

Machine Learning for Data Science (CS4786)

Lecture 13

Clustering + K-Means Clustering

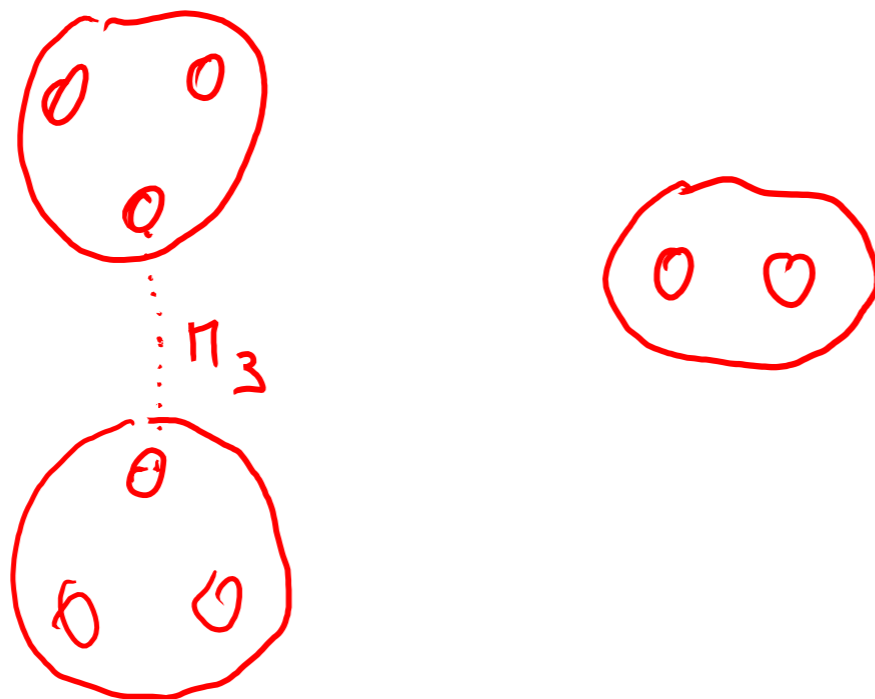
Thorsten Joachims

What does Single-Link HAC optimize?

- Single-link HAC finds the clustering C_1, \dots, C_K that maximizes the following objective:

$$M_3 = \min_{x_s, x_t: c(x_s) \neq c(x_t)} \text{dissimilarity}(x_s, x_t)$$

[Maximize smallest between-cluster distance]



Observation

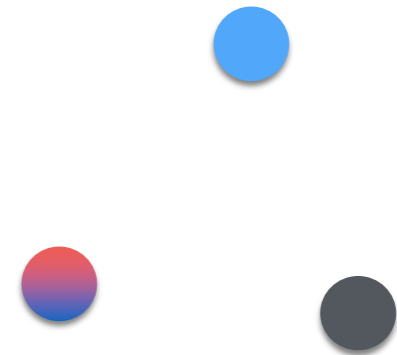
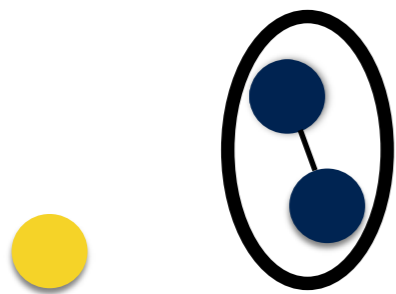
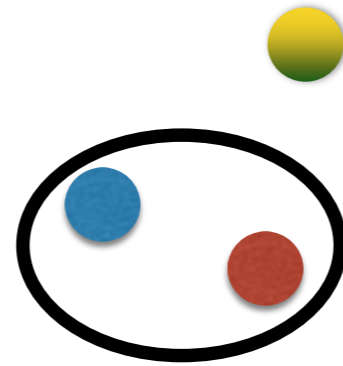
Say c is the clustering produced by single-link HAC. Then it always holds

$$\min_{\mathbf{x}_s, \mathbf{x}_t: c(\mathbf{x}_s) \neq c(\mathbf{x}_t)} \text{dissimilarity}(\mathbf{x}_s, \mathbf{x}_t) > \text{merged distances in tree}$$

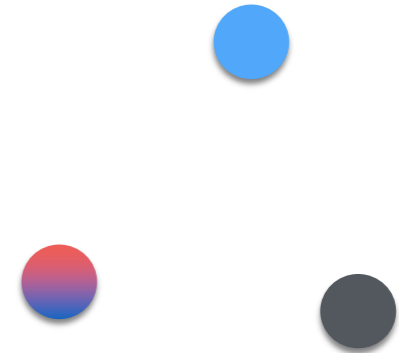
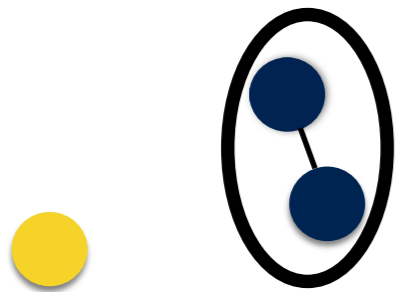
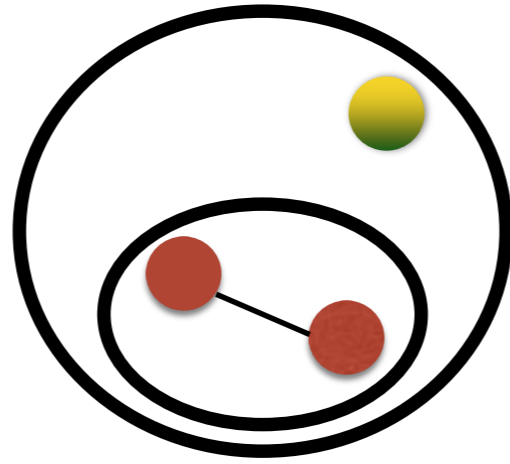
Demo



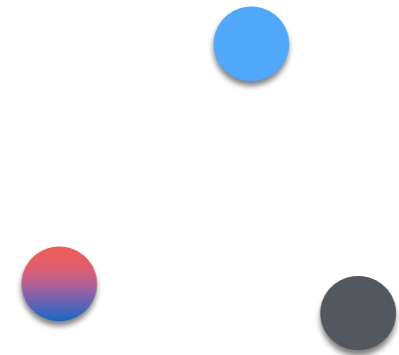
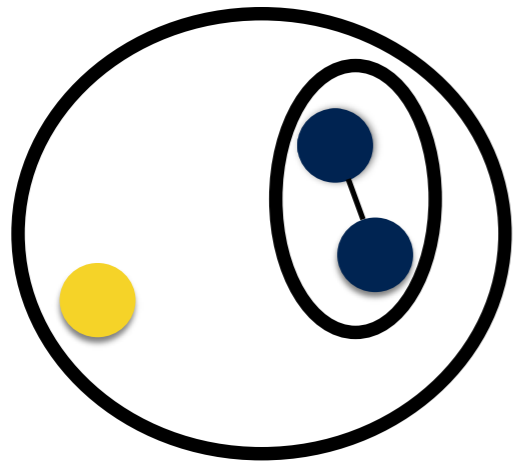
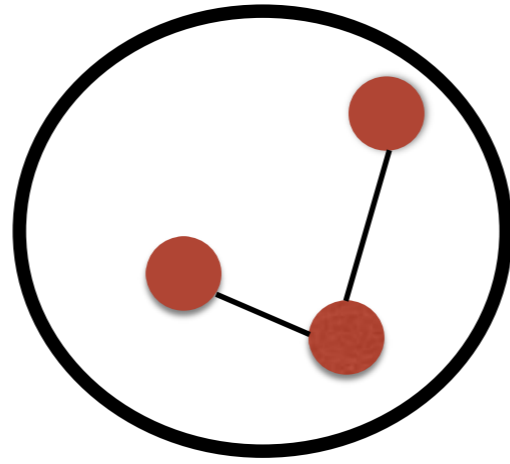
Demo



