

Until now:

Supervised learning:

$$\vec{x}^{(i)} = (x_1^{(i)}, \dots, x_p^{(i)}) ; y^{(i)} \quad i = 1, \dots, N$$

Given \vec{x}^{in} $\overset{\text{infer}}{\underset{\uparrow}{f}}(\vec{x}^{in}) = \hat{y}$

inferred from $(\vec{x}^{(i)}, y^{(i)}) \quad i = 1, \dots, N$

Switch now
to

Unsupervised learning:

- There is no 'observable', y
- There is no ground truth.
- There are no examples.
- There is no training.

We have a series of samples:

$$\{\vec{x}^{(i)}\}_{i=1}^N$$

Data

$$\vec{x}^{(i)} = (x_1^{(i)}, \dots, x_p^{(i)}) \quad i=1, \dots, n$$

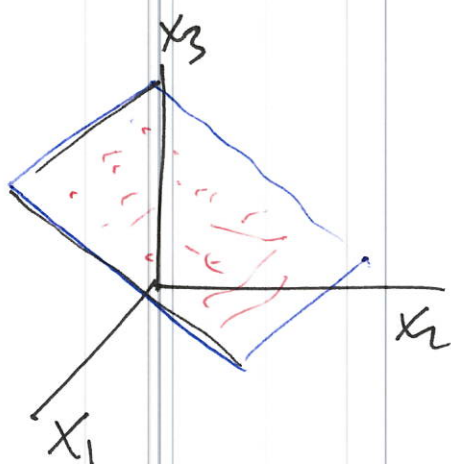
of different types

$$\begin{cases} \vec{x}^{(i)} \in \mathbb{R}^p \\ \vec{x}^{(i)} \in \{c_1, \dots, c_k\}^p \end{cases}$$

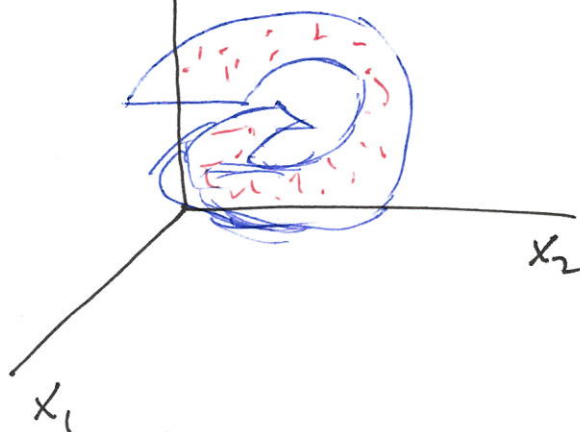
Types of tasks in unsupervised learning.

(1) Clustering: find groups of points that are more similar to each other within the group than to points outside the group.

(2) Dimensionality reduction:



$p \gg$ $\vec{x}^{(i)} \in \mathbb{R}^3$



Visualization .

Clustering :

Key ingredient : Similarity vs Dissimilarity
"Distance".
(not always a true metric)

$$D(\vec{x}^{(i)}, \vec{x}^{(j)})$$

Typical: If $\vec{x} \in \mathbb{R}^p$

$$D(\vec{x}^{(i)}, \vec{x}^{(j)}) = \|\vec{x}^{(i)} - \vec{x}^{(j)}\|^2$$

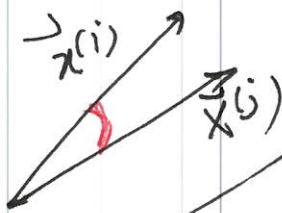
(i) dissimilarity:

$$= |\vec{x}^{(i)} - \vec{x}^{(j)}|$$

Distances.

(ii) similarities

cosine similarity: $S(\vec{x}^{(i)}, \vec{x}^{(j)}) = \frac{\vec{x}^{(i)} \cdot \vec{x}^{(j)}}{\|\vec{x}^{(i)}\| \|\vec{x}^{(j)}\|}$



statistical similarity:

$$g(\vec{x}^{(i)}, \vec{x}^{(j)}) = \frac{\text{cov}(\vec{x}^{(i)}, \vec{x}^{(j)})}{\sqrt{\text{cov}(\vec{x}^{(i)}, \vec{x}^{(i)}) \cdot \text{cov}(\vec{x}^{(j)}, \vec{x}^{(j)})}}$$

B. If \vec{x} are categorical:

x belonging to K classes

Distance matrix is based on distance assigned to feature classes.

