

# Divisive clustering algorithms

# Divisive (top-down) methods: main idea

- A top-down approach is used to successively partition the data objects into a tree-like structure.
- A flat clustering algorithm  $\mathcal{A}$  (e.g. k-means algorithm) may be used for the partitioning in a given step. The algorithm  $\mathcal{A}$  can be any arbitrary clustering algorithm, not necessarily a distance based algorithm.

# Divisive (top-down) methods: main idea

- Provide flexibility in terms of choosing the trade-off between the balance in the tree structure and the balance in the number of objects in each cluster. For example
  - **Strategy 1:** split the heaviest node (a cluster with the maximum number of objects). Will result in leaf nodes (clusters) with a similar number of objects in them.
  - **Strategy 2:** split each cluster into the same number of subclusters. Will result in a balanced tree structure with the same number of children at each node, but the leaf nodes (most granular clusters) will have varying numbers of objects.

# Generic divisive clustering algorithm

**Input:** dataset:  $\mathcal{D}$ ; flat algorithm:  $\mathcal{A}$

1. **Initialise** tree  $\mathcal{T}$  to contain a single (root) vertex with entire dataset  $\mathcal{D}$
2. **Repeat:**
  1. Select a leaf node  $L$  in  $\mathcal{T}$  based on pre-defined criterion;
  2. Use algorithm  $\mathcal{A}$  split  $L$  into  $L_1, \dots, L_k$ ;
  3. Add  $L_1, \dots, L_k$  as children of  $L$  in  $\mathcal{T}$ ;
3. **Until** termination criterion
4. **Return** current clustering **or** hierarchy or clusterings

# Generic divisive clustering algorithm

**Input:** dataset:  $\mathcal{D}$ ; flat algorithm:  $\mathcal{A}$

1. **Initialise** tree  $\mathcal{T}$  to contain a single (root) vertex with entire dataset  $\mathcal{D}$

2. **Repeat:**

1. Select a leaf node  $L$  in  $\mathcal{T}$  based on pre-defined criterion;

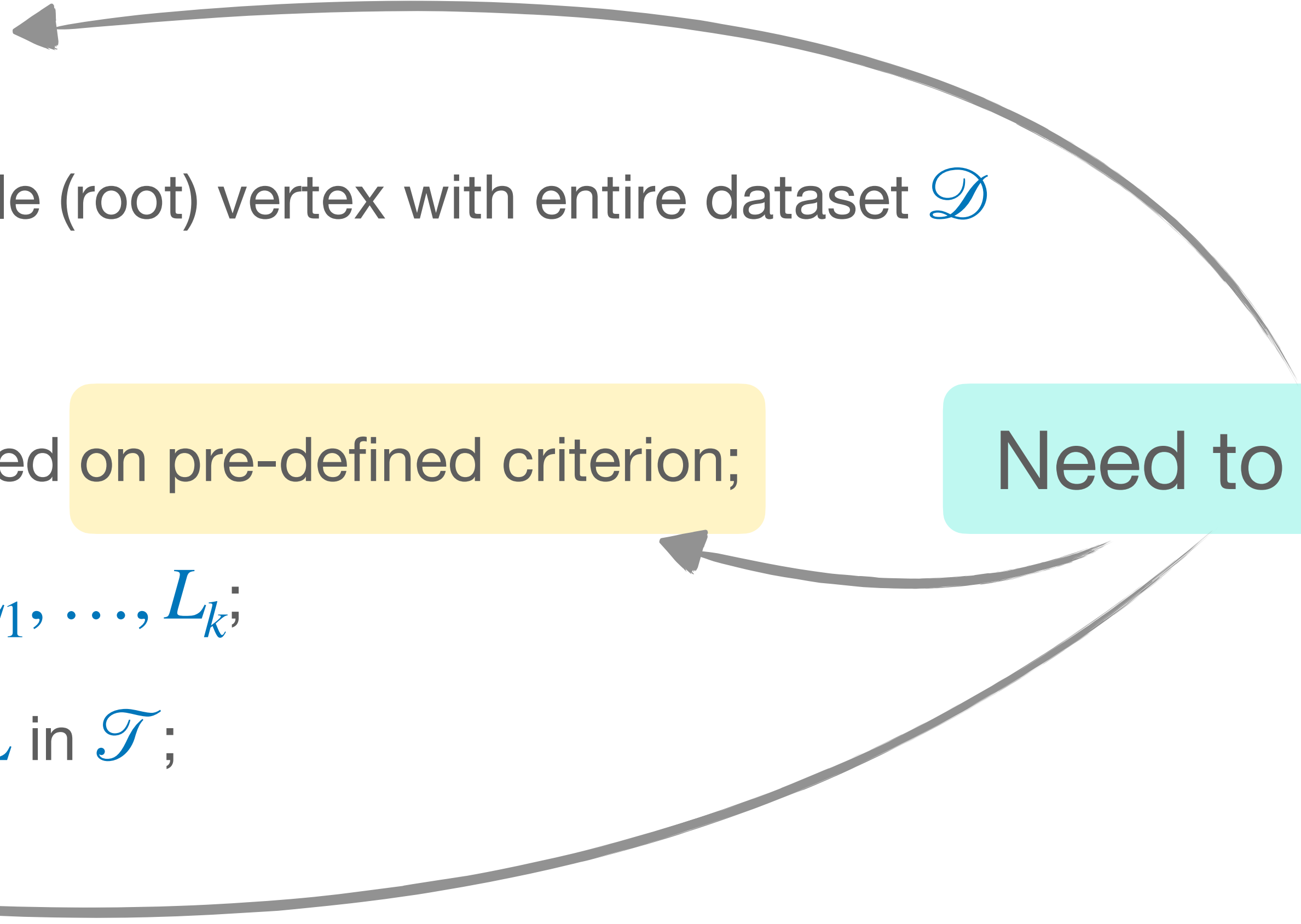
2. Use algorithm  $\mathcal{A}$  split  $L$  into  $L_1, \dots, L_k$ ;