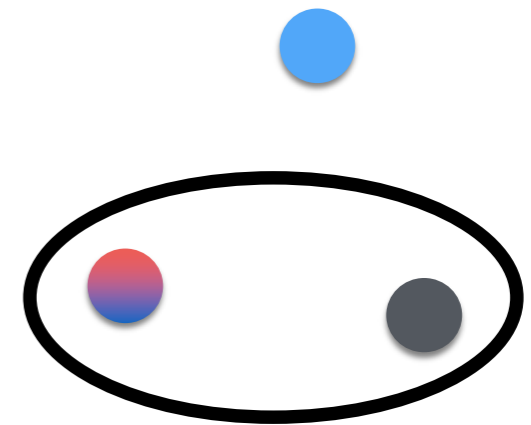
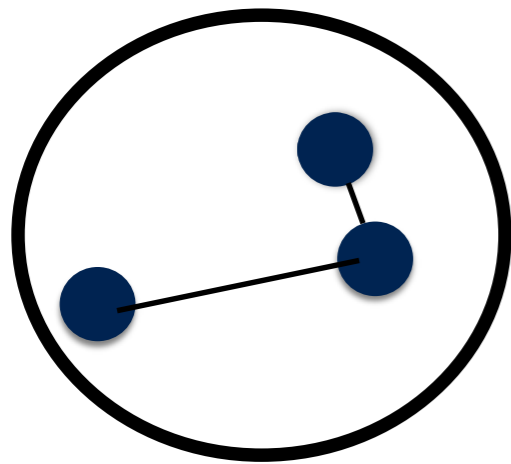
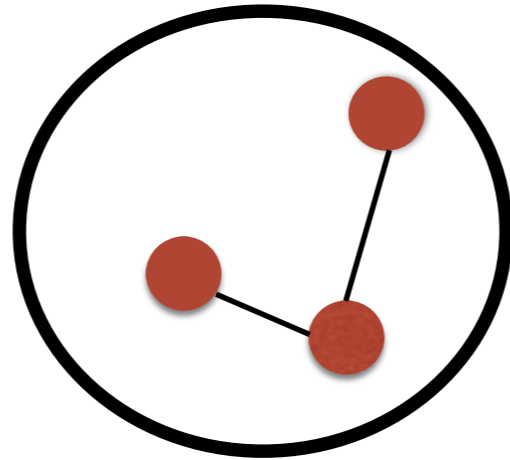
Demo



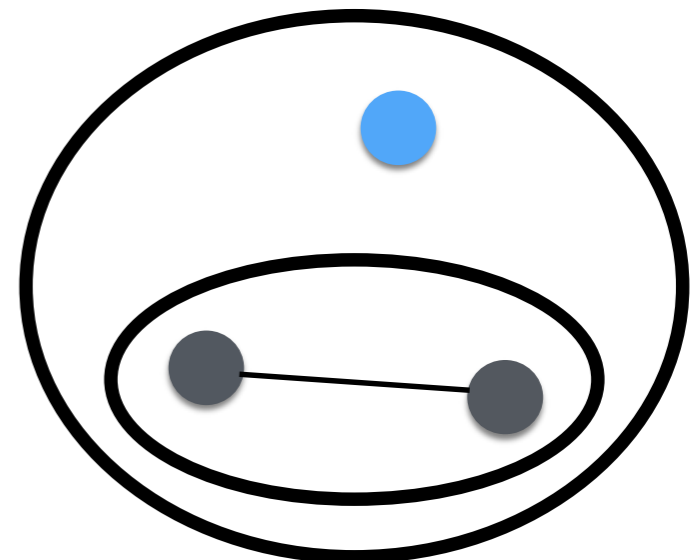
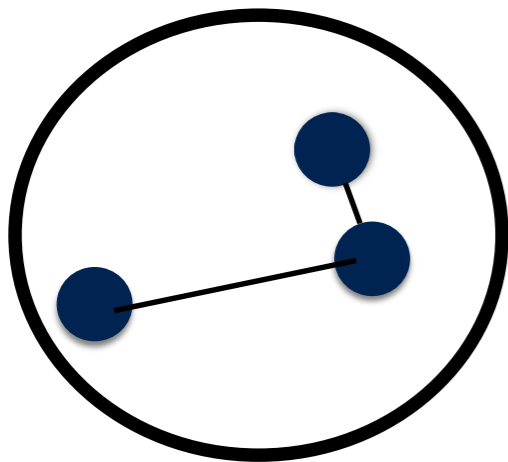
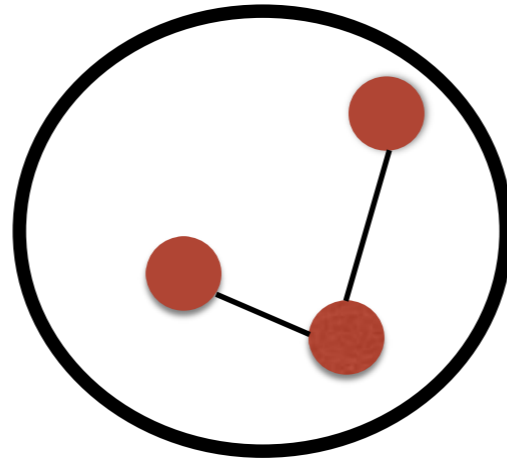
Demo



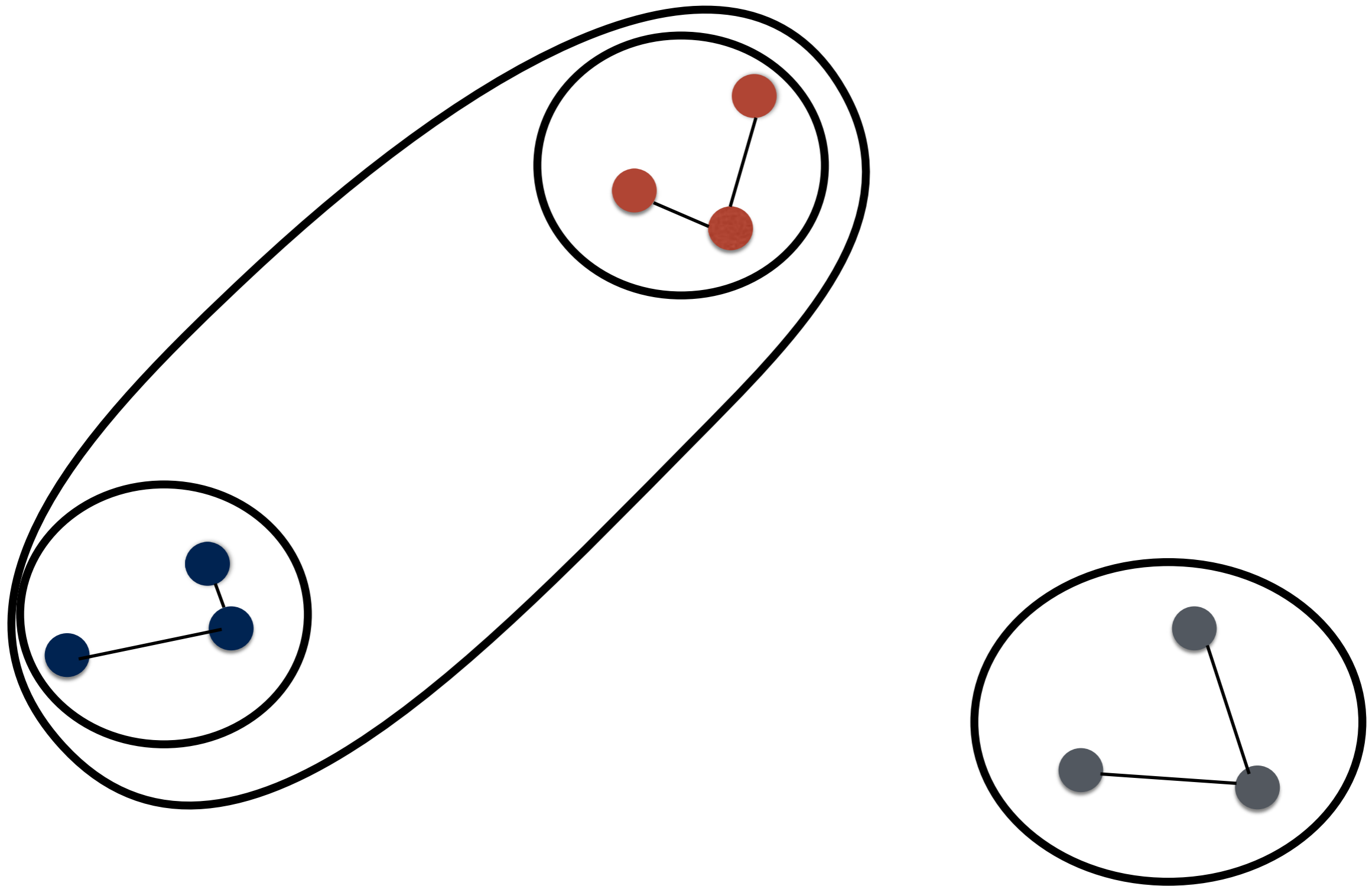
Demo



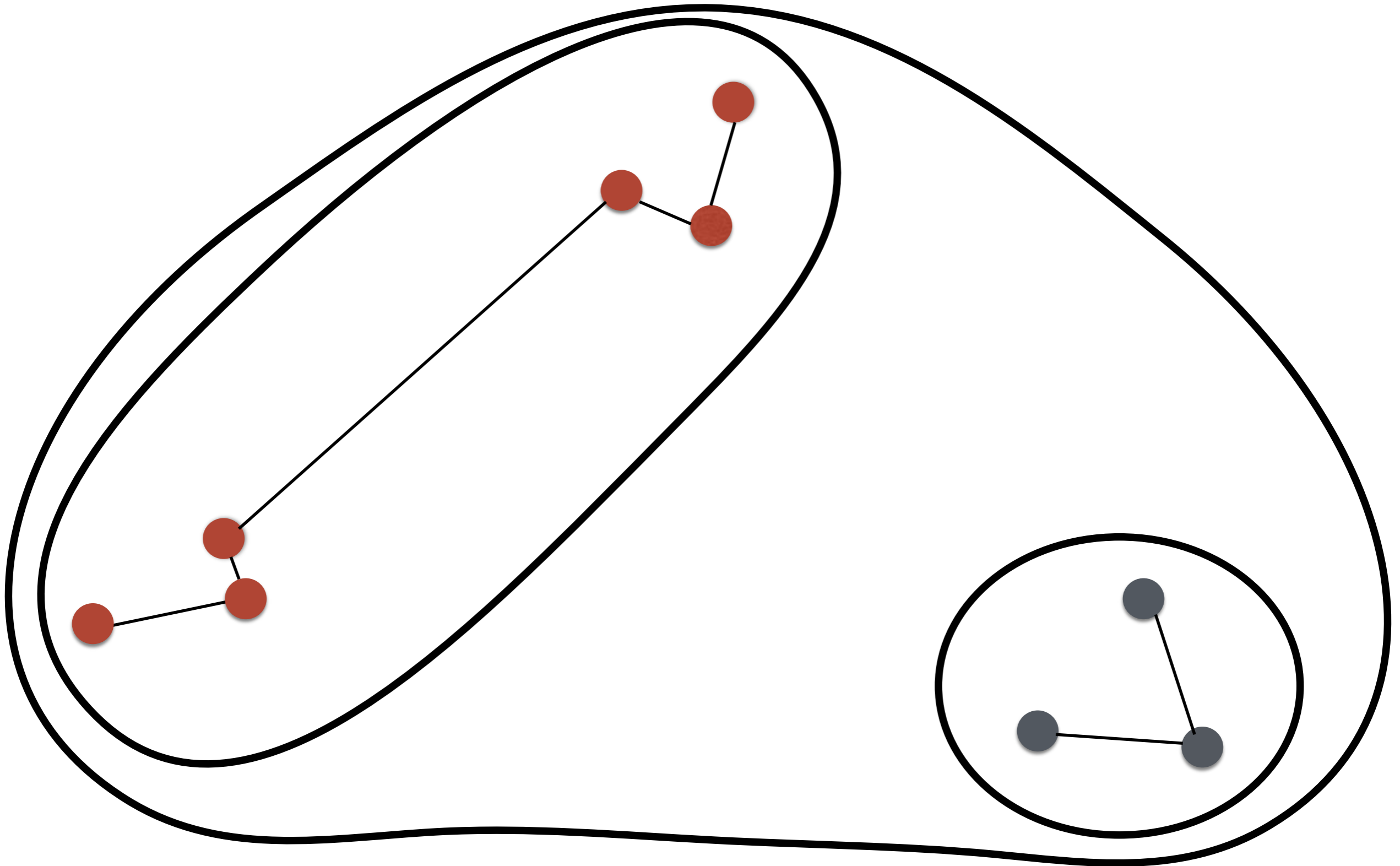
Demo



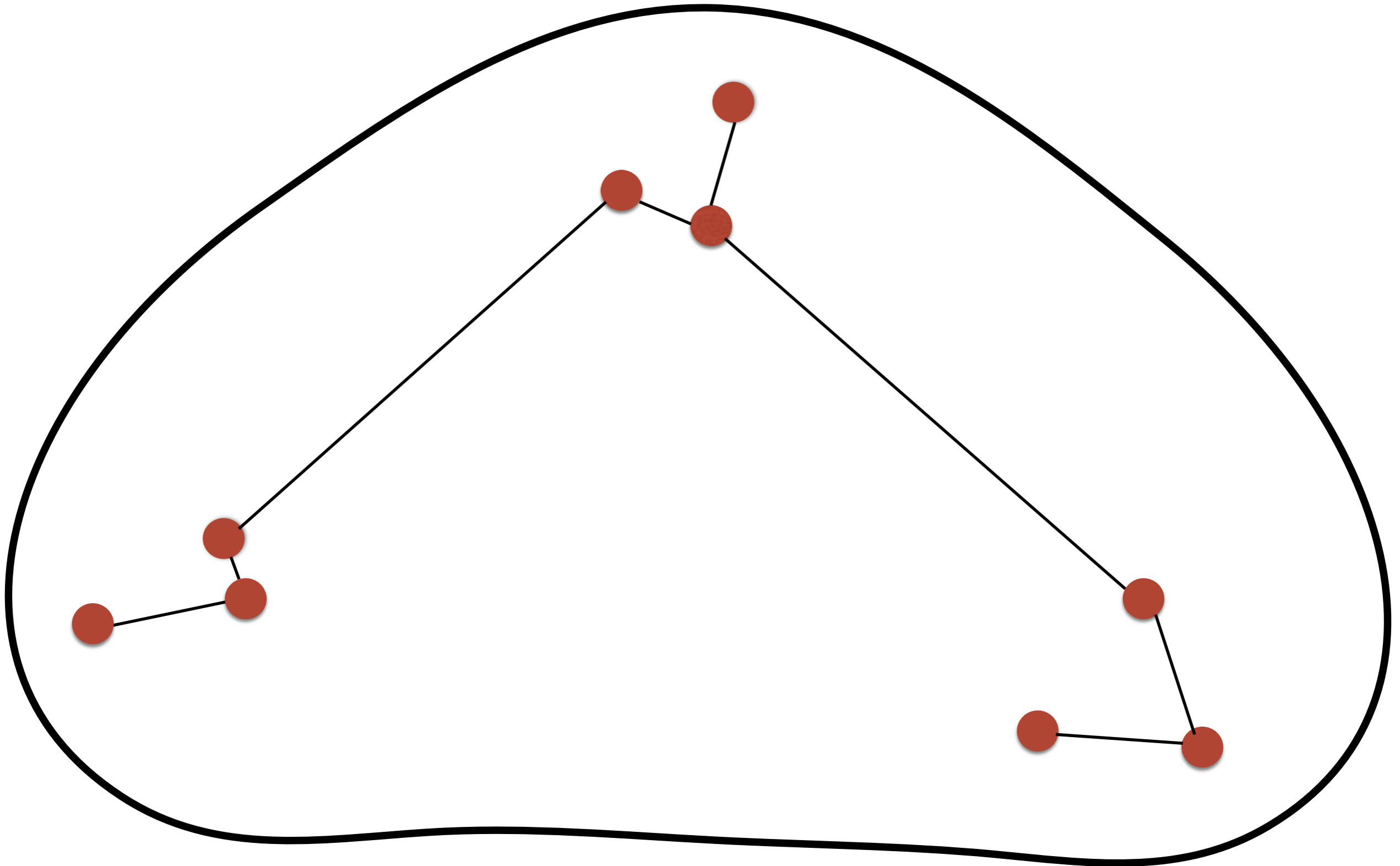
Demo



Demo



Demo



What does Single-Link HAC optimize?