	C1	C2	C3
C1	0	D_{12}	D_{13}
C2		0	D_{23}
C3			0

$K=3$

C. Ordinal variables: ranked.

$A \rightarrow B \rightarrow C \rightarrow D$

$K=4$

$i = 1 \quad 2 \quad 3 \quad 4$

$$j = \frac{i - \frac{1}{2}}{K} \left[\frac{1}{8} \quad \frac{3}{8} \quad \frac{5}{8} \quad \frac{7}{8} \right]$$

~~Heuristics for clustering:~~

~~1. K-means~~

~~$$\left\{ \vec{x}^{(i)} \mid \begin{matrix} N \\ i=1 \end{matrix} \right\} \vec{x}^{(i)} \in \mathbb{R}^p$$~~

~~$$D(\vec{x}^{(i)}, \vec{x}^{(j)}) = \|\vec{x}^{(i)} - \vec{x}^{(j)}\|^2$$~~

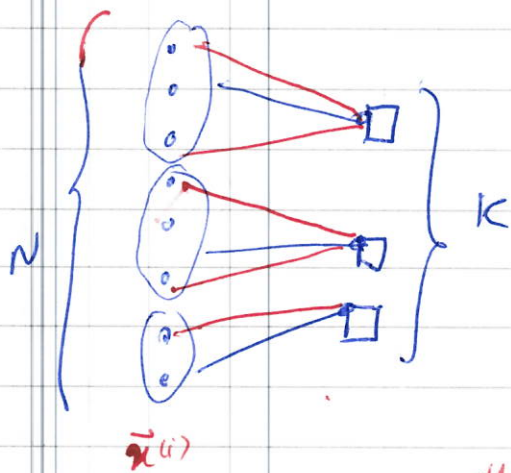
Unsupervised learning:

(1) Clustering problem

General statement: Given $\{\vec{x}^{(i)}\}_{i=1}^N$ our N samples
and a dissimilarity ("distance")
 $D(\vec{x}^{(i)}, \vec{x}^{(j)})$,

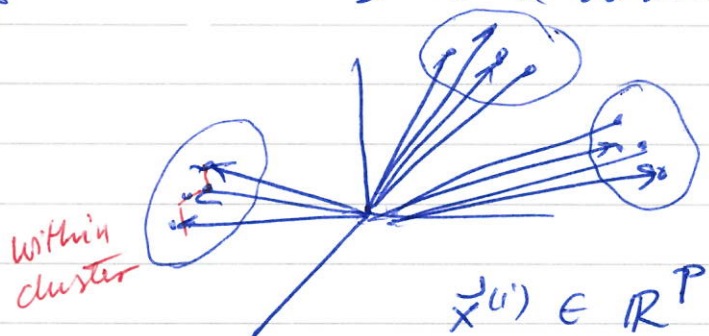
the objective is to find a partition of
the N samples into K ~~classes~~ clusters
such that the dissimilarity is smaller
within cluster ~~and~~ than across clusters.

In pictures: we have a mapping
of the N samples to K clusters



Geometrically:

If the samples are
continuous variables and
 D is a distance:

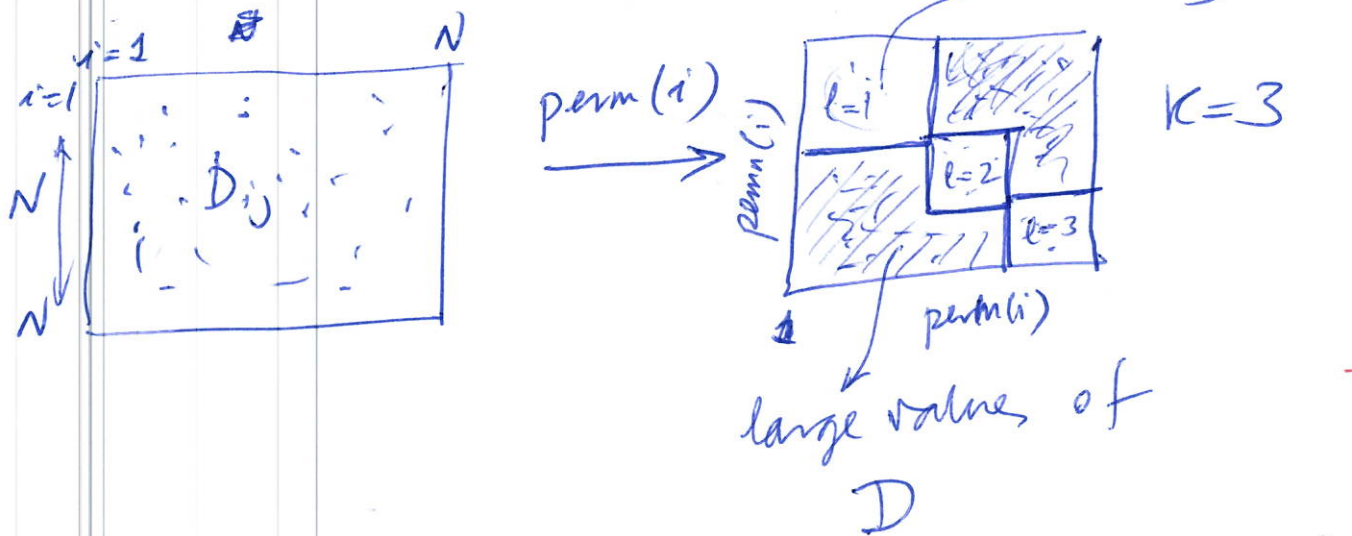


such that

$$\left\{ \begin{array}{l} \sum_{\text{within cluster}} D_{ij} \ll \\ \sum_{\text{across clusters}} D_{ij} \gg \end{array} \right\}$$

Another look at this problem:

We have a $D_{N \times N}$ matrix and we want to reorder the rows and columns (i.e. find the permutation) that makes blocks in your D matrix:



This will mean that for this clustering:

Within cluster dissimilarity

$$W(C) = \sum_{l=1}^K \sum_{i,j \in C_l} d(\vec{x}^{(i)}, \vec{x}^{(j)}) \frac{1}{2}$$

is small

and

Total dissimilarity

$$B(C) = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d(\vec{x}^{(i)}, \vec{x}^{(j)}) - W(C)$$

is large

Between cluster dissimilarity

Finding the best H (or the best clustering) is

a Combinatorial optimisation problem:

- There are usually NP-hard, all possible arrangements have to be tried to find the global optimum.
- But the number of clusterings to try ~~is~~ explodes:
For N samples, K clusters.

$$S(N, K) = \frac{1}{K!} \sum_{l=1}^K (-1)^{K-l} \binom{K}{l} l^N$$

$$S(10, 4) = 34,105$$

$$S(19, 4) = 10^{10}$$

⋮

Impractical to do this enumeration.

Therefore, we need to come up with heuristics to optimise.

k-means : What is the k-means heuristic?

We want to find an assignment between the N samples and k ~~clusters~~ clusters

First: Decide on k

Task: Find an assignment matrix:

$$H_{N \times k} = \begin{matrix} & \begin{matrix} c_1 & & & & c_k \end{matrix} \\ \begin{matrix} i=1 \\ \vdots \\ i=N \end{matrix} & \begin{bmatrix} 0 & 0 & 1 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 0 & \dots & 0 \\ \vdots & & & & \\ 0 & \dots & \dots & 0 & 1 \end{bmatrix} \end{matrix}$$

This summarizes a hard assignment that is exhaustive.

$$\{C_\ell\}_{\ell=1}^k \quad (i) \quad C_\ell \cap C_{\ell'} = \emptyset \quad \ell \neq \ell'$$

$$(ii) \quad \bigcup_{\ell=1}^k C_\ell = \{\vec{x}^{(i)}\}_{i=1}^N$$

As stated above,

This is a combinatorial problem : check all the permutations to see which one is best

k-means algorithm :

Given $\{\vec{x}^{(i)}\} \in \mathbb{R}^p$

Compute $D(\vec{x}^{(i)}, \vec{x}^{(j)}) = \|\vec{x}^{(i)} - \vec{x}^{(j)}\|^2$ Distance matrix

For a given clustering: $\{C_\ell\}_{\ell=1}^K = \mathcal{C}$

~~can't minimize~~
we have the within cluster normalized distance

$$W(\mathcal{C}) = \frac{1}{2} \sum_{\ell=1}^K \frac{1}{|C_\ell|} \sum_{i,j \in C_\ell} \|\vec{x}^{(i)} - \vec{x}^{(j)}\|^2$$

~~equivalent to finding $H_{N \times K}$~~

where $|C_\ell|$ is the cardinality of C_ℓ

$W(\mathcal{C})$ can be rewritten in terms of H :

$$\text{Tr} \left[\frac{1}{2} (H^T H)^{-1} \left[\begin{matrix} H^T & D_{N \times N} & H \\ k \times N & \text{"} & N \times K \end{matrix} \right] \right] = W(\mathcal{C})$$