3. Add  $L_1, \dots, L_k$  as children of  $L$  in  $\mathcal{T}$ ;

3. **Until** termination criterion

4. **Return** current clustering or hierarchy or clusterings

Need to specify

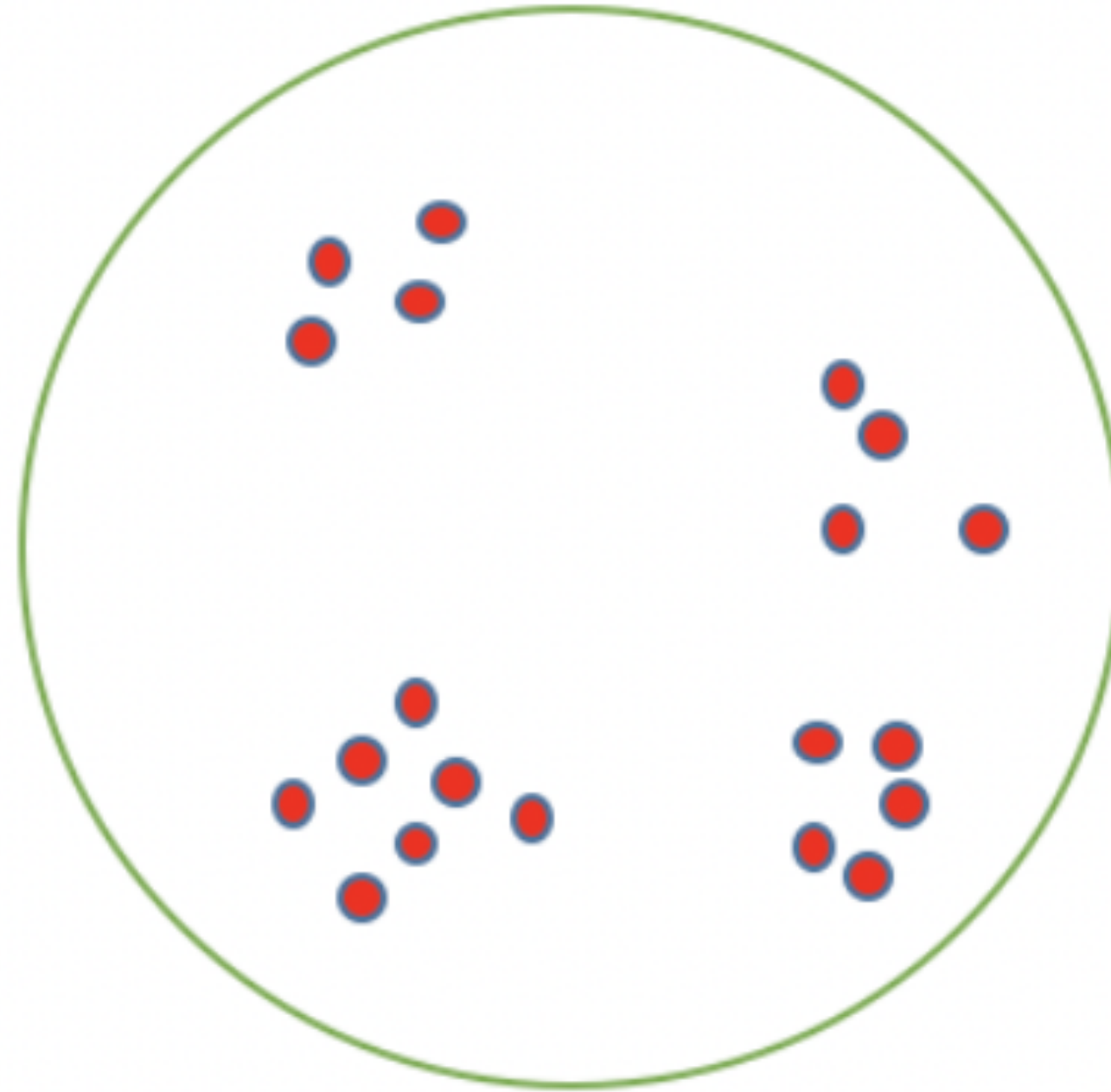


# Bisecting k-Means algorithm

**Input:** dataset:  $\mathcal{D}$ ; number of clusters:  $s$

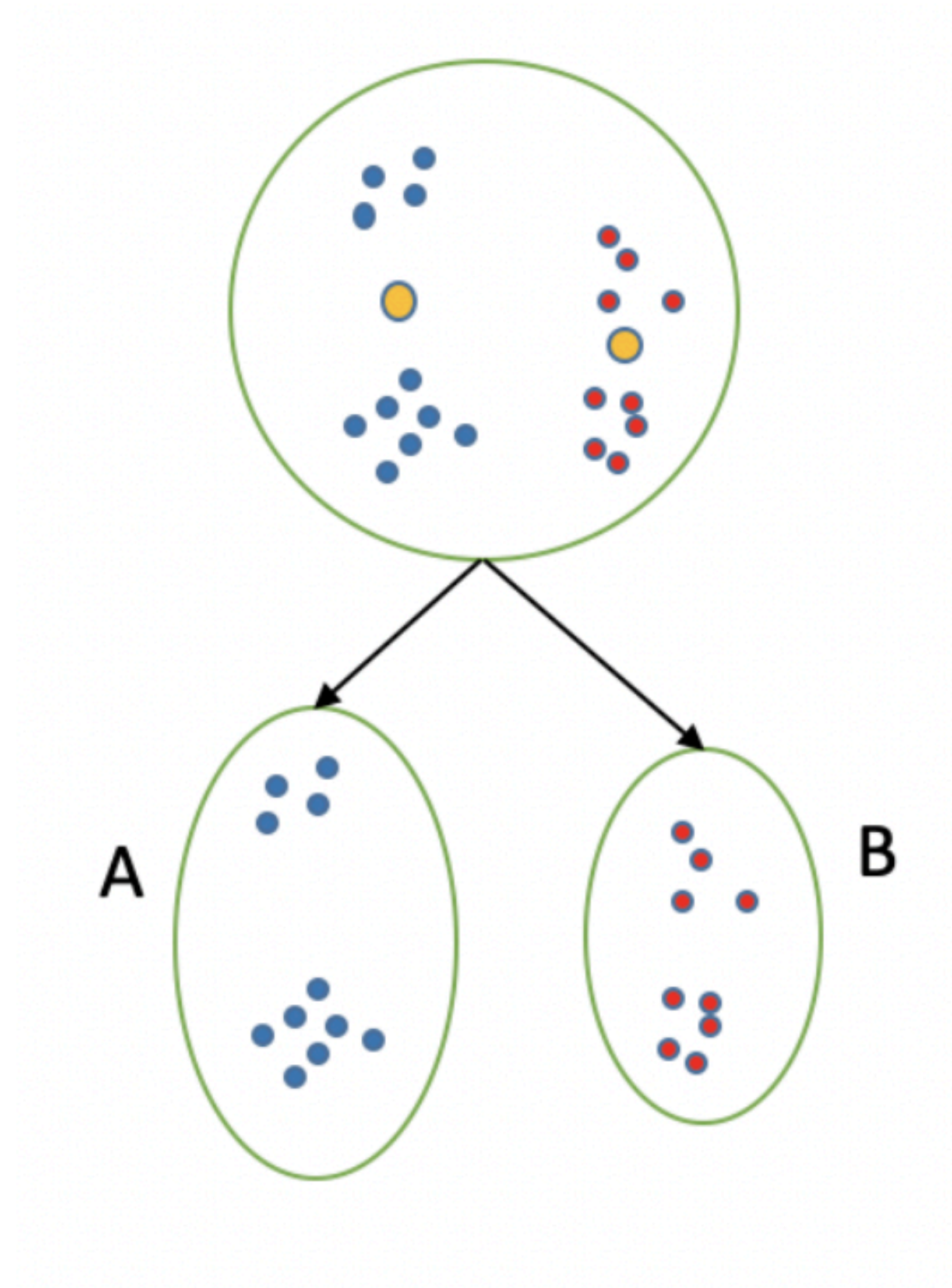
1. **Initialise** tree  $\mathcal{T}$  to contain a single (root) vertex with entire dataset  $\mathcal{D}$
2. **Repeat:**
  1. Select a leaf node (cluster)  $L$  in  $\mathcal{T}$  that has the largest sum of square distance  $\sum_{\bar{X}, \bar{Y} \in L} dist(\bar{X}, \bar{Y})^2$ ;
  2. Split  $L$  into 2 clusters  $L_1, L_2$  using k-means algorithm;
  3. Add  $L_1, L_2$  as children of  $L$  in  $\mathcal{T}$ ;
3. **Until** the number of leaf clusters is  $s$
4. **Return** the leaf clusters

