

Tópicos Especiais em Mecânica Estatística

-Mathematical and Computational Methods of Machine Learning-

Lisan Durão, P.h.D

(GetNinjas S/A)

Aula 6 – Unsupervised Learning I: Clustering

0. Previous Lectures

Supervised Learning:

=

Learning with

Teacher

1. Data : $\mathcal{D} = \{(\vec{x}_i, y_i)\}_{i=1}^N$

2. Model : $f : \mathbb{R}^D \rightarrow \mathbb{R}$

3. Learning: minimize $C(\vec{x}, y, \theta)$

Feature: $\vec{x} \in \mathbb{R}^D \rightarrow$ Euclidean Vector Space.

Target: $y \in \mathbb{R}$

$(\vec{x}, y) \sim$ Random Variables $\Rightarrow P(\vec{x}, y)$ joint
density probability

Goal $\mathcal{M} \rightarrow P(y | \vec{x})$

$$\mu(\vec{x}) = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{y | \vec{x}} \left[\hat{h}(y, \theta) \right]$$

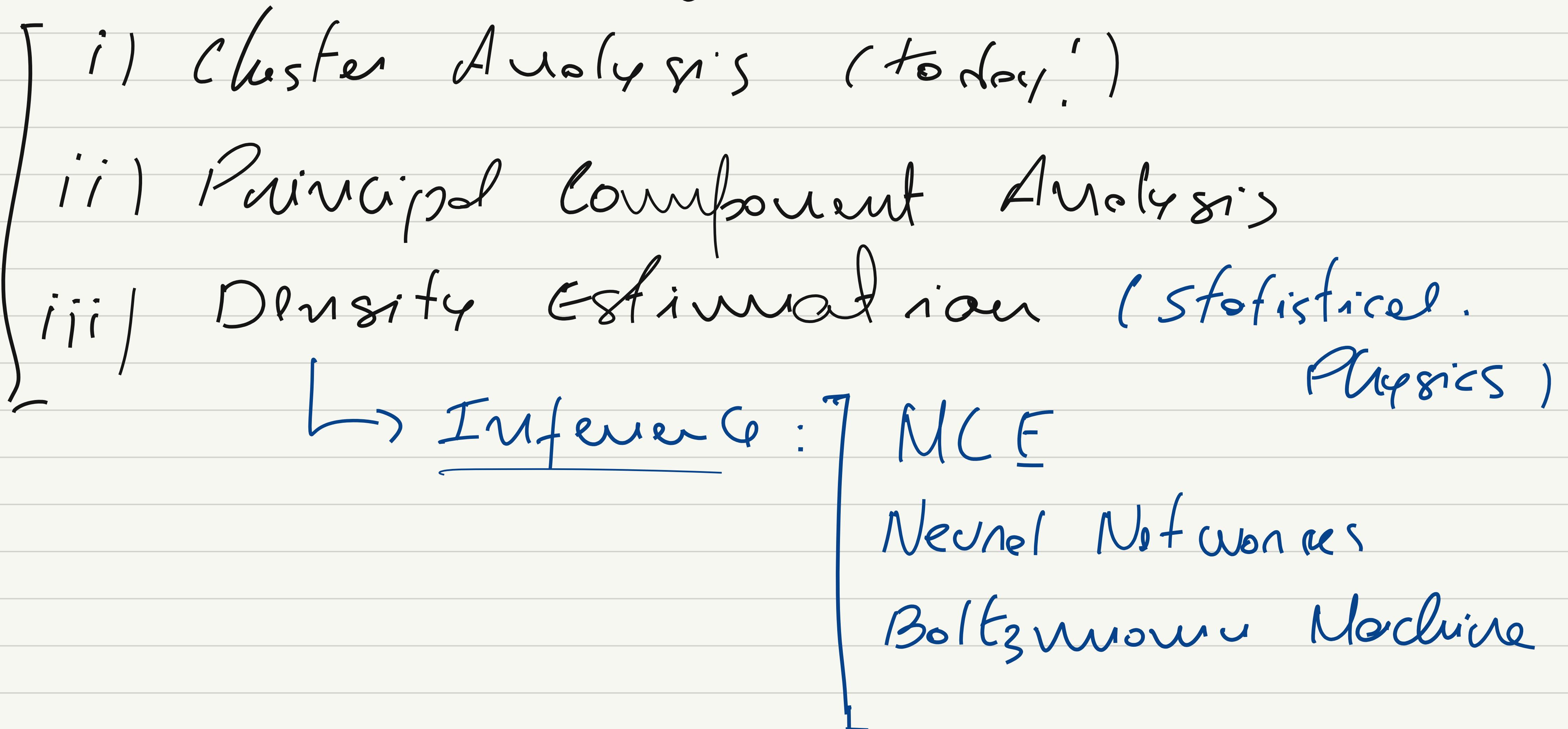
1. Unsupervised Learning

$$\mathcal{D} = \{(\vec{x}_i)\}_{i=1}^N, \quad \vec{x}_i \in \mathbb{R}^D \quad \epsilon \quad \mathcal{D} \in \mathbb{R}^{N \times D}$$

No labels $P(\vec{x}) \rightarrow$ Hidden Underlying
 N -instances

GOAL : Infer the properties $P(\vec{x})$
without a teacher.