- Single-link HAC finds the clustering C_1, \dots, C_K that maximizes the following objective:

$$M_3 = \min_{\mathbf{x}_s, \mathbf{x}_t: c(\mathbf{x}_s) \neq c(\mathbf{x}_t)} \text{dissimilarity}(\mathbf{x}_s, \mathbf{x}_t)$$

[Maximize smallest between-cluster distance]

SINGLE LINK OBJECTIVE

Proof:

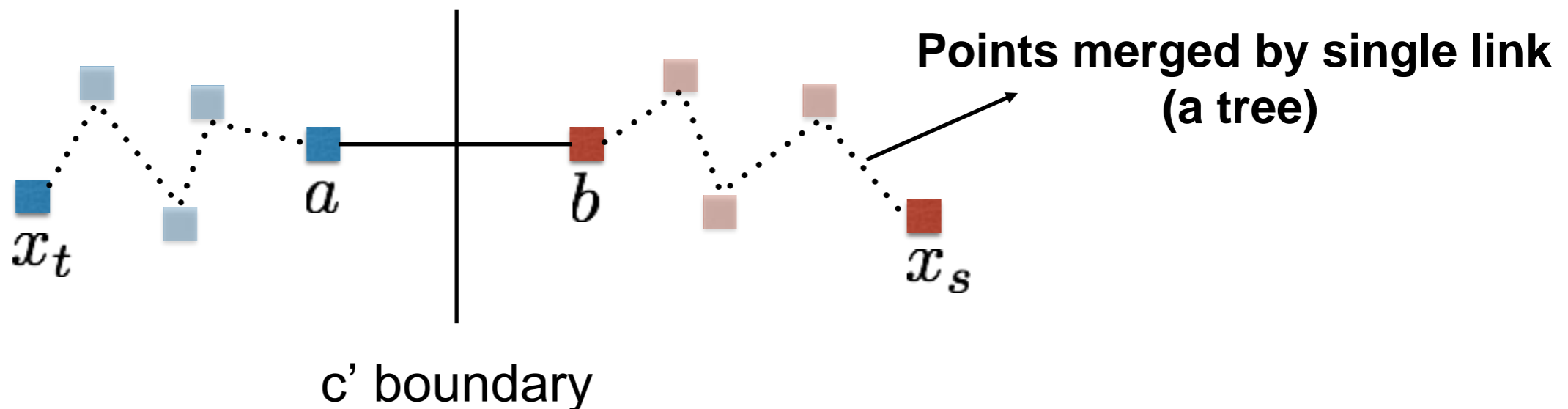
Say c is solution produced by single-link clustering

Key observation:

$\min_{t,s:c(x_t) \neq c(x_s)} \text{dissimilarity}(x_t, x_s) > \text{Merged distances in tree}$

Say $c' \neq c$ then,

$\exists t, s$ s.t. $c'(x_t) \neq c'(x_s)$ but $c(x_t) = c(x_s)$



Summary

- Clustering: Find partitioning of the data points.
- HAC: Repeatedly merge clusters bottom up.
 - Design: Distance measure and merge criterion.
 - Efficiency: $O(n^2 \log n)$
- Single-link optimizes M_3 criterion.

K-MEANS CLUSTERING

- For all $j \in [K]$, initialize cluster centroids $\hat{\mathbf{r}}_j^0$ randomly and set $m = 1$
- Repeat until convergence (or until patience runs out)
 - 1 For each $t \in \{1, \dots, n\}$, set cluster identity of the point

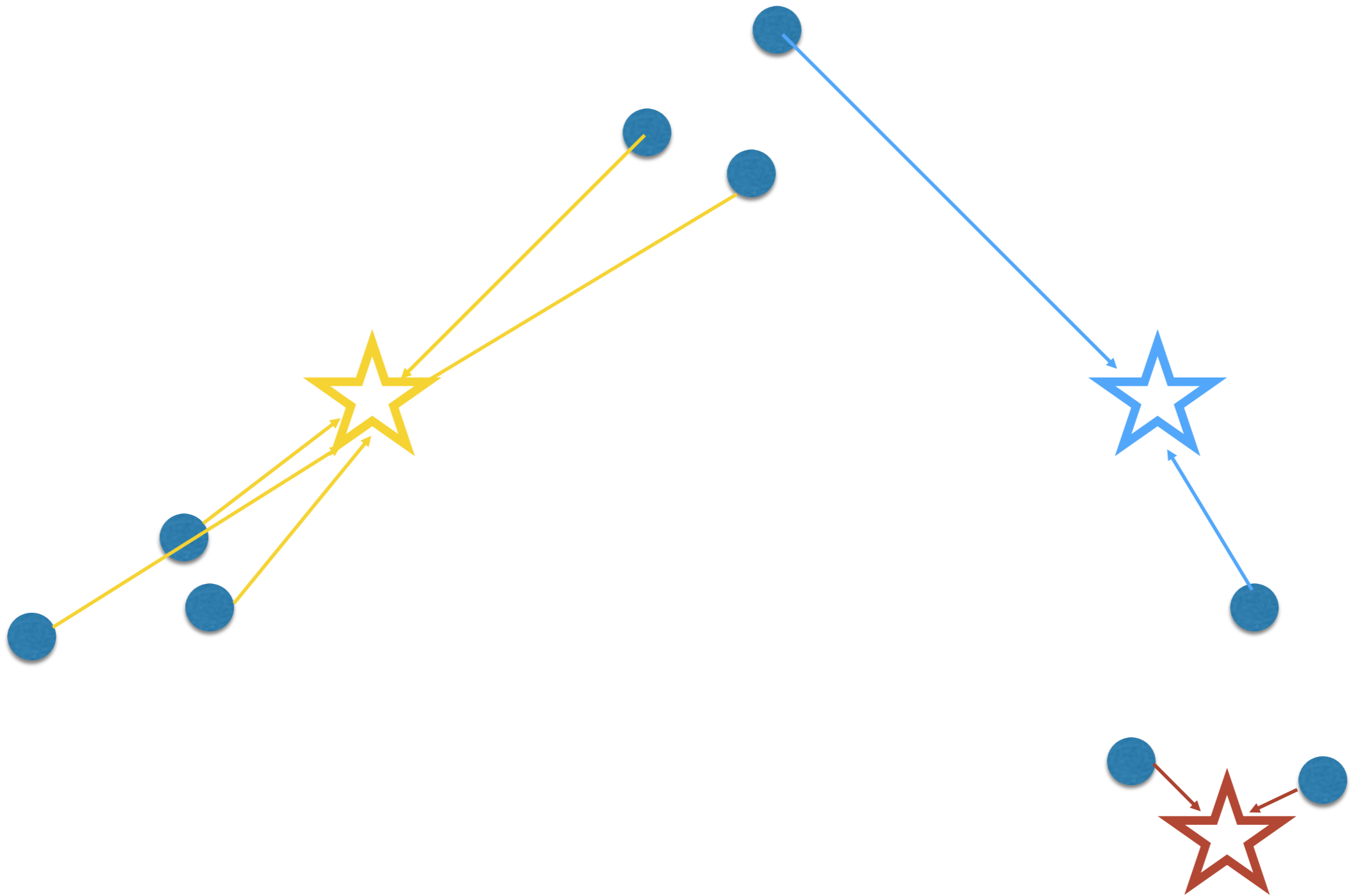
$$\hat{c}^m(\mathbf{x}_t) = \operatorname{argmin}_{j \in [K]} \|\mathbf{x}_t - \hat{\mathbf{r}}_j^{m-1}\|$$

- 2 For each $j \in [K]$, set new representative as

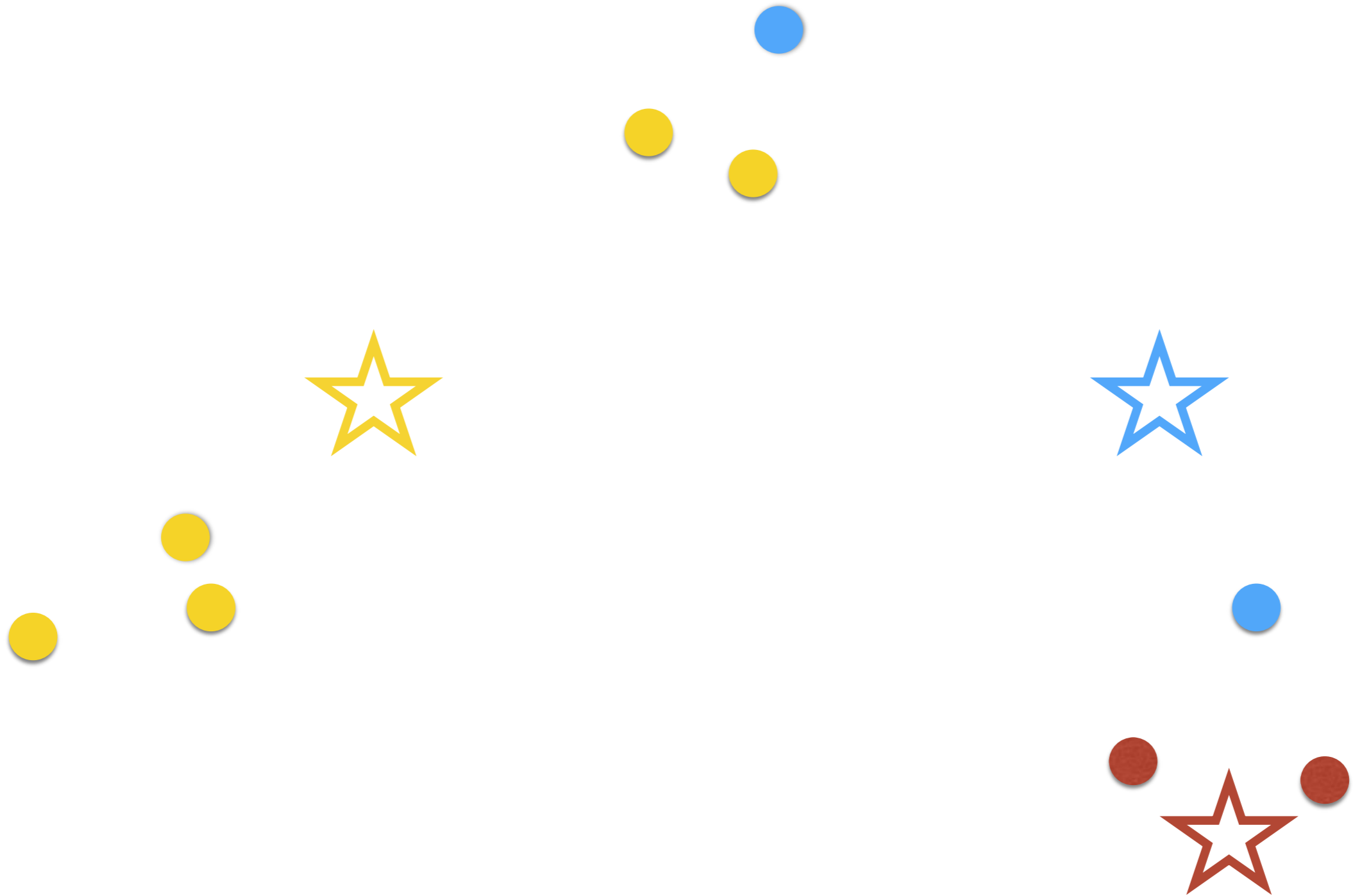
$$\hat{\mathbf{r}}_j^m = \frac{1}{|\hat{C}_j^m|} \sum_{\mathbf{x}_t \in \hat{C}_j^m} \mathbf{x}_t$$

- 3 $m \leftarrow m + 1$

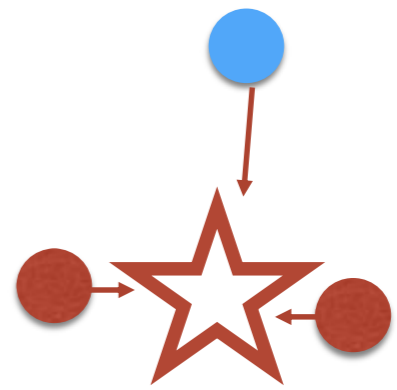
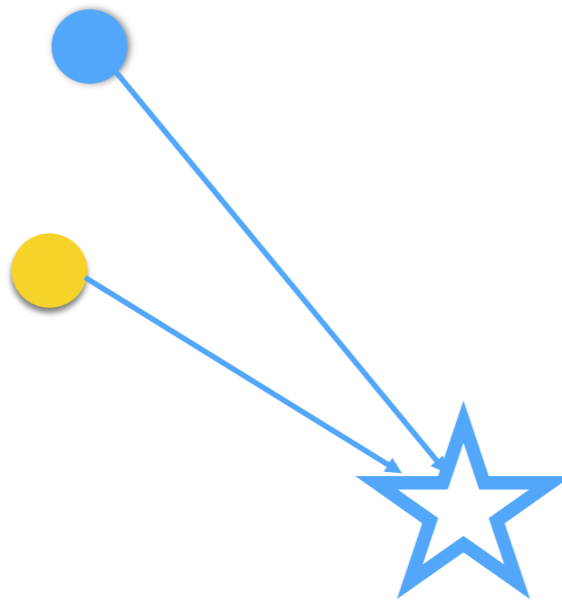
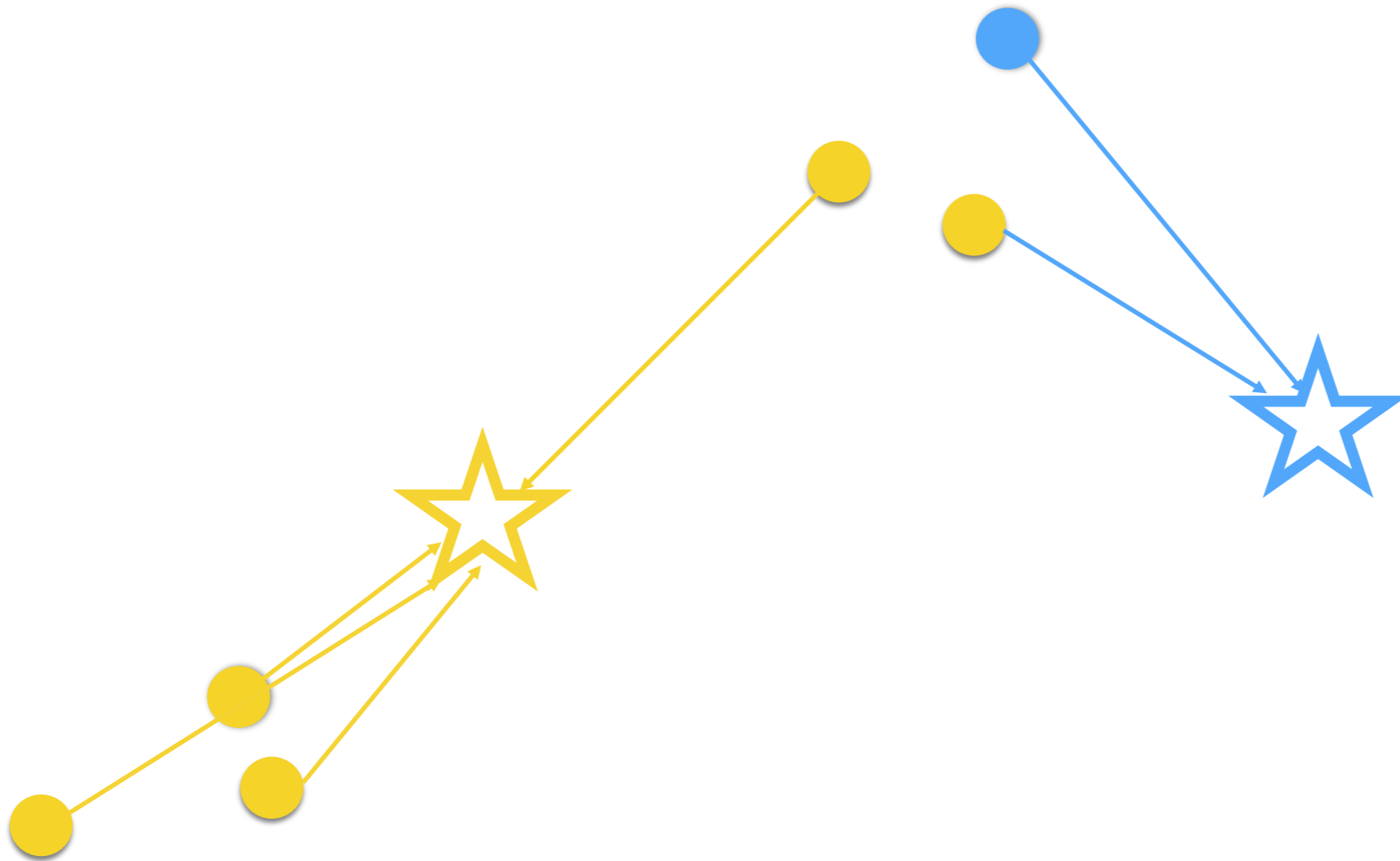
Demo



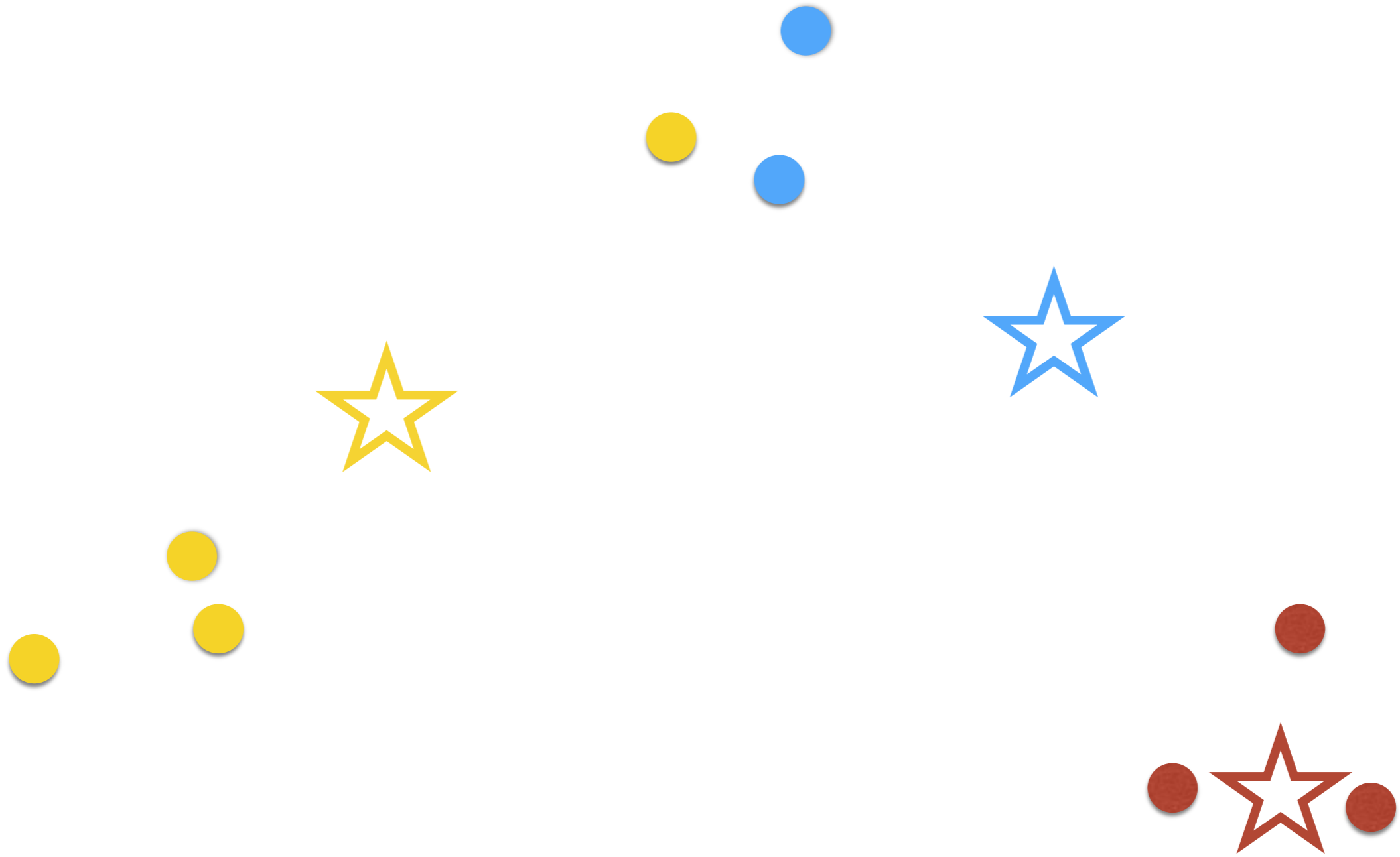
Demo



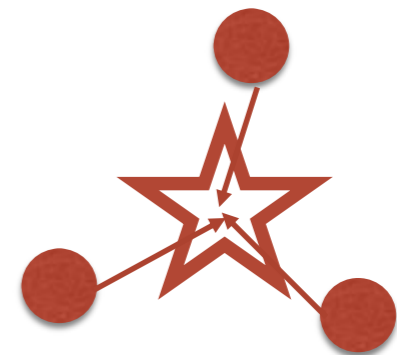
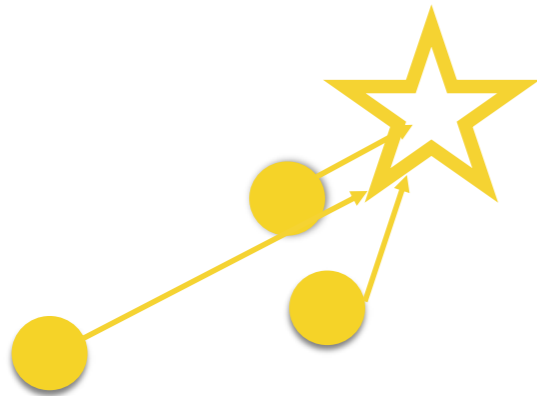
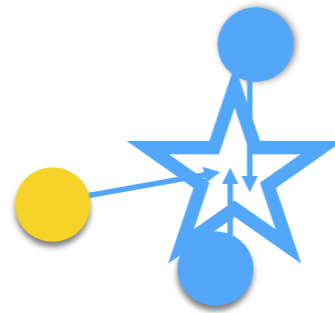
Demo



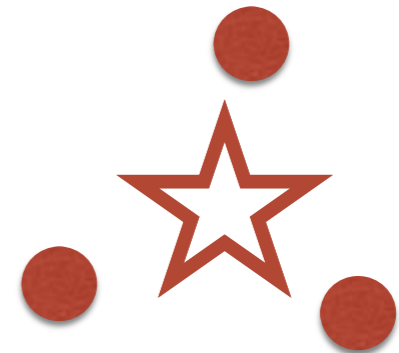
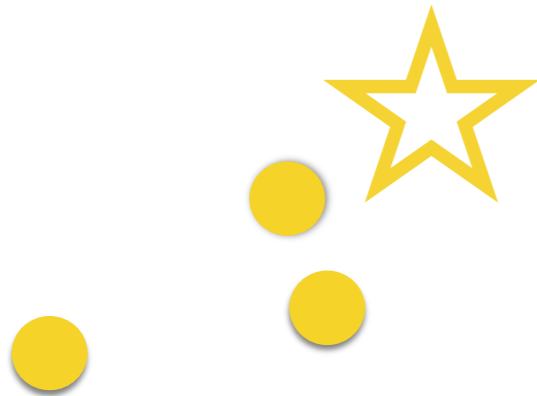
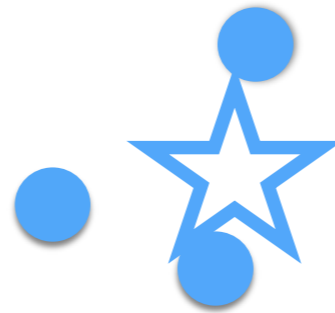
Demo



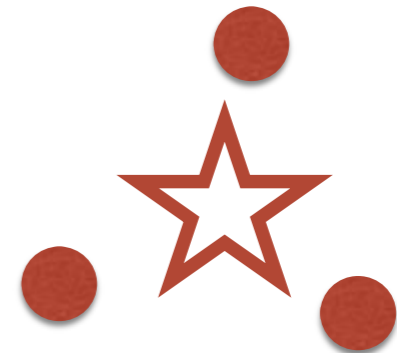
Demo



Demo



Demo



K-MEANS CLUSTERING

- For all $j \in [K]$, initialize cluster centroids $\hat{\mathbf{r}}_j^0$ randomly and set $m = 1$
- Repeat until convergence (or until patience runs out)
 - 1 For each $t \in \{1, \dots, n\}$, set cluster identity of the point

$$\hat{c}^m(\mathbf{x}_t) = \operatorname{argmin}_{j \in [K]} \|\mathbf{x}_t - \hat{\mathbf{r}}_j^{m-1}\|$$

- 2 For each $j \in [K]$, set new representative as

$$\hat{\mathbf{r}}_j^m = \frac{1}{|\hat{C}_j^m|} \sum_{\mathbf{x}_t \in \hat{C}_j^m} \mathbf{x}_t$$

- 3 $m \leftarrow m + 1$

Question: Does the K-Means algorithm always terminate, or can it run into an infinite loop?

K-MEANS CONVERGENCE

- K-means algorithm converges to local minima of objective

$$O(c; \mathbf{r}_1, \dots, \mathbf{r}_K) = \sum_{j=1}^K \sum_{c(\mathbf{x}_t)=j} \|\mathbf{x}_t - \mathbf{r}_j\|_2^2$$



- Proof:

$$\hat{c}^m(x_t) = \underset{j \in \{1..K\}}{\operatorname{arg\,min}} \left(x_t - r_j^{m-1} \right)^2$$

K-MEANS CONVERGENCE

- K-means algorithm converges to local minima of objective

$$O(c; \mathbf{r}_1, \dots, \mathbf{r}_K) = \sum_{j=1}^K \sum_{c(\mathbf{x}_t)=j} \|\mathbf{x}_t - \mathbf{r}_j\|_2^2$$

- Proof:

Clustering assignment improves objective:

$$O(\hat{c}^{m-1}; \mathbf{r}_1^{m-1}, \dots, \mathbf{r}_K^{m-1}) \geq O(\hat{c}^m; \mathbf{r}_1^{m-1}, \dots, \mathbf{r}_K^{m-1})$$

(By definition of $\hat{c}^m(\mathbf{x}_t)$)

Computing centroids improves objective:

$$O(\hat{c}^m; \mathbf{r}_1^{m-1}, \dots, \mathbf{r}_K^{m-1}) \geq O(\hat{c}^m; \mathbf{r}_1^m, \dots, \mathbf{r}_K^m)$$

(By the fact about centroid)

Fact: Centroid is Minimizer

$$\forall \mathbf{r}_j, \sum_{t \in C_j} \left\| \mathbf{x}_t - \overbrace{\frac{1}{|C_j|} \sum_{s \in C_j} \mathbf{x}_s}^{\text{Centroid}} \right\|^2 \leq \sum_{t \in C_j} \|\mathbf{x}_t - \mathbf{r}_j\|^2$$

$$0 = \frac{\partial \sum_{t \in C_j} (\mathbf{x}_t - \mathbf{r}_j)^2}{\partial \mathbf{r}_j} = \sum_{t \in C_j} 2(\mathbf{x}_t - \mathbf{r}_j) = -2|C_j| \mathbf{r}_j + 2 \sum_{t \in C_j} \mathbf{x}_t$$

$$\Leftrightarrow 2|C_j| \mathbf{r}_j = 2 \sum_{t \in C_j} \mathbf{x}_t$$

$$\mathbf{r}_j = \frac{1}{|C_j|} \sum_{t \in C_j} \mathbf{x}_t \leftarrow \text{minimum } \mathbf{r}_j$$

Time Complexity

- Assume computing distance between two instances is $O(d)$ where d is the dimensionality of the vectors.
- Reassigning clusters for n points: $O(Kn)$ distance computations, or $O(Knd)$.
- Computing centroids: Each instance gets added once to some centroid: $O(nd)$.
- Assume these two steps are each done once for i iterations: $O(iKnd)$.
- Linear in all relevant factors, assuming a fixed number of iterations.

Buckshot Algorithm

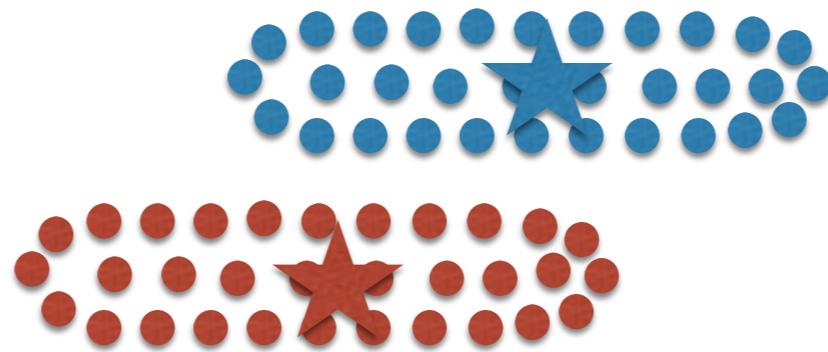
Problem

- Results can vary based on random seed selection, especially for high-dimensional data.
- Some seeds can result in poor convergence rate, or convergence to sub-optimal clusterings.

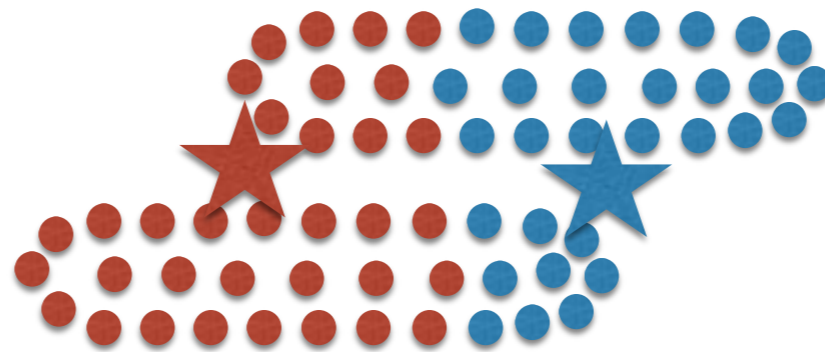
Idea: Combine HAC and K-means clustering.

- First randomly take a sample of instances of size $n^{1/2}$
- Run average-link HAC on this sample
- Use the results of HAC as initial seeds for K-means.
- Overall algorithm is efficient and avoids problems of bad seed selection.

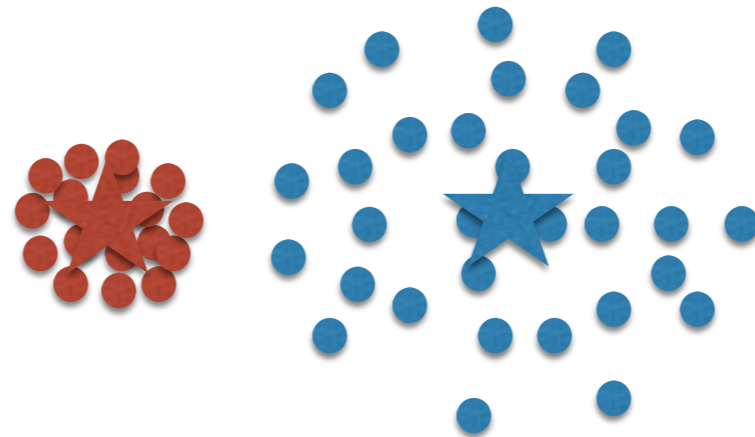
K-means: pitfalls



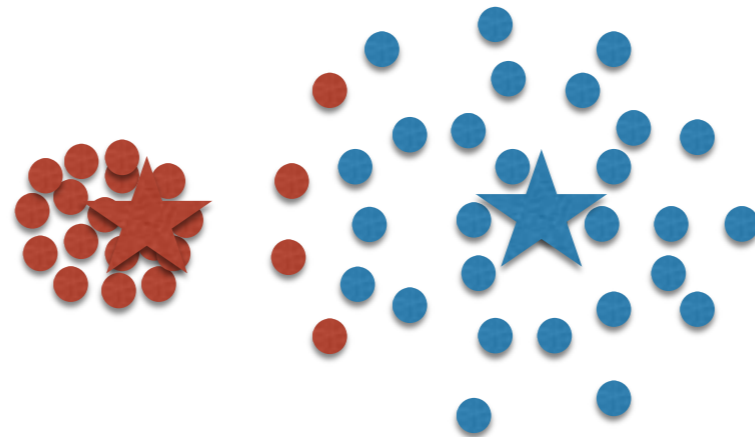
K-means: pitfalls



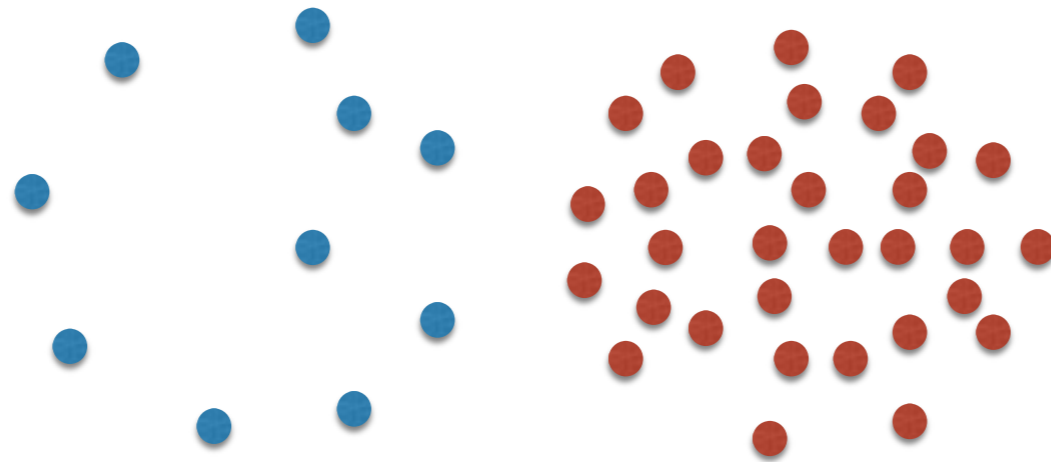
K-means: pitfalls



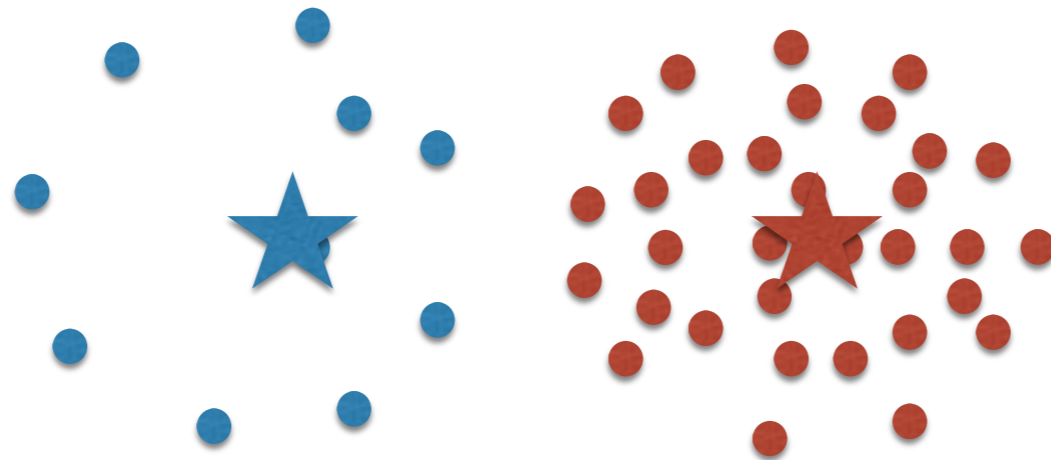
K-means: pitfalls



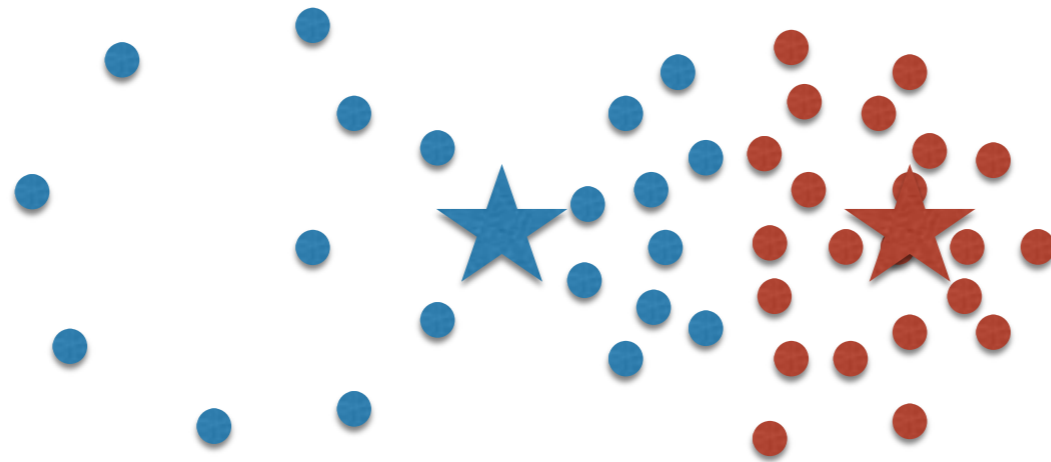
K-means: pitfalls



K-means: pitfalls



K-means: pitfalls



K-means: pitfalls

- Looks for spherical clusters
- Of same radius
- And with roughly equal number of points

K-means: pitfalls

- Can we design algorithm that can address these shortcomings?