

Lab – Catch-up / Clustering #2

*Data Mining, Spring 2018*

*Carolina (berm@itu.dk) | Daniel (dafr@itu.dk) | Mathias (stsa@itu.dk)*

*Today's Lab: Catch-up and optional exercises*

# Catchup Lab

- Today's lab is mostly a catch up lab
  - Finish the previous labs
  - Or work on your individual assignment due on April 2nd.
  - Walk through of pseudo code
- Optional exercises:
  - Implement the GSP (Generalized Sequential Patterns) algorithm
  - Implement the FSG (Frequent Subgraph) algorithm



# *Preprocessing*

# Catch-up Lab

- Normalization
- Cleaning data
- Transforming data
- Data reduction

# Normalization

## Normalization

- **Min-max normalization:** to  $[new\_min_A, new\_max_A]$

$$v' = \frac{v - min_A}{max_A - min_A} (new\_max_A - new\_min_A) + new\_min_A$$

- Ex. Let income range \$12,000 to \$98,000 normalized to  $[0.0, 1.0]$ . Then \$73,000 is mapped to  $\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$

- **Z-score normalization** ( $\mu$ : mean,  $\sigma$ : standard deviation):

$$v' = \frac{v - \mu_A}{\sigma_A}$$

- Ex. Let  $\mu = 54,000$ ,  $\sigma = 16,000$ . Then  $\frac{73,600 - 54,000}{16,000} = 1.225$

- **Normalization by decimal scaling**

$$v' = \frac{v}{10^j} \quad \text{Where } j \text{ is the smallest integer such that } \text{Max}(|v'|) < 1$$

# Sample Code: Min-max

```
/**
 * Normalize x.
 * @param x The value to be normalized.
 * @return The result of the normalization.
 */
public double normalize(double x) {
    return ((x - dataLow)
            / (dataHigh - dataLow))
        * (normalizedHigh - normalizedLow) + normalizedLow;
}
```

```
public static double arrayMax(double[] arr) {
    double max = Double.NEGATIVE_INFINITY;

    for(double cur: arr)
        max = Math.max(max, cur);

    return max;
}
```



You can find this and other examples on [finding the max and min value in an array](#) in Java here.

# Decimal scaling

- **Normalization by decimal scaling**

$$v' = \frac{v}{10^j} \quad \text{Where } j \text{ is the smallest integer such that } \text{Max}(|v'|) < 1$$

## Decimal Scaling Normalization

Suppose that the recorded values of  $F$  range from  $-986$  to  $917$ . The maximum absolute value of  $F$  is  $986$ . To normalize by decimal scaling, we therefore divide each value by  $1,000$  (i.e.,  $j = 3$ ) so that  $-986$  normalizes to  $-0.986$  and  $917$  normalizes to  $0.917$ .



## *Classification*

# k-NN

- Step 1: Determine parameter  $K$  = number of nearest neighbors
- Step 2: Calculate the distance between the query-instance and all the training examples.
- Step 3: Sort the distance and determine nearest neighbors based on the  $k$ -th minimum distance.
- Step 4: Gather the category  $Y$  of the nearest neighbors.
- Step 5: Use simple majority of the category of nearest neighbors as the prediction value of the query instance.

In [this link](#) you can find examples of KNN using different distances.  
In [this article](#), you can find implementations and discussions of the different distances.

# ID3

## Input:


- Data partition,  $D$ , which is a set of training tuples and their associated class labels;
- *attribute\_list*, the set of candidate attributes;
- *Attribute\_selection\_method*, a procedure to determine the splitting criterion that “best” partitions the data tuples into individual classes. This criterion consists of a *splitting\_attribute* and, possibly, either a *split-point* or *splitting\_subset*.

**Output:** A decision tree.

## Method:

- (1) create a node  $N$ ;
- (2) **if** tuples in  $D$  are all of the same class,  $C$ , **then**
- (3)     return  $N$  as a leaf node labeled with the class  $C$ ;
- (4) **if** *attribute\_list* is empty **then**
- (5)     return  $N$  as a leaf node labeled with the majority class in  $D$ ; // majority voting
- (6) apply **Attribute\_selection\_method**( $D$ , *attribute\_list*) to **find** the “best” *splitting\_criterion*;
- (7) label node  $N$  with *splitting\_criterion*;
- (8) **if** *splitting\_attribute* is discrete-valued **and**  
      multiway splits allowed **then** // not restricted to binary trees
- (9)     *attribute\_list*  $\leftarrow$  *attribute\_list* – *splitting\_attribute*; // remove *splitting\_attribute*
- (10) **for each** outcome  $j$  of *splitting\_criterion*  
      // partition the tuples and grow subtrees for each partition
- (11)     let  $D_j$  be the set of data tuples in  $D$  satisfying outcome  $j$ ; // a partition
- (12)     **if**  $D_j$  is empty **then**
- (13)         attach a leaf labeled with the majority class in  $D$  to node  $N$ ;
- (14)     **else** attach the node returned by **Generate\_decision\_tree**( $D_j$ , *attribute\_list*) to node  $N$ ;
- endfor**
- (15) return  $N$ ;

Pseudocode  
explained  
with  
examples  
here and  
here



ID3 ( Learning Sets  $S$ , Attributes Sets  $A$ , Attributesvalues  $V$ )  
Return Decision Tree.

Begin

Load learning sets first, create decision tree root node 'rootNode', add learning set  $S$  into root node as its subset.

For rootNode, we compute  
Entropy(rootNode.subset) first

If Entropy(rootNode.subset)=0, then  
rootNode.subset consists of records  
all with the same value for the  
categorical attribute, return a leaf  
node with **decision**  
**attribute:attribute value**;

If Entropy(rootNode.subset)≠0, then  
compute information gain for each  
attribute left(have not been used in  
splitting). find attribute  $A$  with  
Maximum(Gain( $S$ , $A$ )). Create child  
nodes of this rootNode and add to  
rootNode in the decision tree.

For each child of the rootNode, apply  
ID3( $S$ , $A$ , $V$ ) recursively until reach  
node that has entropy=0 or reach  
leaf node.

End ID3.

# ID3: formulas

$$\text{Entropy}(S) = \sum - p(I) \cdot \log_2 p(I)$$

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum [ p(S|A) \cdot \text{Entropy}(S|A) ]$$

# ID3: example - decision making factors to play tennis at outside

Day	Outlook	Temp.	Humidity	Wind	Decision
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

## ID3: example - entropy

Find the [example here](#)

Initially, we need to calculate the overall entropy.

Decision column consists of 14 instances and includes two labels: yes and no. There are 9 decisions labeled yes, and 5 decisions labeled no.

$$\text{Entropy}(\text{Decision}) = - p(\text{Yes}) \cdot \log_2 p(\text{Yes}) - p(\text{No}) \cdot \log_2 p(\text{No})$$

$$\text{Entropy}(\text{Decision}) = - (9/14) \cdot \log_2(9/14) - (5/14) \cdot \log_2(5/14) = 0.940$$

Now, we need to find the most dominant factor for decisioning.

## ID3: example of formulas – wind factor on decision

$$\text{Gain}(\text{Decision}, \text{Wind}) = \text{Entropy}(\text{Decision}) - \sum [ p(\text{Decision}|\text{Wind}) \cdot \text{Entropy}(\text{Decision}|\text{Wind}) ]$$

Wind attribute has two labels: weak and strong. We would reflect it to the formula.

$$\begin{aligned} \text{Gain}(\text{Decision}, \text{Wind}) = & \text{Entropy}(\text{Decision}) - [ p(\text{Decision}|\text{Wind}=\text{Weak}) \cdot \text{Entropy}(\text{Decision}|\text{Wind}=\text{Weak}) \\ & ] - [ p(\text{Decision}|\text{Wind}=\text{Strong}) \cdot \text{Entropy}(\text{Decision}|\text{Wind}=\text{Strong}) ] \end{aligned}$$

We need to calculate (Decision|Wind=Weak) and (Decision|Wind=Strong) respectively.

There are 8 instances for weak wind. Decision of 2 items are no, whereas 6 items are yes as illustrated below.

$$\text{Entropy}(\text{Decision}|\text{Wind}=\text{Weak}) = - p(\text{No}) \cdot \log_2 p(\text{No}) - p(\text{Yes}) \cdot \log_2 p(\text{Yes})$$

$$\text{Entropy}(\text{Decision}|\text{Wind}=\text{Weak}) = - (2/8) \cdot \log_2(2/8) - (6/8) \cdot \log_2(6/8) = 0.811$$

## ID3: example of formulas – wind factor on decision

There are 6 instances for strong wind. Decision is divided into two equal parts.