# Bisecting k-Means algorithm: example



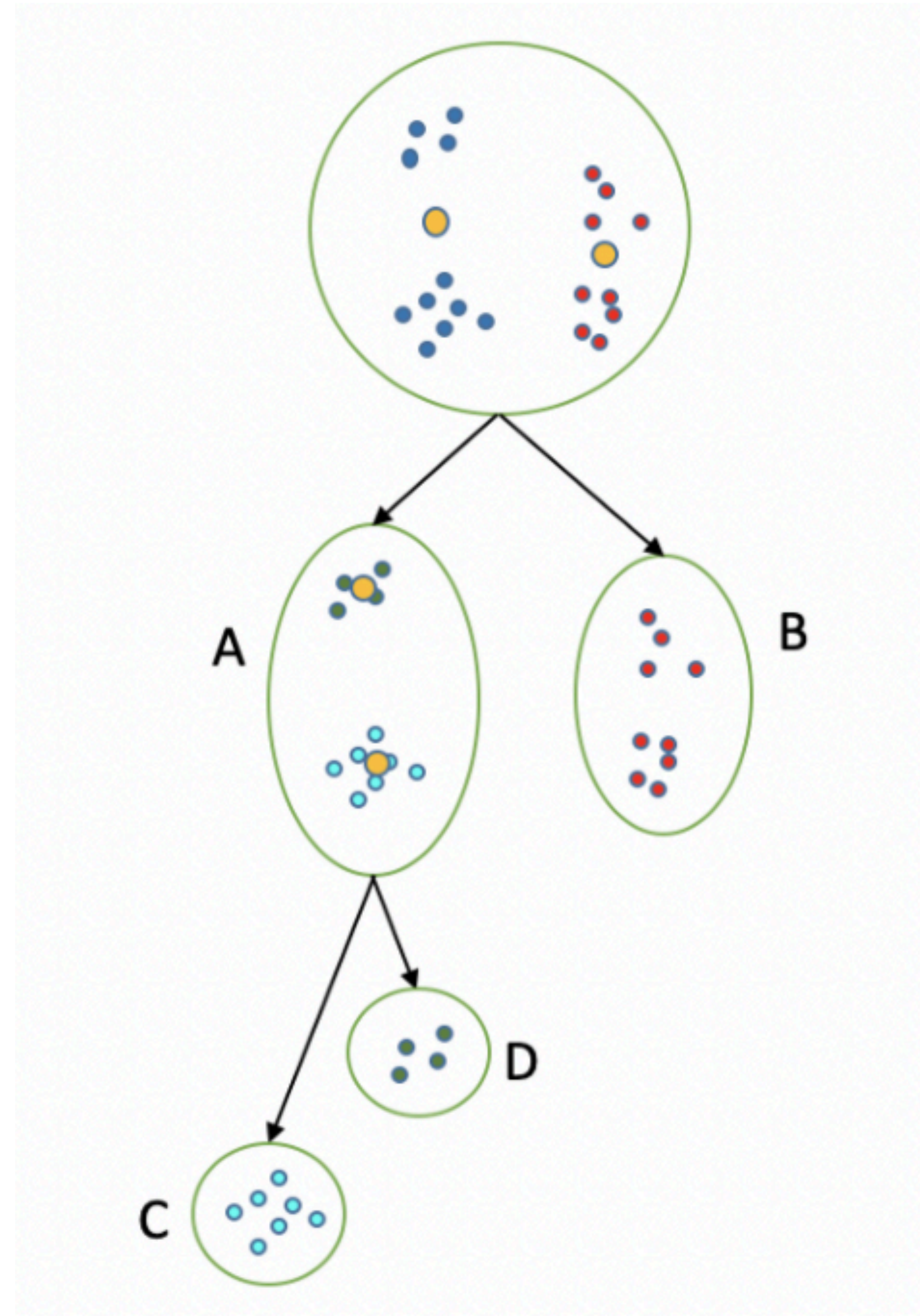


# Bisecting k-Means algorithm: example





# Bisecting k-Means algorithm: example





# Bisecting k-Means algorithm: example