$\text{diag}(C_\ell) = H^T H =$
 $\begin{bmatrix} \square & 0 & 0 \\ 0 & \square & 0 \\ 0 & 0 & \text{red } \square \end{bmatrix}$
 $|C_3|$

 $\begin{bmatrix} \square & \text{blue } \square & \text{blue } \square \\ \text{blue } \square & \square & \text{blue } \square \\ \text{blue } \square & \text{blue } \square & \square \end{bmatrix}$
 $k \times k$

 $k=3$
 $\sum_{i,j \in C_3} \|\vec{x}^{(i)} - \vec{x}^{(j)}\|^2$

$$T = \frac{1}{2} \underbrace{\vec{1}^T}_{1 \times N} \underbrace{D}_{N \times N} \underbrace{\vec{1}}_{N \times 1} = \frac{1}{2} \underbrace{\vec{1}^T}_{1 \times K} \underbrace{[H^T D H]}_{K \times K} \underbrace{\vec{1}}_{K \times 1}$$

$$\vec{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}_{N \times 1}$$

$$\underbrace{H}_{N \times K} \underbrace{\vec{1}}_{K \times 1} = \underbrace{\vec{1}}_{N \times 1}$$

Objective: Minimise $W(C)$

\Downarrow

maximise $T - W(C) = \text{between cluster distance.}$

$\min_H \text{Tr} \left[\frac{1}{2} (H^T H)^{-1} (H^T D H) \right]$

$D_{N \times N}(\vec{x}^{(i)})$ is the distance matrix

Algorithm:

• Step 0: Assign at random ~~the~~ every sample to a cluster.

Initialization

• Step 1: compute the centroid of each of the k clusters

$$\vec{m}_\ell = \frac{1}{|C_\ell|} \sum_{i \in C_\ell} \vec{x}^{(i)}$$

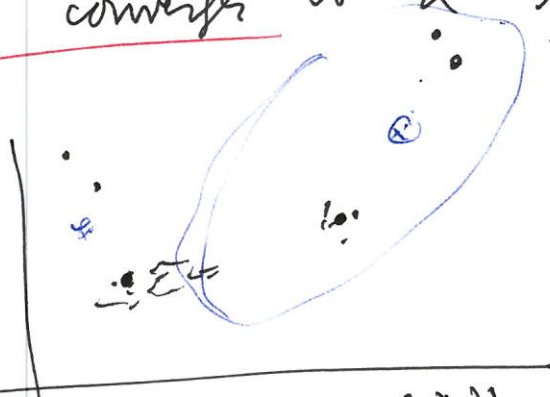
• Step 2: Reassign each $\vec{x}^{(i)}$ to the closest centroid.

$\forall i$, ~~Find~~ ^{Consider} $l_i^{(t)}$
 Evaluate $l_i^{(t+1)} = \arg \min_l \| \bar{x}^{(i)} - \bar{m}_l \|^2$
 Iterate step 1 and 2
 Until convergence:
 i.e., $W(C)$ does not improve much.
 or assignments are not changing.

Heuristic: Gradient method
 because at every step $W(C)$
 is decreased.

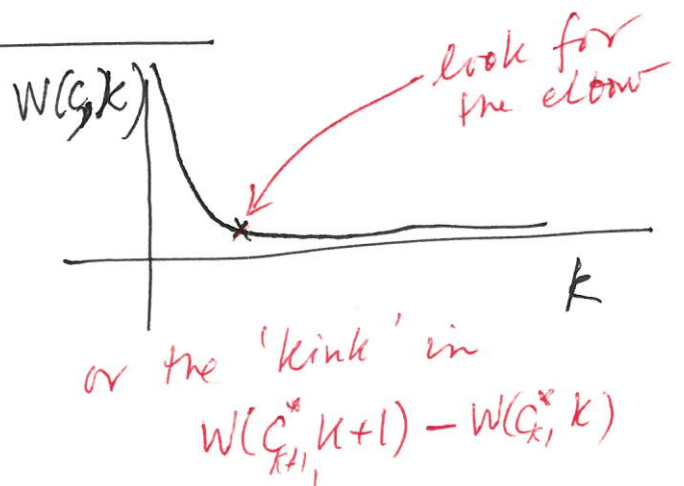
It converges to a local optimum.

(1) Outliers:



Instead:
 k-medoids.

(2) Try different k 's:
 How to choose k



Comments: