

Machine Learning

CSE 142

Xin (Eric) Wang

Friday, November 12, 2021

**T
o
d
a
y**

- Clustering (k-means, k-medoids)

Guest Lecture next Monday

- **Predicting behavior of road users**, by Vihan Jain
- **Bio:** Vihan Jain is a Staff Software Engineer at Waymo in Mountain View where he is tech-leading a team which develops and deploys Machine Learned models that predict the behavior of road users as the autonomous driving vehicle navigates through a scene. Previously, he was at Google Research where he worked on multi-modal learning, natural language compositionality and video semantic understanding. He has also worked on long-term value modeling in recommendation systems, designing configurable simulation platforms for studying recommendations, wide and deep learning and TensorFlow infrastructure. Prior to moving to the US, he worked with Ads Infrastructure at Google Canada and did an internship at Google India during the summer of 2012. He graduated from Indian Institute of Technology Roorkee as a gold medalist in 2013. In his free time, he likes to travel and play/follow several sports.

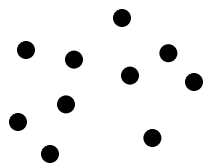
Clustering vs. classification

- Classification vs. clustering
 - In a **classifier**, possible class labels are provided
 - { dog, cat, elephant, mouse, ... }, { spam, ham }, etc.
 - Given in the training data (for supervised classification)
 - In a **clustering** problem, possible labels are the cluster labels learned from the training set
 - { cluster 1, cluster 2, cluster 3, ... }
 - Not given in the training data
- Terminology: In both cases, people often refer to the assigning of labels or clusters to data points (during the learning/training process, or afterwards in testing) as **classification**
 - Even if it's a **clustering** problem!

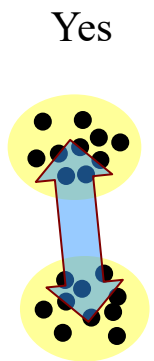
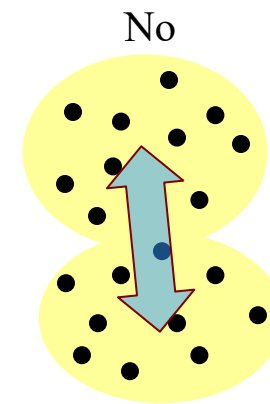
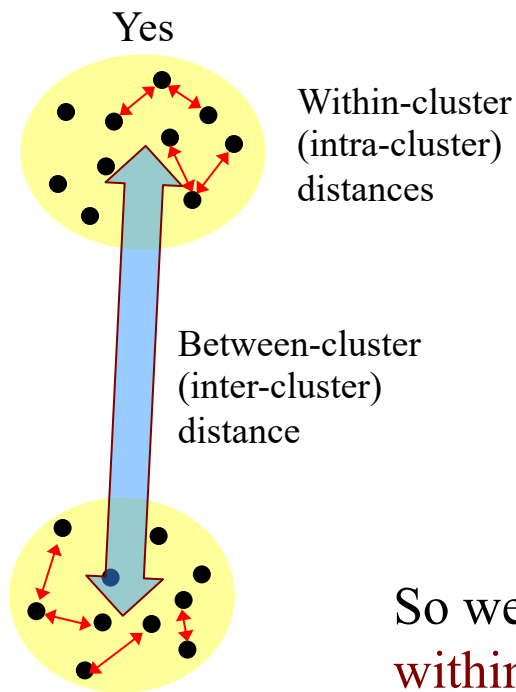
Clustering

- The goal of clustering is to find clusters (groupings) that are **compact** with respect to the distance metric
- What do we mean by *compactness*?

Is this cluster compact?



It depends....



So we'd like to have a measure of **within-cluster** and **between-cluster** distribution or *scatter*

Scatter matrix

This is a widely-used concept in ML!

Covariance matrix:

$$\left[\begin{array}{l} \text{Sample covariance: } \hat{\Sigma}_{ij} = \frac{1}{k} \sum_k (x_{ik} - \hat{\mu}_i)(x_{jk} - \hat{\mu}_j) = \frac{1}{k} S_{ij} \\ \text{If } \mathbf{X}_z \text{ is a matrix that holds all the zero-centered} \\ \text{samples as } \underline{\text{column}} \text{ vectors, then} \end{array} \right. \quad \hat{\Sigma} = \frac{1}{k} \boxed{\mathbf{X}_z \mathbf{X}_z^T} = \frac{1}{k} \mathbf{S}$$

S is the
Scatter matrix

Alternatively, if \mathbf{X}_z is a matrix that holds all the
zero-centered samples as row vectors, then

$$\mathbf{S} = \boxed{\mathbf{X}_z^T \mathbf{X}_z}$$

It depends on how we define \mathbf{X}_z !

For the **scatter matrix** (and thus the covariance matrix), \mathbf{X}_z is **zero-mean**

– That is, the mean data point $\bar{\mathbf{x}}$ (or $\boldsymbol{\mu}$ or $\boldsymbol{\mu}_x$) is first subtracted from every data point \mathbf{x}_i

By the way, the **Gram matrix** is not zero-mean...

Scatter matrix

- If the data D is partitioned into K subsets $\{D_1, D_2, \dots, D_K\}$ then the scatter matrix can be written as

$$\mathbf{S} = \left(\sum_{j=1}^K \mathbf{S}_j \right) + \mathbf{B}$$

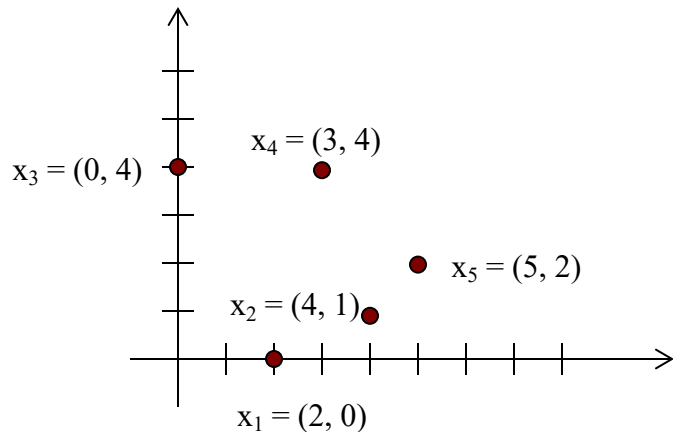
Subset scatter matrices
Within-cluster scatter matrices

Scatter matrix of the
partition means
Between-cluster scatter matrix

The diagram illustrates the decomposition of the total scatter matrix \mathbf{S} . The equation $\mathbf{S} = \left(\sum_{j=1}^K \mathbf{S}_j \right) + \mathbf{B}$ is shown. A blue arrow points from the text 'Subset scatter matrices' and 'Within-cluster scatter matrices' to the summation term $\sum_{j=1}^K \mathbf{S}_j$. Another blue arrow points from the text 'Scatter matrix of the partition means' and 'Between-cluster scatter matrix' to the term \mathbf{B} .

To compute \mathbf{B} , replace every point in D with the mean of its partition D_i and compute the scatter matrix

Example: Scatter matrix



$$\mathbf{X} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \mathbf{x}_3 \quad \mathbf{x}_4 \quad \mathbf{x}_5] = \begin{bmatrix} 2 & 4 & 0 & 3 & 5 \\ 0 & 1 & 4 & 4 & 2 \end{bmatrix}$$

The **Gram matrix** is ...

$$\mathbf{G} = \mathbf{X}^T \mathbf{X} = \begin{bmatrix} 4 & 8 & 0 & 6 & 10 \\ 8 & 17 & 4 & 16 & 22 \\ 0 & 4 & 16 & 16 & 8 \\ 6 & 16 & 16 & 25 & 23 \\ 10 & 22 & 8 & 23 & 29 \end{bmatrix}$$

($k \times k$), where k is the number of data points

$$\bar{\mathbf{x}} = \frac{1}{5} \sum_{i=1}^5 \mathbf{x}_i = \frac{1}{5} \begin{bmatrix} 14 \\ 11 \end{bmatrix} = \begin{bmatrix} 2.8 \\ 2.2 \end{bmatrix}$$

$$\mathbf{X}_z = [\mathbf{x}_1 - \bar{\mathbf{x}} \quad \mathbf{x}_2 - \bar{\mathbf{x}} \quad \mathbf{x}_3 - \bar{\mathbf{x}} \quad \mathbf{x}_4 - \bar{\mathbf{x}} \quad \mathbf{x}_5 - \bar{\mathbf{x}}] = \begin{bmatrix} -0.8 & 1.2 & -2.8 & 0.2 & 2.2 \\ -2.2 & -1.2 & 1.8 & 1.8 & -0.2 \end{bmatrix}$$

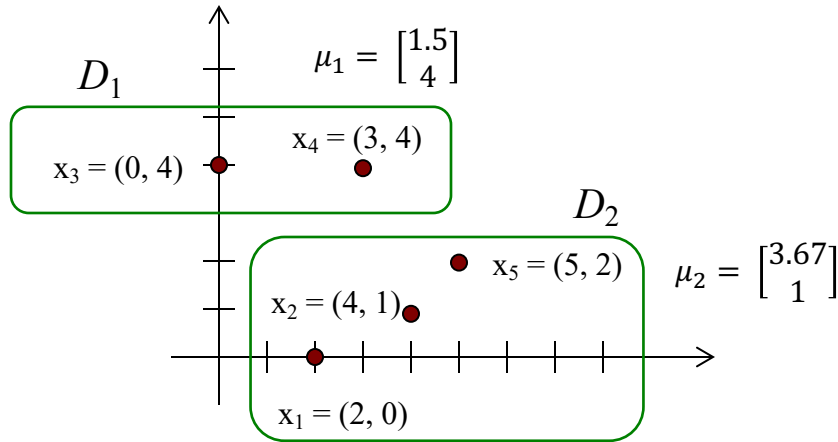
The **Scatter matrix** is ...

$$\mathbf{S} = \mathbf{X}_z \mathbf{X}_z^T = \begin{bmatrix} 14.8 & -4.8 \\ -4.8 & 12.8 \end{bmatrix}$$

($N \times N$), where N is the dimensionality of the data points

Example: Scatter matrix via partitions

$$S = \left(\sum_{j=1}^K S_j \right) + B$$



Partition means

$$\begin{bmatrix} 1.5 & 1.5 & 3.67 & 3.67 & 3.67 \\ 4 & 4 & 1 & 1 & 1 \end{bmatrix} \quad \mu_B = \begin{bmatrix} 2.8 \\ 2.2 \end{bmatrix}$$

Zero-mean partition means

$$B_Z = \begin{bmatrix} -1.3 & -1.3 & 13/15 & 13/15 & 13/15 \\ 1.8 & 1.8 & -1.2 & -1.2 & -1.2 \end{bmatrix}$$

Between-cluster scatter matrix

$$B = B_Z B_Z^T = \begin{bmatrix} 5.633 & -7.8 \\ -7.8 & 10.8 \end{bmatrix}$$

Scatter matrix of D_1

$$S_1 = \left[\mathbf{x}_3 - \begin{bmatrix} 1.5 \\ 4 \end{bmatrix} \quad \mathbf{x}_4 - \begin{bmatrix} 1.5 \\ 4 \end{bmatrix} \right] \left[\mathbf{x}_3 - \begin{bmatrix} 1.5 \\ 4 \end{bmatrix} \quad \mathbf{x}_4 - \begin{bmatrix} 1.5 \\ 4 \end{bmatrix} \right]^T = \begin{bmatrix} -1.5 & 1.5 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} -1.5 & 1.5 \\ 0 & 0 \end{bmatrix}^T = \begin{bmatrix} 4.5 & 0 \\ 0 & 0 \end{bmatrix}$$

Scatter matrix of D_2

$$S_2 = \begin{bmatrix} -5/3 & 1/3 & 4/3 \\ -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} -5/3 & 1/3 & 4/3 \\ -1 & 0 & 1 \end{bmatrix}^T = \begin{bmatrix} 4.67 & 3 \\ 3 & 2 \end{bmatrix}$$

$$S = S_1 + S_2 + B = \begin{bmatrix} 4.5 & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} 4.67 & 3 \\ 3 & 2 \end{bmatrix} + \begin{bmatrix} 5.633 & -7.8 \\ -7.8 & 10.8 \end{bmatrix} = \begin{bmatrix} 14.8 & -4.8 \\ -4.8 & 12.8 \end{bmatrix}$$

Clustering

- The goal of clustering is to find clusters (groupings) that are **compact** with respect to the distance metric
- Good clustering is characterized by low **within-cluster** variance and high **between-cluster** variance
- The **scatter matrix**, $\mathbf{S} = \mathbf{X}_z \mathbf{X}_z^T$, gives us a measure of variance
- If the **data** D is partitioned into K subsets/partitions $\{D_1, D_2, \dots, D_K\}$, then the scatter matrix can be written as

$$\mathbf{S} = \left(\sum_{j=1}^K \mathbf{S}_j \right) + \mathbf{B}$$

Within-cluster scatter matrices
(minimize!)

Between-cluster scatter matrix
(maximize!)

Scatter

- The **scatter** of D is defined as the **trace** of the scatter matrix
 - The *trace* is the sum of the diagonal elements of a square matrix

$$\begin{aligned}\text{Scat}(D) &= \text{Tr}(\mathbf{S}) = \text{Tr}(\mathbf{X}_z \mathbf{X}_z^T) \\ &= \text{Tr} \left(\begin{bmatrix} 14.8 & -4.8 \\ -4.8 & 12.8 \end{bmatrix} \right) = 14.8 + 12.8 = 27.6\end{aligned}$$

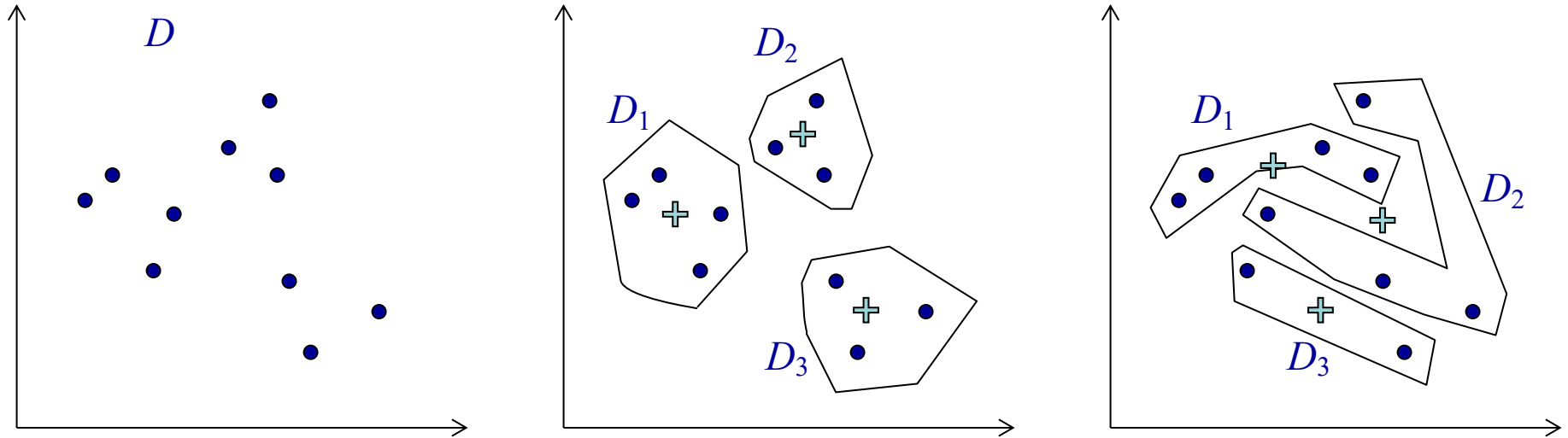
- Since \mathbf{S} can be decomposed into partitions, so can $\text{Scat}(D)$

$$\text{Scat}(D) = \underbrace{\sum_{j=1}^K \text{Scat}(D_j)}_{\substack{\text{Want to choose} \\ \text{partitions that} \\ \text{minimize this}}} + \underbrace{\sum_{j=1}^K |D_j| \|\mu_j - \mu\|^2}_{\substack{\text{Equivalent to} \\ \text{maximizing this}}}$$

Fixed for a given data set

This is the goal of
k-means clustering

Scatter



$$\text{Scat}(D) = \sum_{j=1}^K \text{Scat}(D_j) + \sum_{j=1}^K |D_j| \|\mu_j - \mu\|^2$$

K-means clustering

- The general **K-means clustering** problem is **NP-complete**, so there is no efficient solution to find the **optimal clustering** (data partition)
- A widely-used **heuristic** algorithm for clustering is also known as the **K-means algorithm**, but it is not optimal
 - It will converge to a solution, but there is no guarantee that the solution is the best one (the global minimum of scatter)
 - But it works quite well in most cases!
- Typically, the **K-means algorithm** would be run several times (with a random starting point) and then the best solution is selected
 - I.e., the solution with the smallest **within-cluster scatter**

K-means algorithm

K is an input parameter

Algorithm $KMeans(D, K)$ – K -means clustering using Euclidean distance Dis_2 .

Input : data $D \subseteq \mathbb{R}^d$; number of clusters $K \in \mathbb{N}$.

Output : K cluster means $\mu_1, \dots, \mu_K \in \mathbb{R}^d$.

randomly initialise K vectors $\mu_1, \dots, \mu_K \in \mathbb{R}^d$;

repeat

 assign each $\mathbf{x} \in D$ to $\arg\min_j Dis_2(\mathbf{x}, \mu_j)$; \leftarrow *1-Nearest neighbor assignment*

for $j = 1$ to K **do**

$D_j \leftarrow \{\mathbf{x} \in D \mid \mathbf{x} \text{ assigned to cluster } j\}$; \leftarrow *Partition defined by assignment*

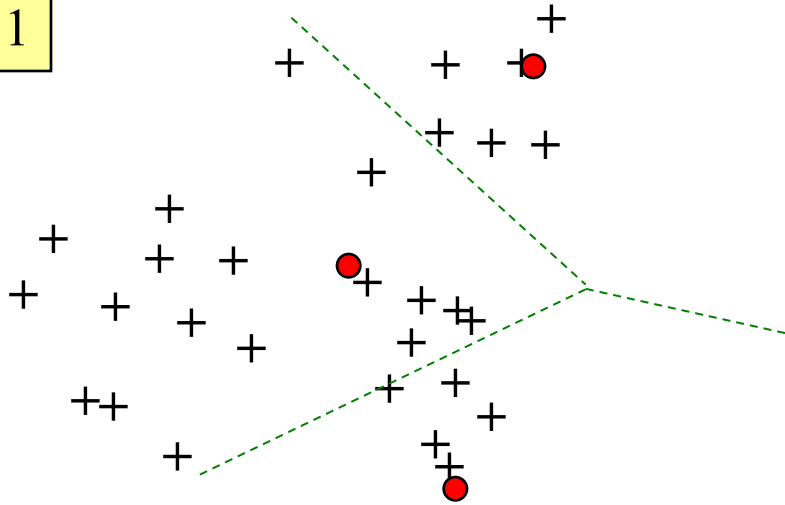
$\mu_j = \frac{1}{|D_j|} \sum_{\mathbf{x} \in D_j} \mathbf{x}$; \leftarrow *Re-compute the cluster mean*

end

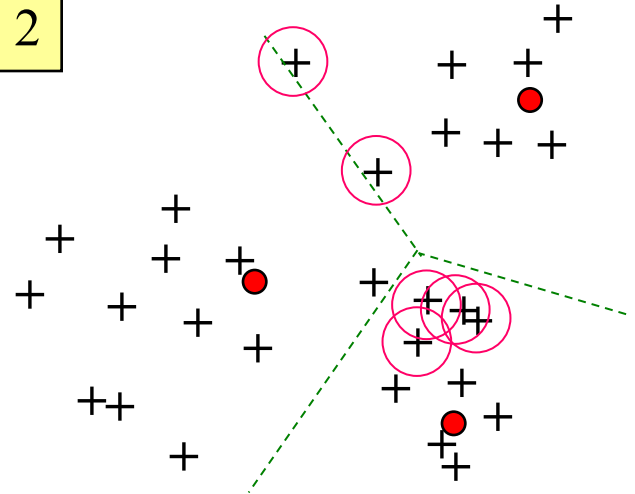
until no change in μ_1, \dots, μ_K ;

return μ_1, \dots, μ_K ;

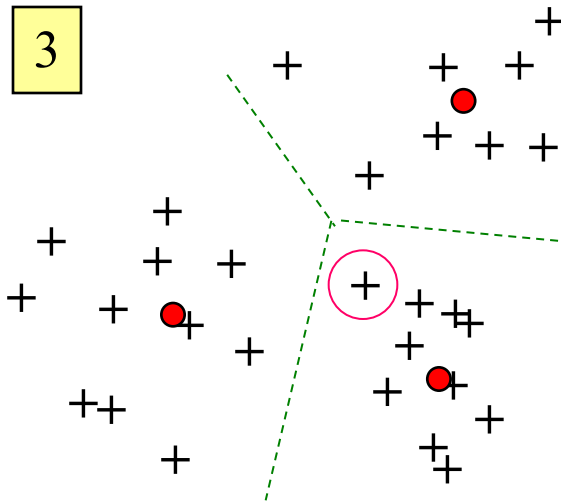
1



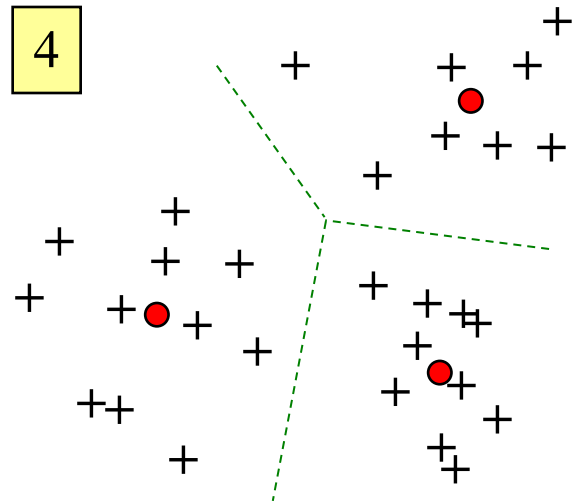
2



3



4



No change – finished!

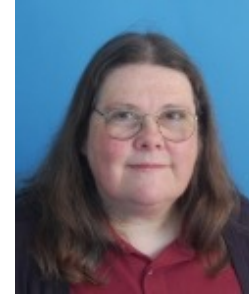
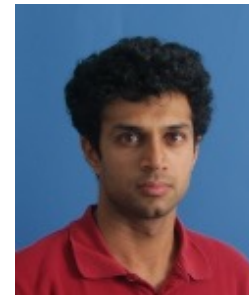
K-means demo

Demo: <https://www.naftaliharris.com/blog/visualizing-k-means-clustering/>

K-means applications

In K-means, we simply take a cluster centroid (exemplar) to be the **mean** of points in the cluster.

- Document clustering
- House clustering based on price, square footage, #bedrooms, etc.
- Image segmentation (pixel clustering based on RGB values and location)
-

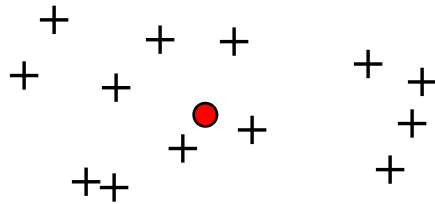


K-medoids algorithm

- In some problems, the cluster **exemplars** are required to be **data points** (from the training data)
 - As opposed to using the **mean** of the cluster points, for example, since the mean is most likely not a point in the data set
- The concept of **medoid** is useful here – the **medoid** of a set of points is the point with the **minimal average dissimilarity** (distance) to all other points in the set
 - Using some distance metric: Euclidian, L1, etc.
 - This is a generalization of the concept of **median** to multiple dimensions
- K-means can be modified to use **data points** as exemplars rather than **means**, by instead computing the cluster **medoids**

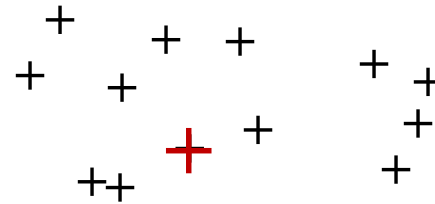
K-medoids algorithm

Cluster mean



Location that minimizes the sum of squared distances to points

Cluster medoid



Point that minimizes the sum of squared distances to points

K-medoids algorithm

Algorithm $KMedoids(D, K, Dis)$ – K -medoids clustering using arbitrary distance metric Dis .

Input : data $D \subseteq \mathcal{X}$; number of clusters $K \in \mathbb{N}$;
distance metric $Dis : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$.

Output : K medoids $\mu_1, \dots, \mu_K \in D$, representing a predictive clustering of \mathcal{X} .
randomly pick K data points $\mu_1, \dots, \mu_K \in D$;

repeat

 assign each $\mathbf{x} \in D$ to $\arg \min_j Dis(\mathbf{x}, \mu_j)$;

for $j = 1$ to K **do**

$D_j \leftarrow \{\mathbf{x} \in D \mid \mathbf{x} \text{ assigned to cluster } j\}$;

$\mu_j = \arg \min_{\mathbf{x} \in D_j} \sum_{\mathbf{x}' \in D_j} Dis(\mathbf{x}, \mathbf{x}')$; *← Re-compute the cluster medoid*

end

until no change in μ_1, \dots, μ_K ;

return μ_1, \dots, μ_K ;

Summary: Distance methods and clustering

- **Similarity** is a function of **distance**
- Euclidian distance may not always be the right choice of distance metric
- **Nearest neighbor** methods assign classes/clusters based on distances to points or **exemplars**, not based on computed **boundaries**
- For good clustering, we want high **within-class** (intra-class) similarity and low **between-class** (inter-class) similarity
- The **scatter matrix** is an important structure in clustering
- The K-means algorithm (and variations) is widely used