$$\text{Entropy}(\text{Decision}|\text{Wind}=\text{Strong}) = - p(\text{No}) \cdot \log_2 p(\text{No}) - p(\text{Yes}) \cdot \log_2 p(\text{Yes})$$

$$\text{Entropy}(\text{Decision}|\text{Wind}=\text{Strong}) = - (3/6) \cdot \log_2(3/6) - (3/6) \cdot \log_2(3/6) = 1$$

Now, we can turn back to  $\text{Gain}(\text{Decision}, \text{Wind})$  equation.

$$\begin{aligned} \text{Gain}(\text{Decision}, \text{Wind}) = & \text{Entropy}(\text{Decision}) - [ p(\text{Decision}|\text{Wind}=\text{Weak}) \cdot \text{Entropy}(\text{Decision}|\text{Wind}=\text{Weak}) \\ & ] - [ p(\text{Decision}|\text{Wind}=\text{Strong}) \cdot \text{Entropy}(\text{Decision}|\text{Wind}=\text{Strong}) ] \end{aligned}$$

$$\text{Gain}(\text{Decision}, \text{Wind}) = 0.940 - [ (8/14) \cdot 0.811 ] - [ (6/14) \cdot 1 ] = 0.048$$

Calculations for wind column is over. Now, we need to apply same calculations for other columns to find the most dominant factor on decision.



# ID3: from decision tree to decision rules

A decision tree can easily be transformed to a set of rules by mapping from the root node to the leaf nodes one by one.

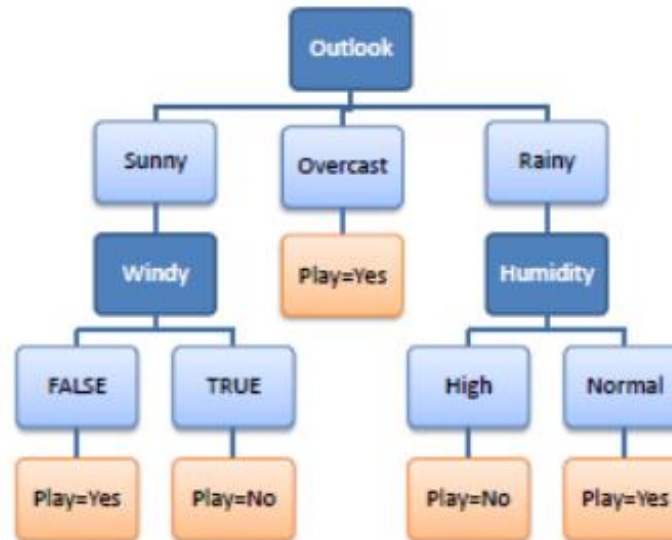
$R_1$ : IF (Outlook=Sunny) AND  
(Windy=FALSE) THEN Play=Yes

$R_2$ : IF (Outlook=Sunny) AND  
(Windy=TRUE) THEN Play=No

$R_3$ : IF (Outlook=Overcast) THEN  
Play=Yes

$R_4$ : IF (Outlook=Rainy) AND  
(Humidity=High) THEN Play=No

$R_5$ : IF (Outlook=Rain) AND  
(Humidity=Normal) THEN  
Play=Yes



# *Pattern Mining*

# Apriori

**Algorithm: Apriori.** Find frequent itemsets using an iterative level-wise approach based on candidate generation.

**Input:**

- $D$ , a database of transactions;
- $min\_sup$ , the minimum support count threshold.

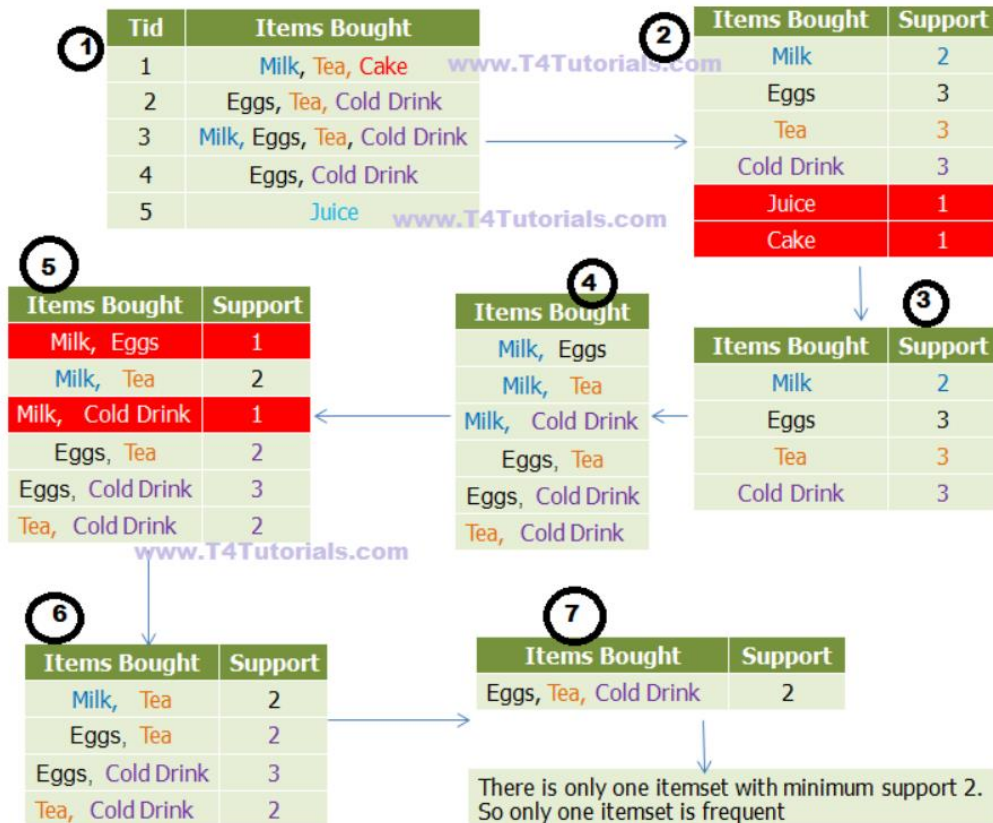
**Output:**  $L$ , frequent itemsets in  $D$ .

**Method:**

```
(1)  $L_1 = \text{find\_frequent\_1-itemsets}(D)$ ;  
(2) for ( $k = 2; L_{k-1} \neq \emptyset; k++$ ) {  
(3)    $C_k = \text{apriori\_gen}(L_{k-1})$ ;  
(4)   for each transaction  $t \in D$  { // scan  $D$  for counts  
(5)      $C_t = \text{subset}(C_k, t)$ ; // get the subsets of  $t$  that are candidates  
(6)     for each candidate  $c \in C_t$   
(7)        $c.\text{count}++$ ;  
(8)   }  
(9)    $L_k = \{c \in C_k | c.\text{count} \geq min\_sup\}$   
(10) }  
(11) return  $L = \cup_k L_k$ ;  
  
procedure apriori_gen( $L_{k-1}$ :frequent  $(k-1)$ -itemsets)  
(1) for each itemset  $l_1 \in L_{k-1}$   
(2)   for each itemset  $l_2 \in L_{k-1}$   
(3)     if ( $(l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2])$   
        $\wedge \dots \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$ ) then {  
(4)        $c = l_1 \bowtie l_2$ ; // join step: generate candidates  
(5)       if has_infrequent_subset( $c, L_{k-1}$ ) then  
(6)         delete  $c$ ; // prune step: remove unfruitful candidate  
(7)       else add  $c$  to  $C_k$ ;  
(8)     }  
(9) return  $C_k$ ;  
  
procedure has_infrequent_subset( $c$ : candidate  $k$ -itemset;  
   $L_{k-1}$ : frequent  $(k-1)$ -itemsets); // use prior knowledge  
(1) for each  $(k-1)$ -subset  $s$  of  $c$   
(2)   if  $s \notin L_{k-1}$  then  
(3)     return TRUE;  
(4) return FALSE;
```

# Apriori - example

Minimum Support: 2



Step 1: Data in the database

[quads id=1]

Step 2: Calculate the support/frequency of all items

Step 3: Discard the items with minimum support less than 2

Step 4: Combine two items

Step 5: Calculate the support/frequency of all items

Step 6: Discard the items with minimum support less than 2

Step 6.5: Combine three items and calculate their support.

Step 7: Discard the items with minimum support less than 2

**Result:**

Only one itemset is frequent (Eggs, Tea, Cold Drink) because this itemset has minimum support 2

# *Clustering*

# k-means

**Algorithm:  $k$ -means.** The  $k$ -means algorithm for partitioning, where each cluster's center is represented by the mean value of the objects in the cluster.

**Input:**

- $k$ : the number of clusters,
- $D$ : a data set containing  $n$  objects.

**Output:** A set of  $k$  clusters.

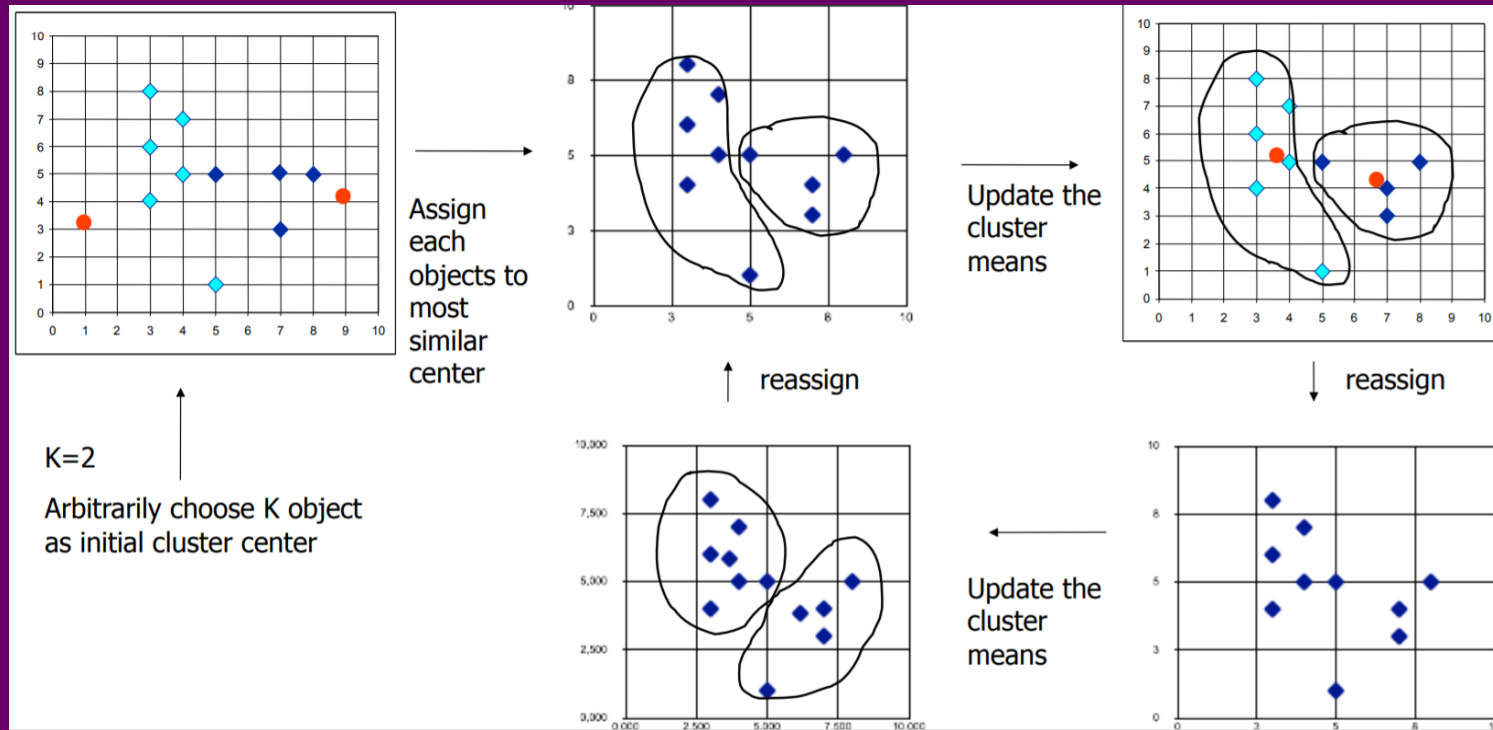
**Method:**

- (1) arbitrarily choose  $k$  objects from  $D$  as the initial cluster centers;
- (2) **repeat**
- (3)     (re)assign each object to the cluster to which the object is the most similar,  
            based on the mean value of the objects in the cluster;
- (4)     update the cluster means, that is, calculate the mean value of the objects for  
            each cluster;
- (5) **until** no change;

Remember to  
normalize each  
dimension first

Example of the implementation with 2 dimensions

# k-means



Example from the lecture

# k-medoids

**Algorithm:** *k-medoids*. PAM, a *k*-medoids algorithm for partitioning based on medoid or central objects.

**Input:**

- $k$ : the number of clusters,
- $D$ : a data set containing  $n$  objects.

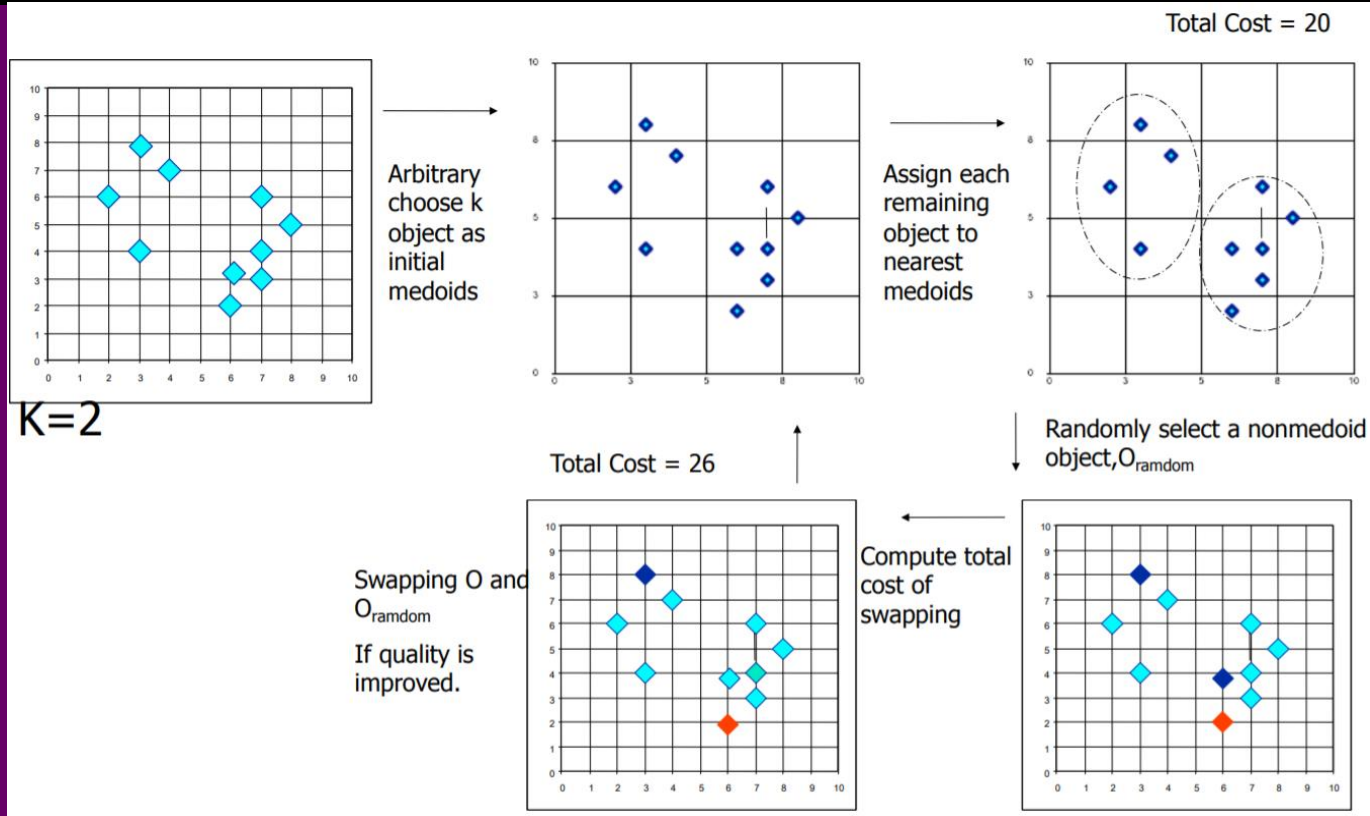
**Output:** A set of  $k$  clusters.

**Method:**

- (1) arbitrarily choose  $k$  objects in  $D$  as the initial representative objects or seeds;
- (2) **repeat**
- (3)     assign each remaining object to the cluster with the nearest representative object;
- (4)     randomly select a nonrepresentative object,  $o_{random}$ ;
- (5)     compute the total cost,  $S$ , of swapping representative object,  $o_j$ , with  $o_{random}$ ;
- (6)     **if**  $S < 0$  **then** swap  $o_j$  with  $o_{random}$  to form the new set of  $k$  representative objects;
- (7) **until** no change;



# k-medoids



[Example](#) from the lecture

*Thanks for listening!*