-Means Algorithm
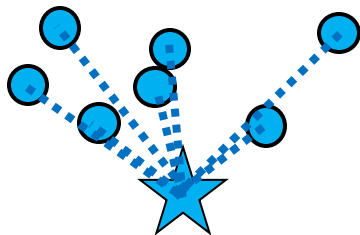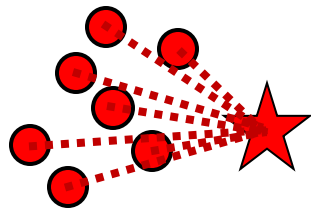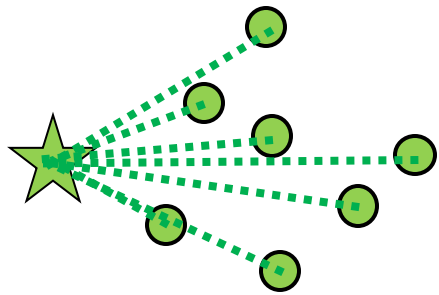
Variants ← Drawbacks ← Proof of Convergence

# OPTIMIZATION OBJECTIVE

**Within-Cluster Sum of Squares (WCSS)**: sum of squared distances between each point and its cluster center

$$\text{WCSS} = \sum_{i=1}^{K} \sum_{x \in C_i} \|x - c_i\|_2^2$$

ith cluster center

Sum over clusters

Set of points in ith cluster

Squared distance between point and center

# K-MEANS AS ALTERNATING MINIMIZATION

**1. Initialization**: Pick K random points as centers
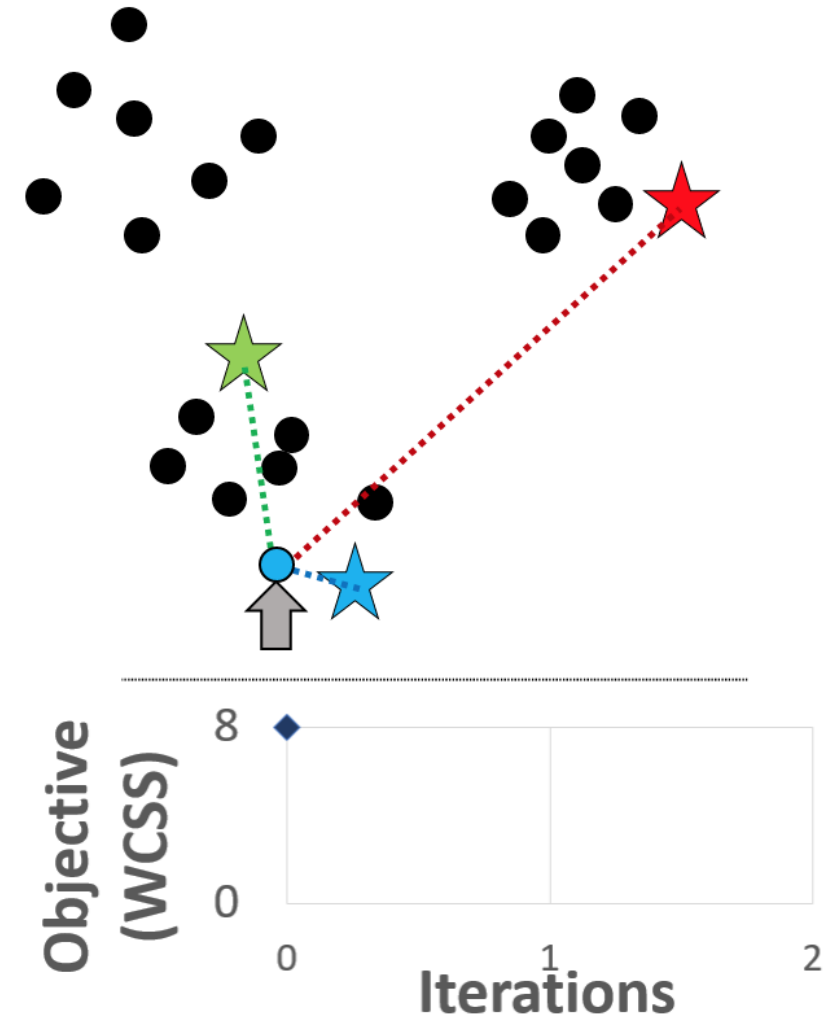
**2. Repeat:**

**a) Assignment:** assign each point to nearest cluster

$$\underset{C_1,\cdots,C_K}{\text{minimize}} \sum_{i=1}^{K} \sum_{x \in C_i} \|x - c_i\|_2^2$$

**b) Update:** move each cluster center to average of its assigned points

$$\underset{c_1,\cdots,c_K}{\text{minimize}} \sum_{i=1}^{K} \sum_{x \in C_i} \|x - c_i\|_2^2$$

**Stop** if no assignments change

# K-MEANS AS ALTERNATING MINIMIZATION

**1. Initialization**: Pick K random points as centers
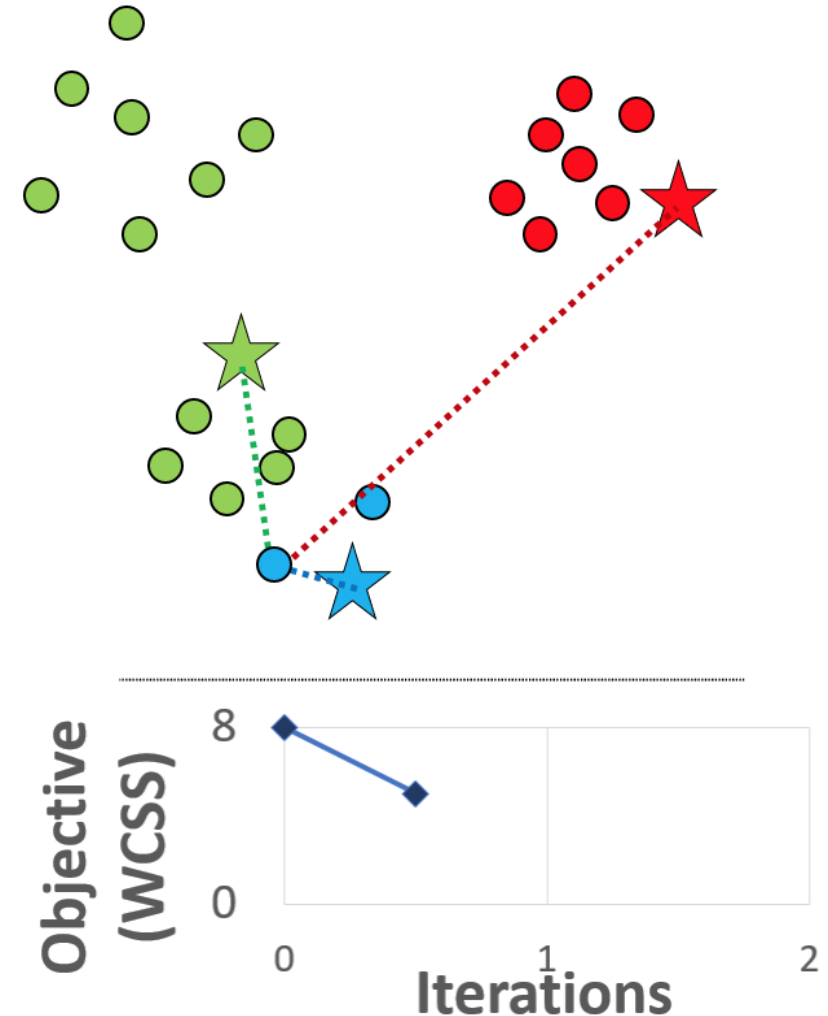
**2. Repeat:**

   **a) Assignment:** assign each point to nearest cluster

$$\underset{C_1,\cdots,C_K}{\text{minimize}} \sum_{i=1}^{K} \sum_{x \in C_i} \|x - c_i\|_2^2$$

**b) Update:** move each cluster center to average of its assigned points

$$\underset{c_1,\cdots,c_K}{\text{minimize}} \sum_{i=1}^{K} \sum_{x \in C_i} \|x - c_i\|_2^2$$

**Stop** if no assignments change

# K-MEANS AS ALTERNATING MINIMIZATION

**1. Initialization**: Pick K random points as centers

**2. Repeat:**

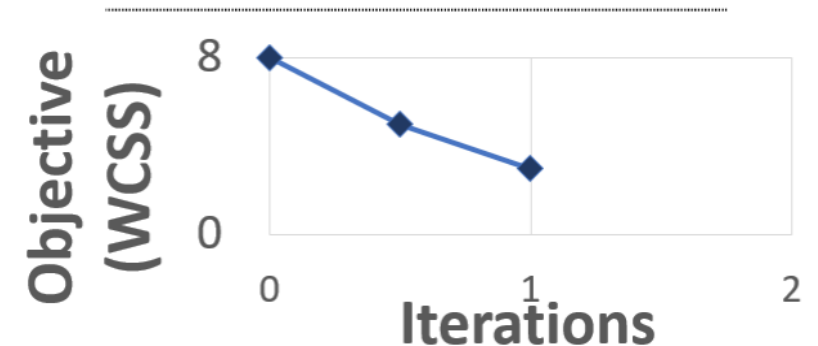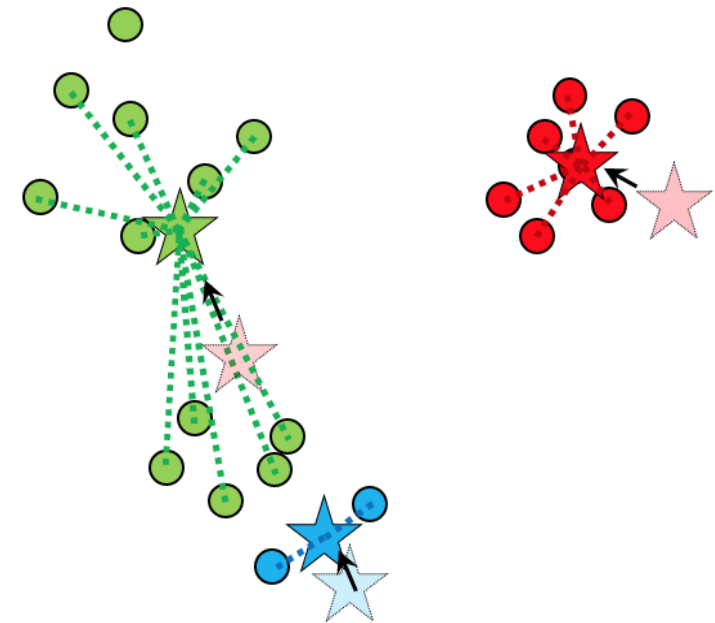  **a) Assignment:** assign each point to nearest cluster

$$\underset{C_1,\cdots,C_K}{\text{minimize}} \sum_{i=1}^{K} \sum_{x \in C_i} \|x - c_i\|_2^2$$

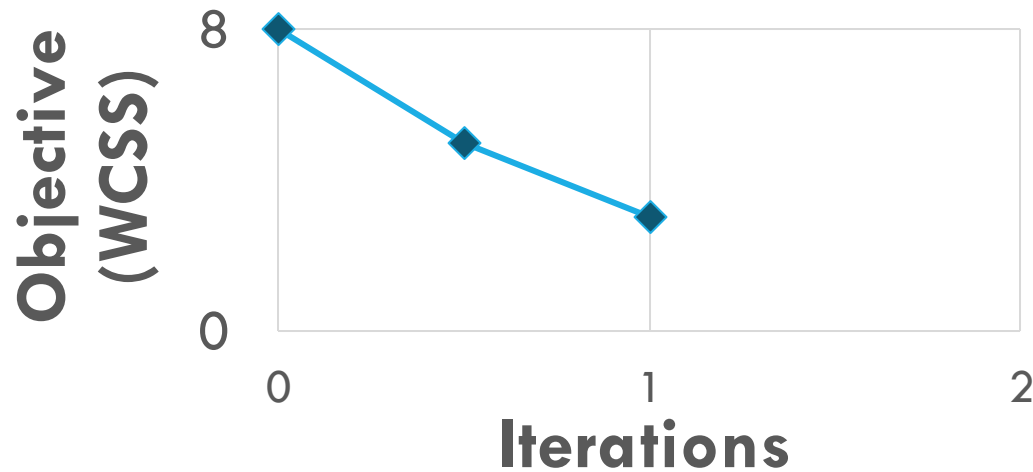  **b) Update:** move each cluster center to average of its assigned points

$$\underset{c_1,\cdots,c_K}{\text{minimize}} \sum_{i=1}^{K} \sum_{x \in C_i} \|x - c_i\|_2^2$$
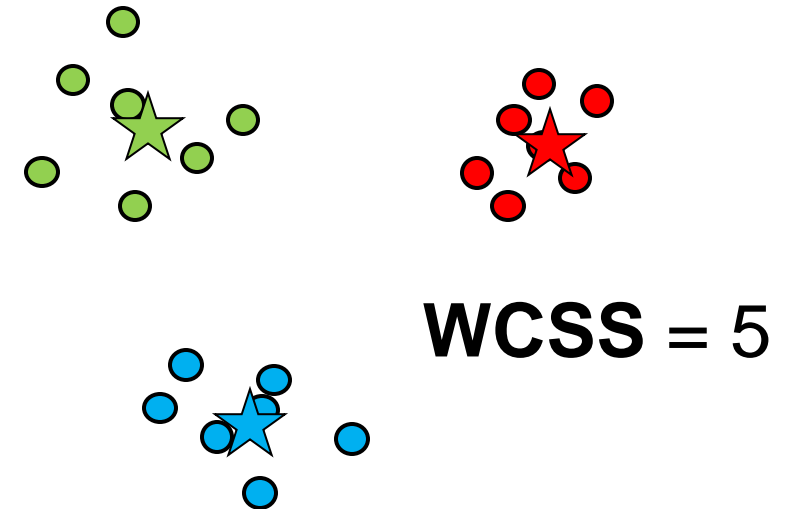
**Stop** if no assignments change

# PROOF OF CONVERGENCE
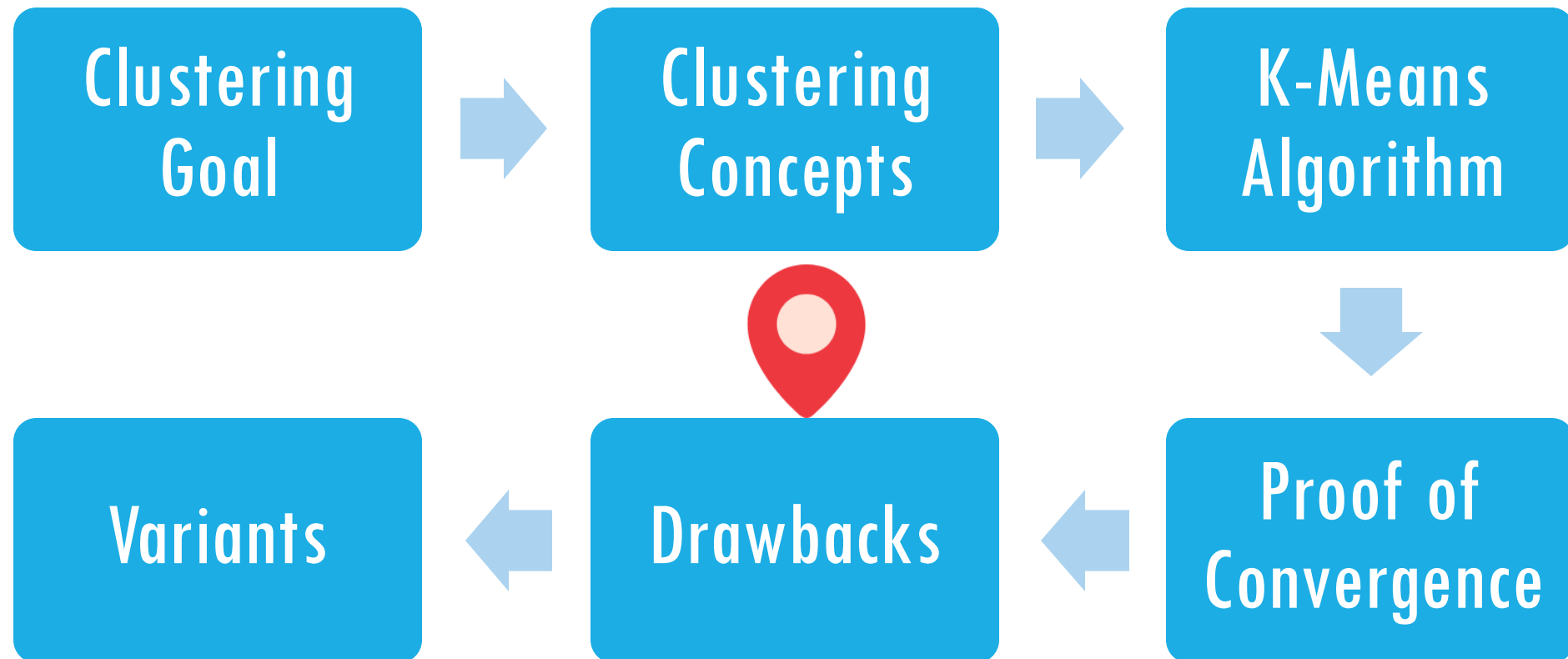
1. The WCSS objective is strictly decreasing



2. There are a finite number of possible clusterings



**WCSS** = 5

➡ **The algorithm must eventually stop!**
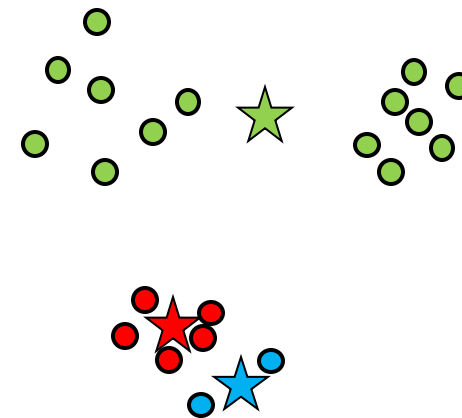
# OUTLINE

# LOCAL, NOT GLOBAL OPTIMUM

- The algorithm only returns a **local**, not a **global** optimum!
  - (Finding the global optimum is NP-hard)
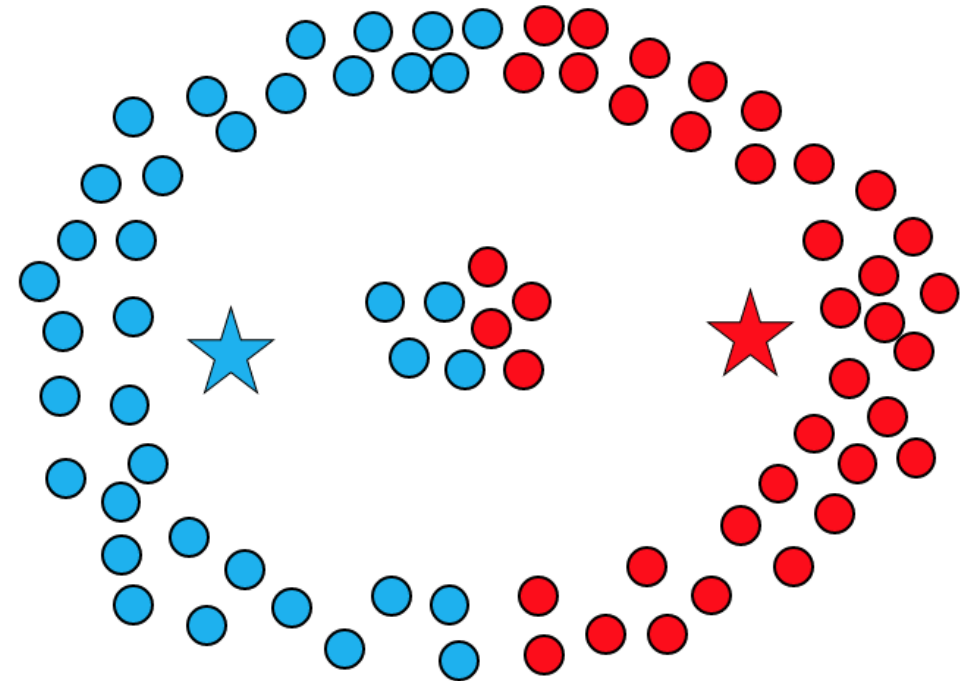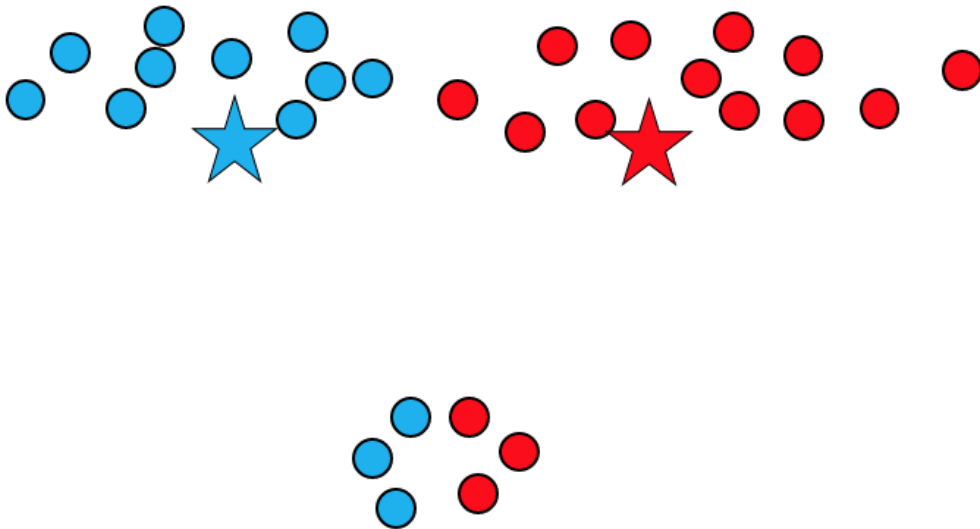
- Initialization is important

Local minima

Global minimum
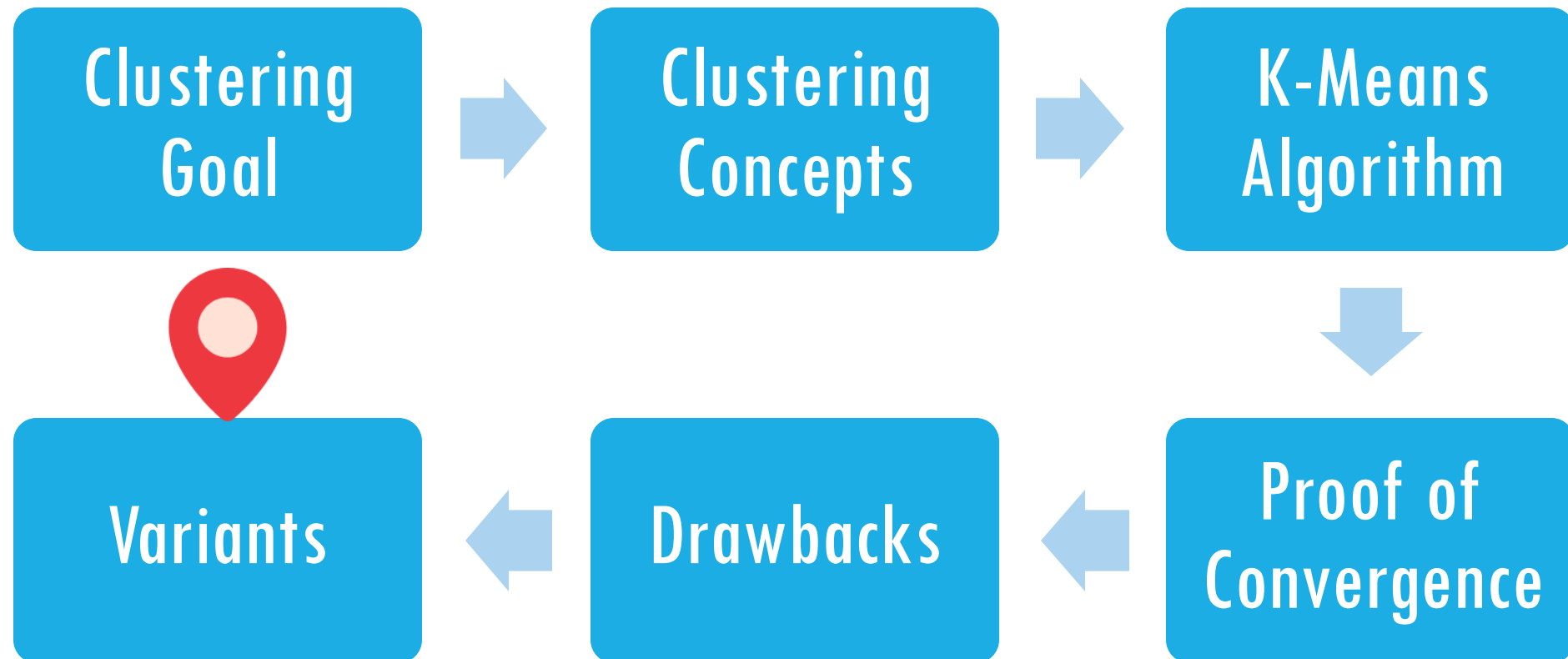
**Example of a local minima**

# NON-SPHERE-LIKE CLUSTERS

- Optimization objective results in roughly sphere shaped clusters
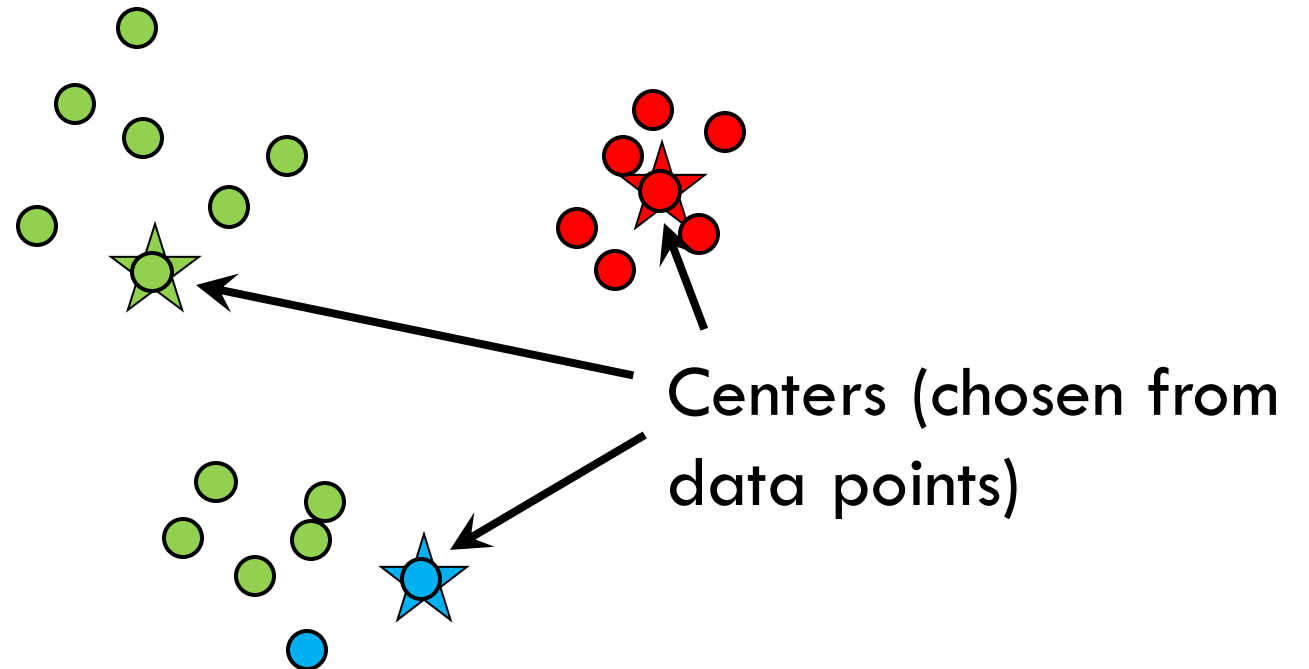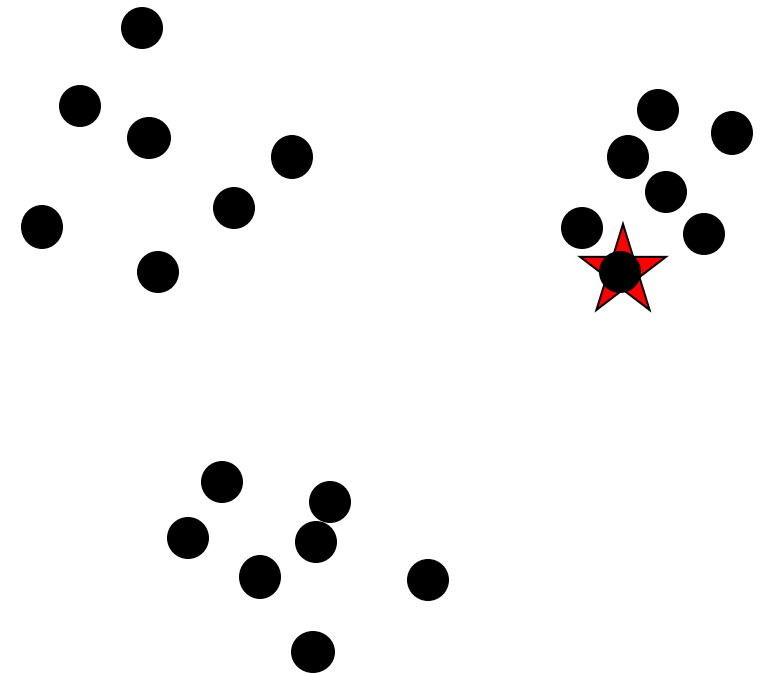
# OUTLINE

# 1. K-MEDOIDS ALGORITHM

- **K-Medoids:** like K-Means, but centers are chosen from data points
- Useful when:
  - We want data points as cluster representatives
  - Complex data types – we can only measure distances between data points
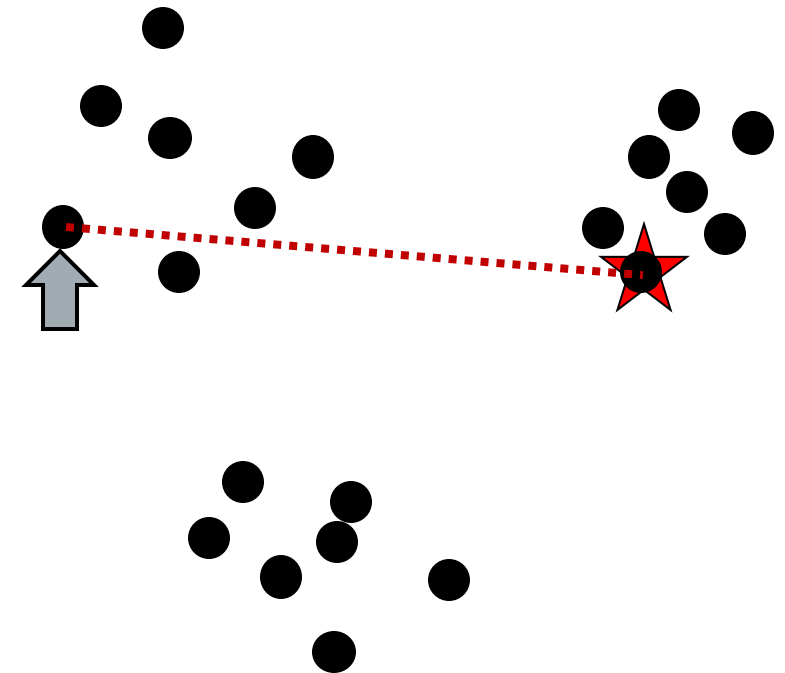
Centers (chosen from data points)

# 2. K-MEANS++ ALGORITHM

- **K-Means++:** only changes the initialization step

- **"Spread out centers":**
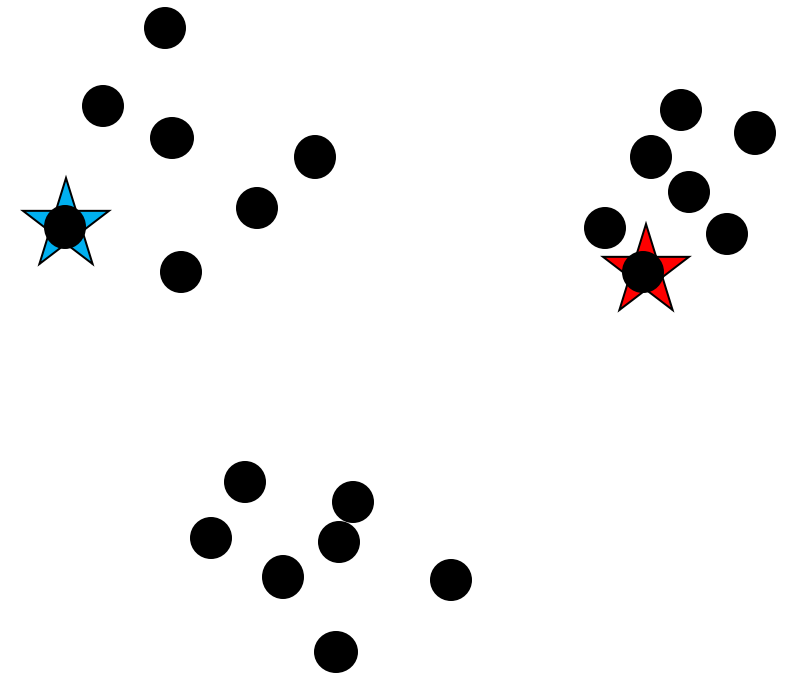    - First center is a uniformly random point

# 2. K-MEANS++ ALGORITHM

- **K-Means++:** only changes the initialization step

- **"Spread out centers":**
  - First center is a uniformly random point
  - Next centers: each point chosen with probability proportional to square of distance to its closest center

# 2. K-MEANS++ ALGORITHM

- **K-Means++:** only changes the initialization step

- **"Spread out centers":**
  - First center is a uniformly random point
  - Next centers: each point chosen with probability proportional to square of distance to its closest center
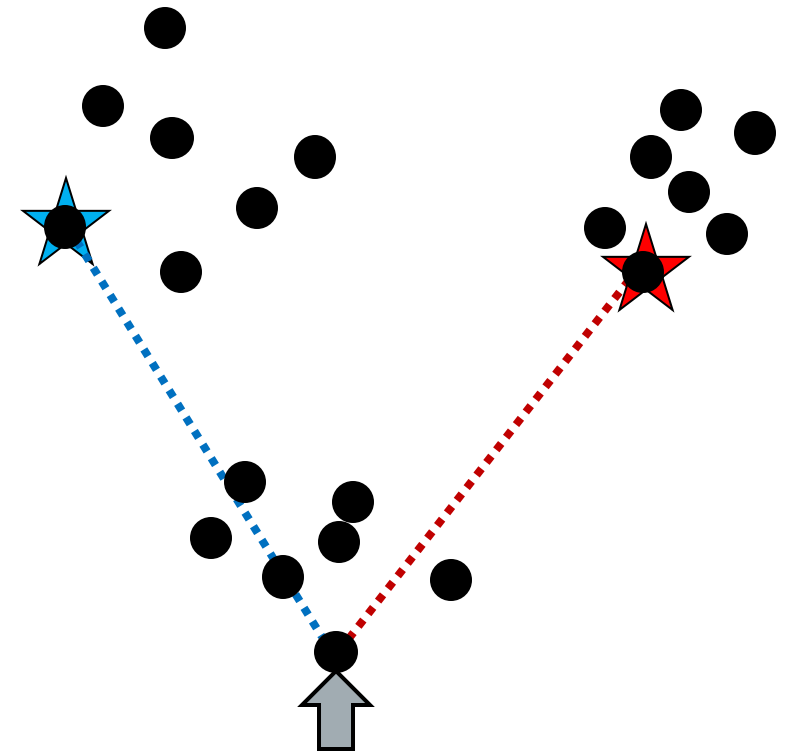
# 2. K-MEANS++ ALGORITHM

- **K-Means++:** only changes the initialization step

- **"Spread out centers":**
  - First center is a uniformly random point
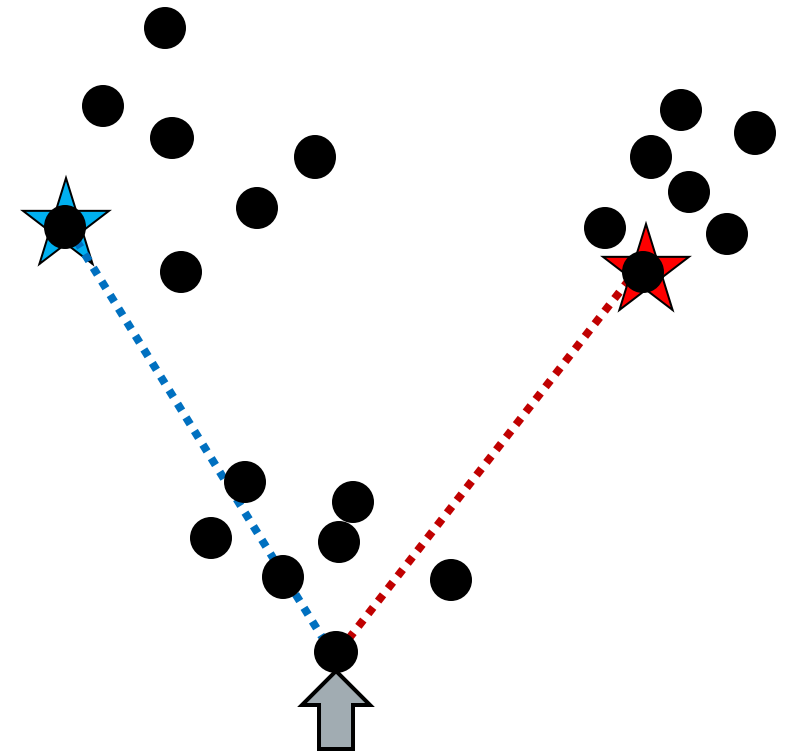  - Next centers: each point chosen with probability proportional to square of distance to its closest center
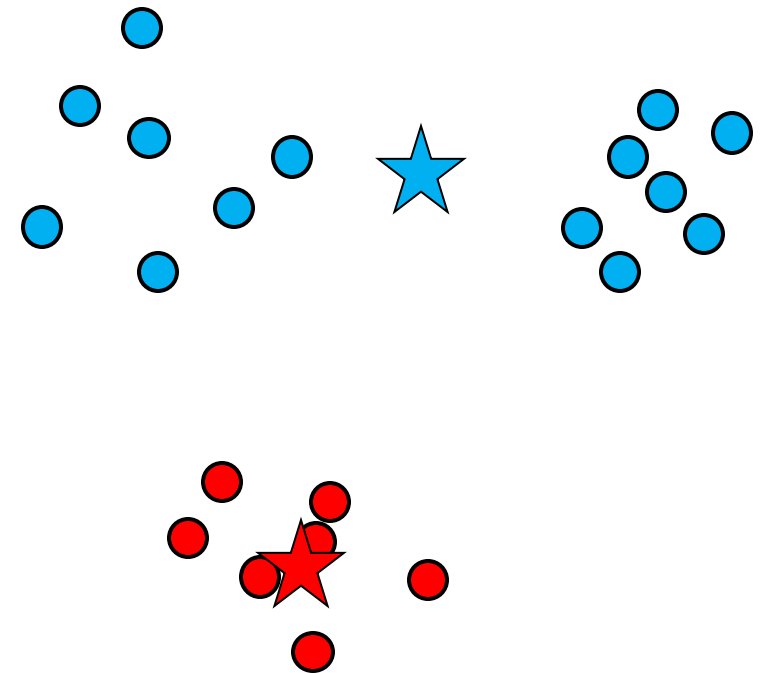
# 2. K-MEANS++ ALGORITHM

- **K-Means++:** only changes the initialization step

- Better practical performance

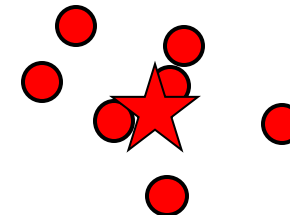- Theoretical guarantee: O(log k) approximation ratio in expectation

# 3. X-MEANS ALGORITHM
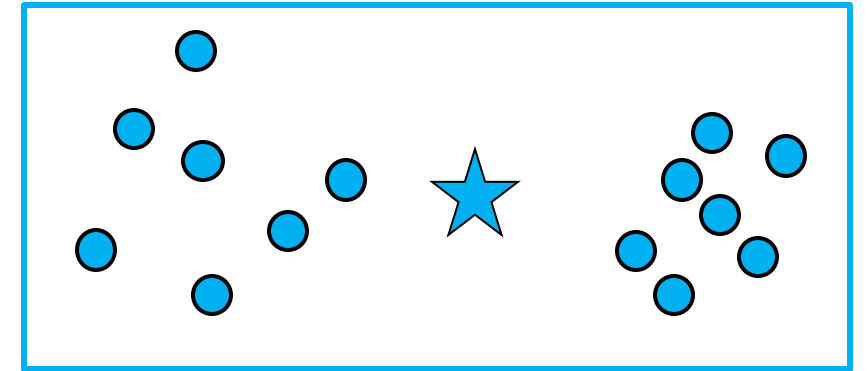
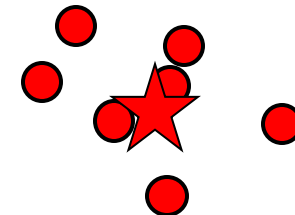- **Automatic way to choose K**

1. Run usual K-Means with K=2

# 3. X-MEANS ALGORITHM

- **Automatic way to choose K**

1. Run usual K-Means with K=2

2. Attempt to split each cluster by running K-Means with K=2 only within that cluster

# 3. X-MEANS ALGORITHM

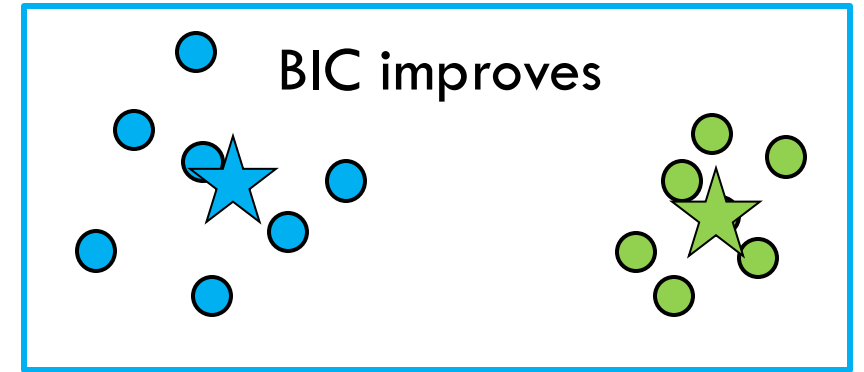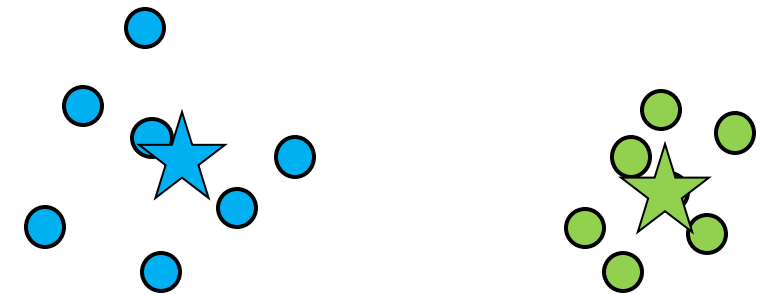- **Automatic way to choose K**

1. Run usual K-Means with K=2
2. Attempt to split each cluster by running K-Means with K=2 only within that cluster
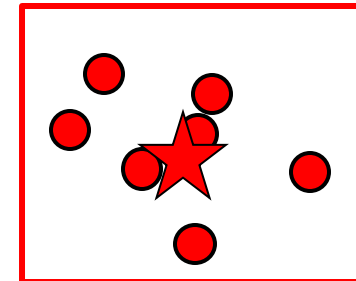   - Use "Bayesian Information Criterion" (BIC) to decide whether to split

BIC improves

# 3. X-MEANS ALGORITHM

- **Automatic way to choose K**

1. Run usual K-Means with K=2

2. Attempt to split each cluster by running K-Means with K=2 only within that cluster

   - Use "Bayesian Information Criterion" (BIC) to decide whether to split

BIC does not improve

Default

4 activities