Unsupervised learning \neq Descriptive Statistics



2. Clustering

- Clustering is a multivariate unsupervised machine learning method that seeks to group unlabelled data into clusters, given a similarity or distance metric. *(Rigid Body \Rightarrow Intentful Movement)*
- Clustering finds many applications throughout: data mining, data compression, data signal. Also is capable to reveal coarse features or high-level structures into unlabelled data.
- Is probably the simplest way to identify hidden structures in a dataset. //
- This is a vast and diverse field of machine learning, and there are a few considerations to be done in order to choose a particular clustering method
 - 1. Distribution of clusters (~~overlapping/noisy~~ or well separated clusters);
 - 2. The geometry of the data (flat or non-flat)
 - 3. Cluster size distribution (multiple sizes or uniform size) *[Performance]*
 - 4. The dimensionality of data .
 - 5. Computational efficiency of desired method.

Storage Options

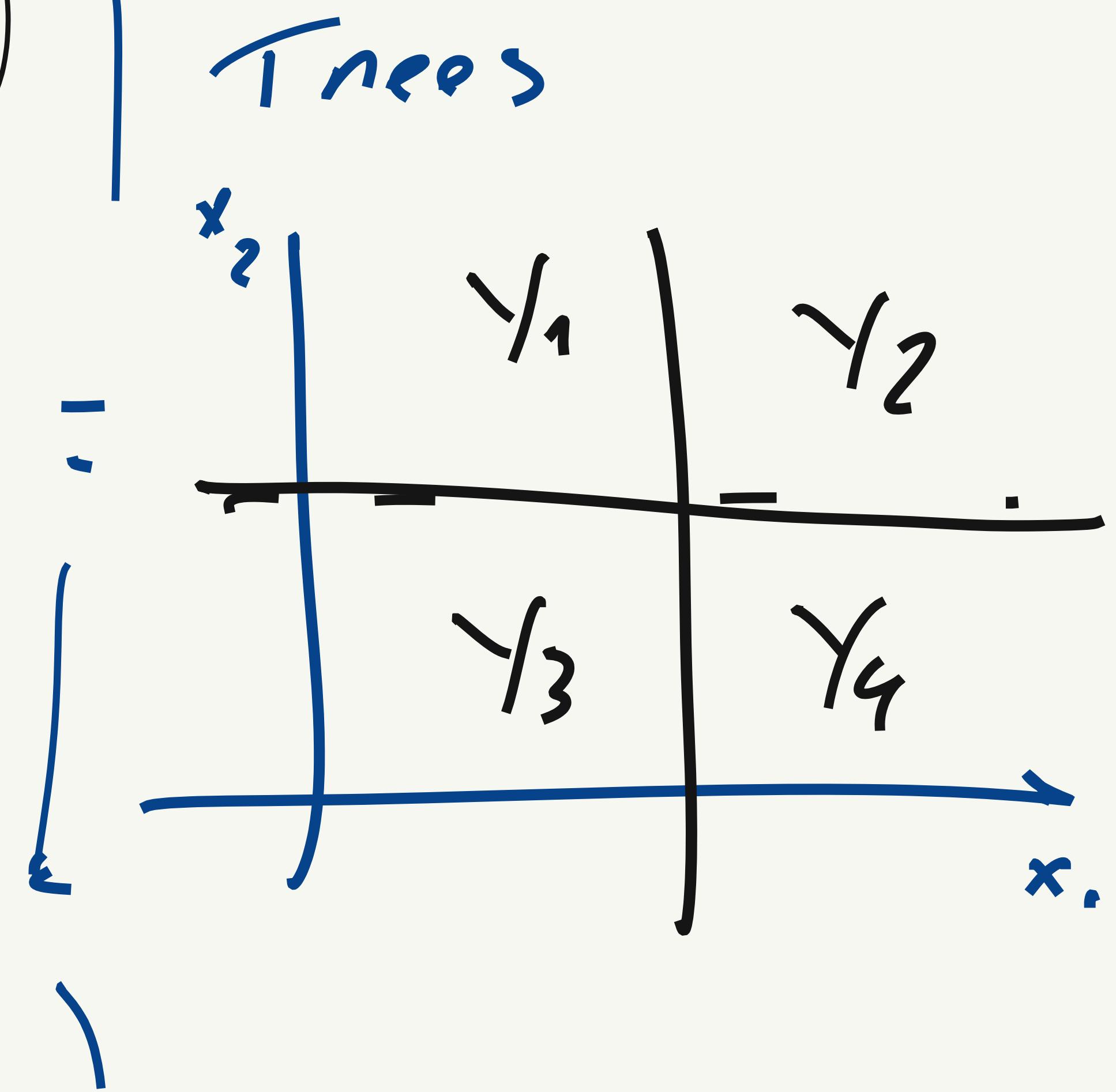
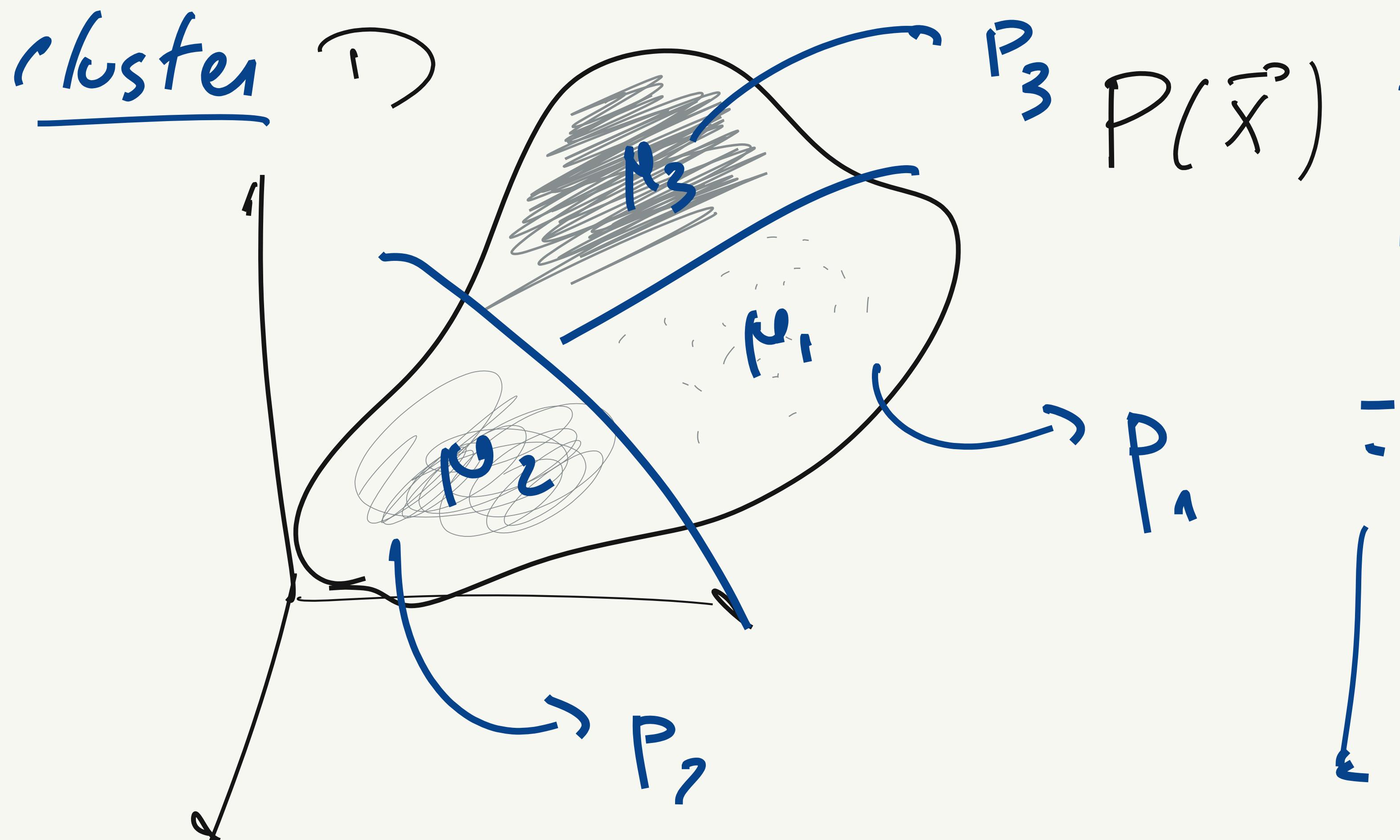
2. Clustering

- There many toy models/business cases that can be solved using clustering algorithms:
 1. Customer segmentation;
 2. Anomaly detection;
 3. Basket Analysis; *(Association Rules)*
 4. Churn Detection;
- In the physical sciences, clustering methods are of great usage in problems like:
 1. Detect celestial emission sources ✓
 2. Building entanglement classifiers ✓
 3. Infer groups of genes and proteins ✓
- We are going to study a few incarnations of the clustering problem along de course: K-Means, DBSCAN, Gaussian Mixtures. In these notes, we deal with K-Means and DBSCAN, also we show how K-Means can be understood as a limit case of Gaussian Mixtures before a deep dive in this model.

DBSCAN

✗

✗



2. Clustering

- Let C be the vector space formed by all datapoints in a sample, a clustering method attempts to produce a partition of this space, so that:

$$C = C_1 \cup C_2 \cup \dots \cup C_k$$

$$C_i \cap C_j = \emptyset \forall i \neq j$$

- Clustering methods are characterized by the presence of a dissimilarity measure. The most common are:

- **Euclidean Distance:** $d^2(\vec{x}_i, \vec{x}_j) = \sum_{k=1}^d (x_{i,k} - x_{j,k})^2;$

- **Mahalanobis Distance:** $d(\vec{x}_i, \vec{x}_j) = \sqrt{\sum_{k=1}^d |x_{i,k} - x_{j,k}|^2};$

- **Cosine Distance:** $d(\vec{x}_i, \vec{x}_j) = 1 - \frac{\vec{x}_i \cdot \vec{x}_j}{\|\vec{x}_i\| \|\vec{x}_j\|};$

- **Jaccard Distance:** $d(\vec{x}_i, \vec{x}_j) = 1 - \frac{\sum_{k=1}^d x_{i,k} x_{j,k}}{\sum_{k=1}^d x_{i,k} + \sum_{k=1}^d x_{j,k} - \sum_{k=1}^d x_{i,k} x_{j,k}};$

- Each distance is more adequate to a type of data, tabular, non-tabular, categorical.

- The fundamental problem of clustering is given N datapoints and K reference points, called centroid of the cluster, we need to find the cluster mean $\{\vec{\mu}\}$ that minimizes the objective function:

$$\mathcal{C}(\{x, \vec{\mu}\}) = \sum_{k=1}^K \sum_{n=1}^N p_{n,k} D(\{x, \vec{\mu}\})$$

- Where D is the chosen dissimilarity measure and $p_{n,k}$ is the probability of a given datapoint to be assigned to a given cluster and usually is a binary variable being 0 or 1.

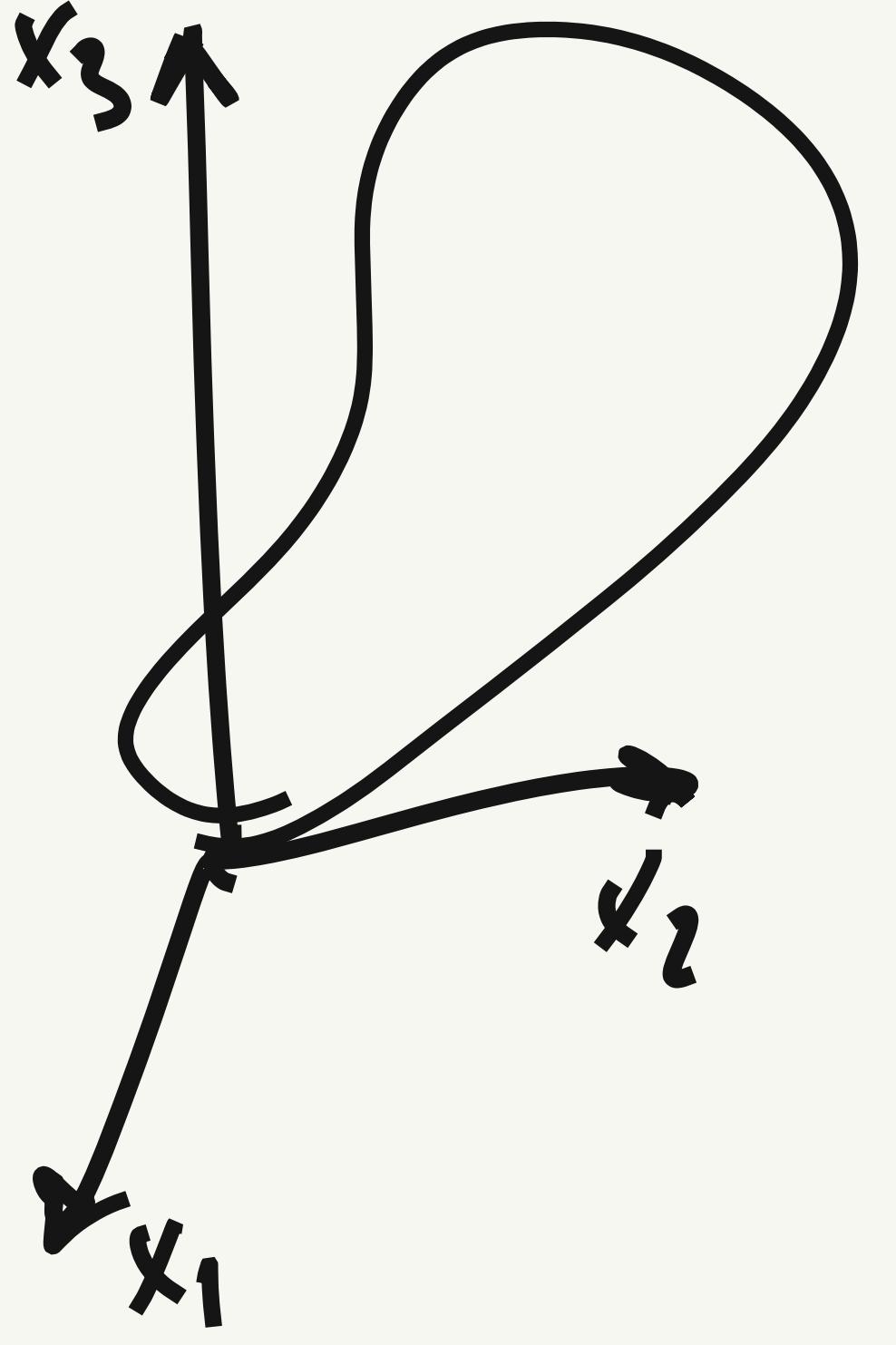
2. Clustering

- Most of the commonly used clustering methods takes as input a dissimilarity matrix, defined over the above defined distances.
- The choice of the metrics, depends upon the type of data and the optimization algorithm.
- Pros and Cons
 - **Euclidean:** Easy to interpret and for optimization. Geometry dependent. Works for quantitative variables.
 - **Mahalanobis:** Not good for optimization. Geometry Independent. Take correlations into account. Works for quantitative variables.
 - **Cosine Distance:** Geometry dependent. Works for any kind of feature represented by vectors. *(more flat)*
 - **Jaccard Distance:** Hard to calculate and interpret. Works for categorical variable
- Due to different scales or even importances, the total dissimilarity has to be a weighted sum of the distances:

$10^9, 10^9, 10^9, 10^9, \dots$

$$D(x_i, x_j) = \sum_{k=1}^K w_k d_k(x_{ik}, x_{jk})$$

d. wⁿ 10^3
↳ weight ↳ 10^3



Let, x_1, x_2, x_3 Variables with
distinct units and dimensions

\rightarrow Preprocessing } scaling
Non uniform
Power transfor.
 \hookrightarrow Uniform
Phase Space.

\rightarrow Non-uniform describ'd. by
discrete prob. non-smooth

The weight of a feature

$$\bar{D} = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N D(x_i, x_j)$$

$$= \sum_{u=1}^k w_u \cdot \bar{d}_u ; \quad \bar{d}_u = \frac{1}{N_u} \sum_{i=1}^{N_u} \sum_{j=1}^{N_u} d_u(x_i, x_j)$$

If we have a uniform prior spec: $w_u = \frac{1}{\bar{d}_u}$

$$w_u = f(\bar{d}_u) , \quad \prod_u p_u(\bar{d}_u)$$

Consider a feature space smooth
and P quantitative features

$$E \cdot D := D_E(x_i, x_j) = \sum w_u (x_{iu} - x_{ju})^2$$

clusters or
discusses by
minimizing
dispersion
of
points

$$\bar{d}_u = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^{z_h} (x_{iu} - x_{ju})^2$$

$$= \partial J_{\text{om}}(\vec{x}_h)$$

* Combinatorial optimization

find clusters that minimize

within cluster scatter point : $W(c) = \frac{1}{2} \sum_{k=1}^K \sum_{c(i)=u} \sum_{c(j)=u} d(x_i, x_j)$

Total scatter point metric :

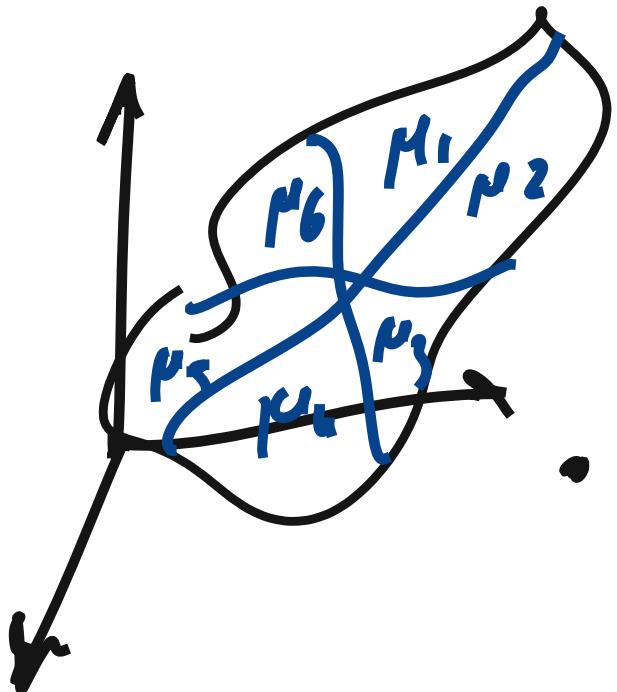
$$\begin{aligned} T(c) &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n d_{ij} = \frac{1}{2} \sum_{k=1}^K \sum_{c(i)=u} \left(\sum_{c(j)=u} d_{ij} + \sum_{c(j) \neq u} d_{ij} \right) \\ &= W(c) + B(c) \Rightarrow \downarrow W(c), \uparrow B(c) \end{aligned}$$

3. K-Means

- Centroid : $\bar{\mu}_u$ — Center-of-mass

of a data subset

$$W(c) = \frac{1}{2} \sum_{k=1}^K \sum_{c(i) = k} \sum_{c(j) = u} (x_{iu} - \bar{x}_{iu})^2$$



- Find the optimal number of centroids
- Find the optimal position centroids.

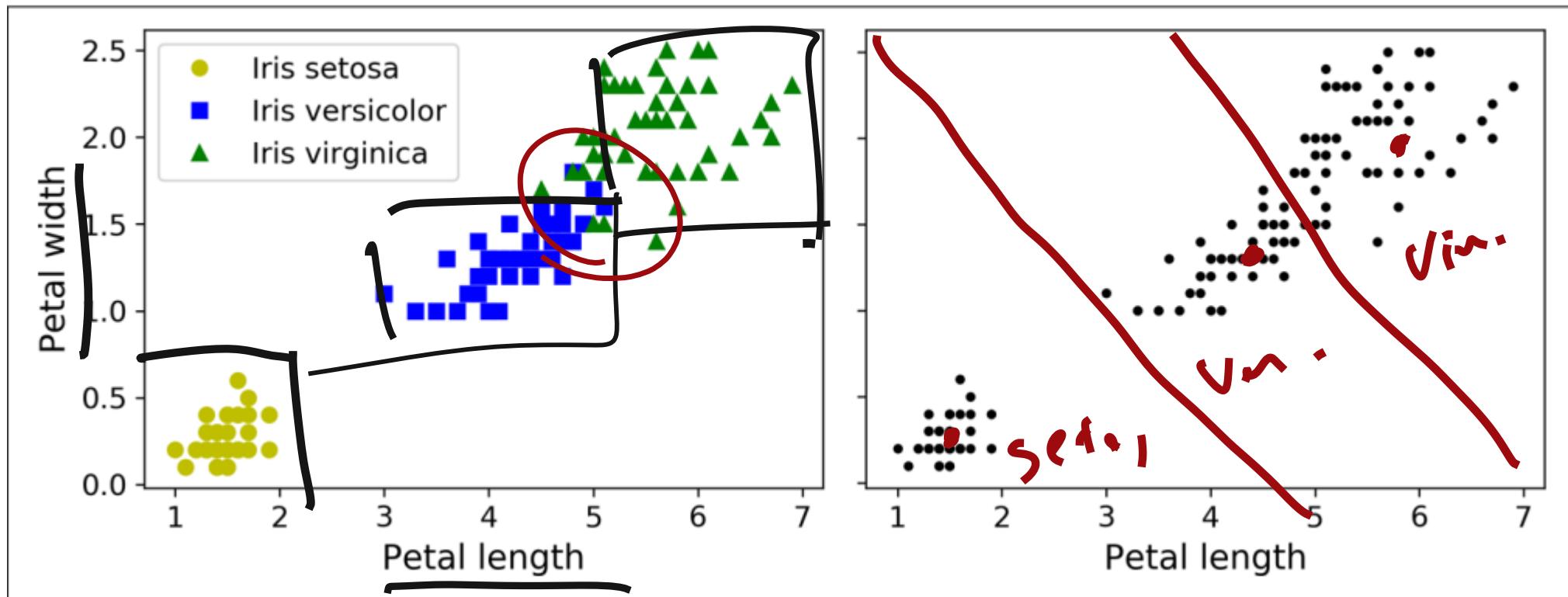
The K-Means method has two steps

1- Expectation : $\bar{\mu}_k = \frac{1}{N_k} \sum_n \pi_{n,k} \vec{x}_n$

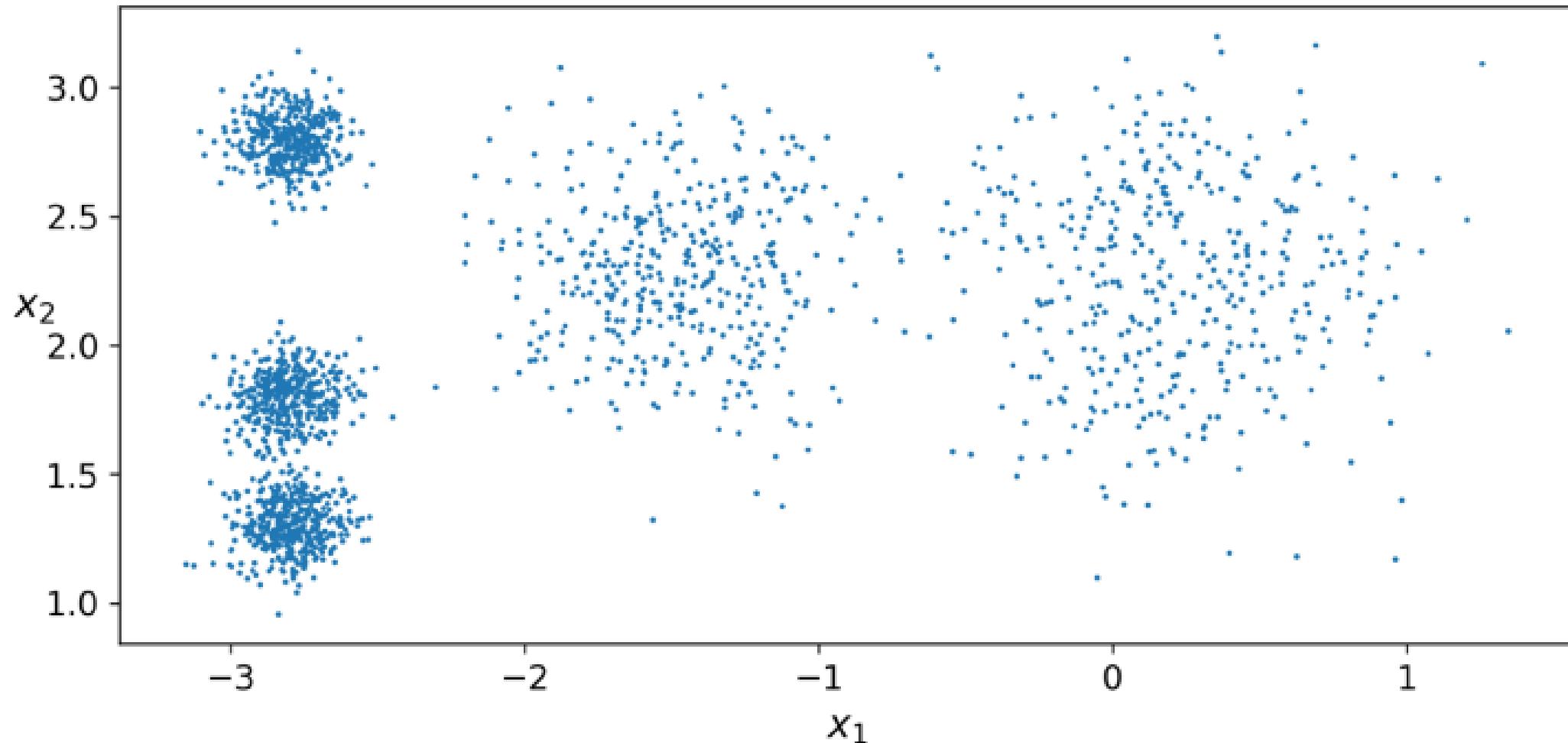
$$\pi_{n,k} = \begin{cases} 1 & \text{if } \vec{x}_n \in C(k) \\ 0 & \text{if not} \end{cases}$$

2- Minimization : Given $\bar{\mu}_k$ wcc

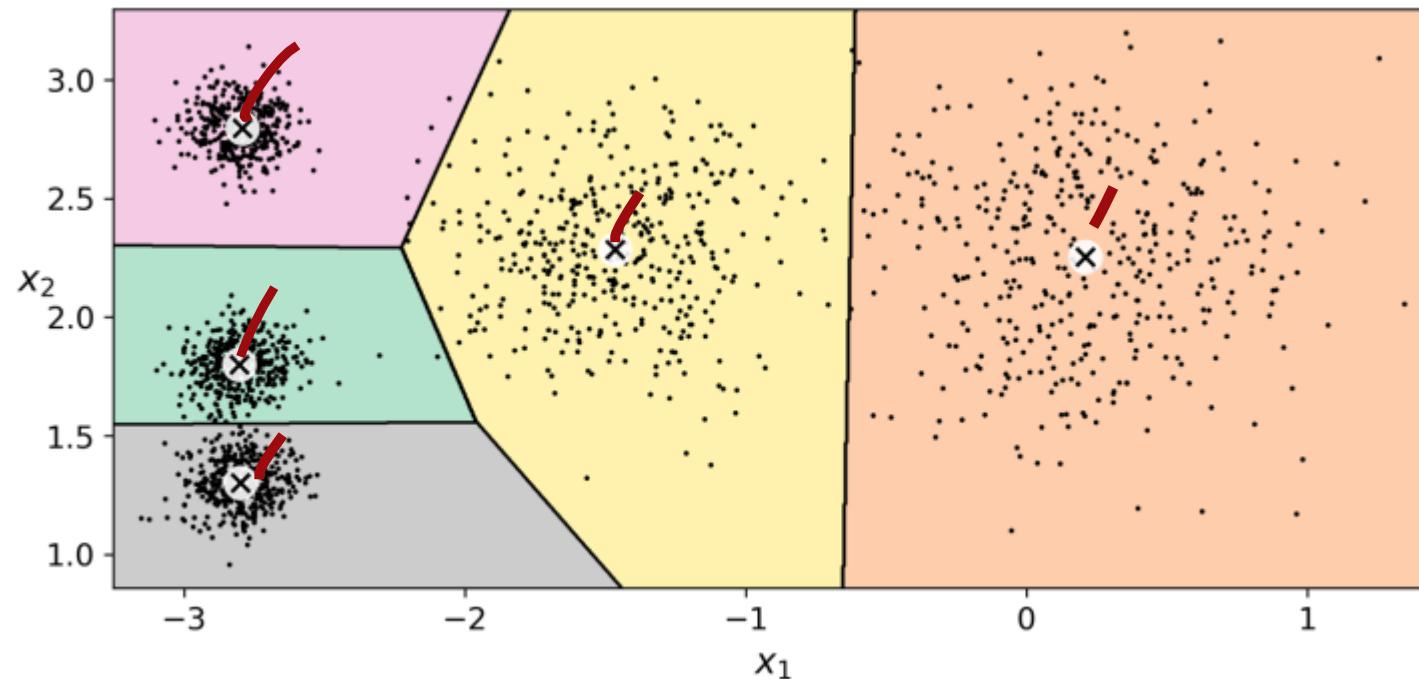
3. K-Means



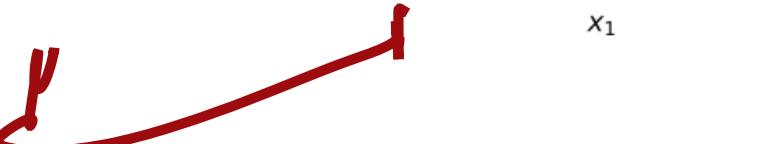
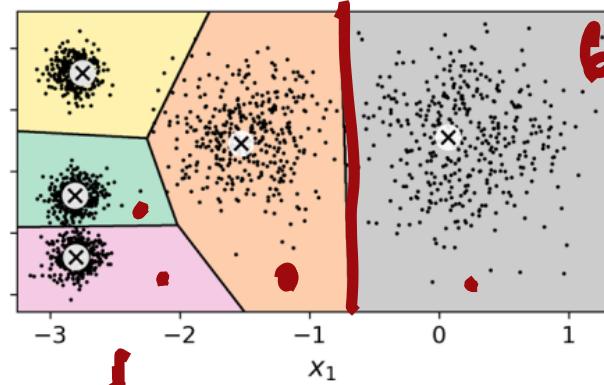
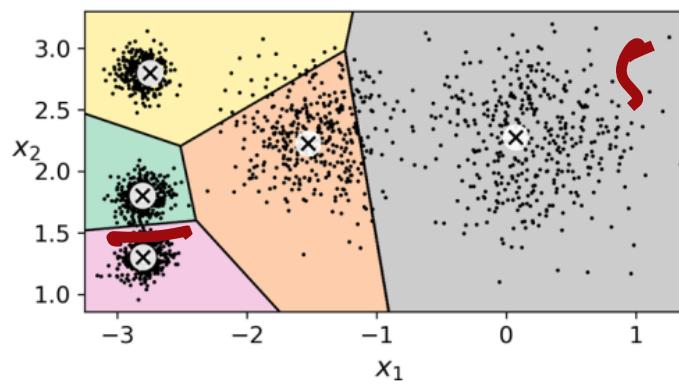
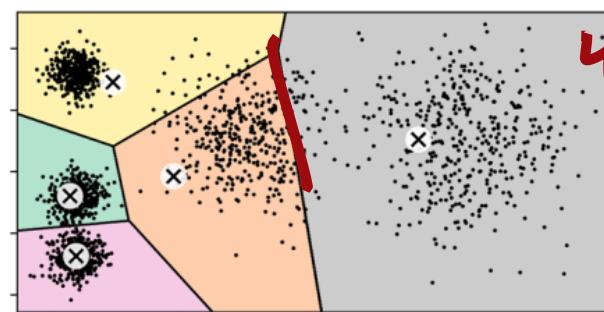
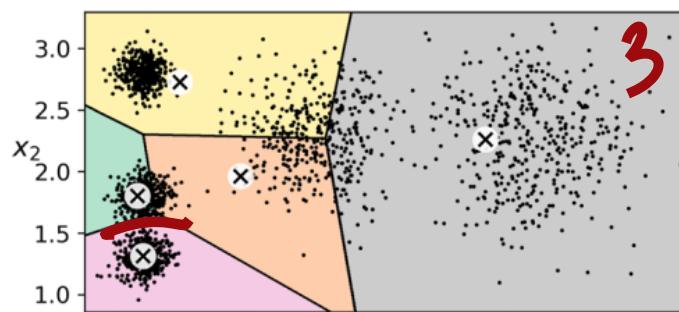
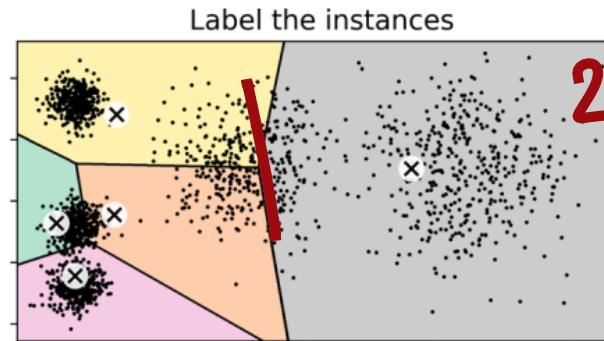
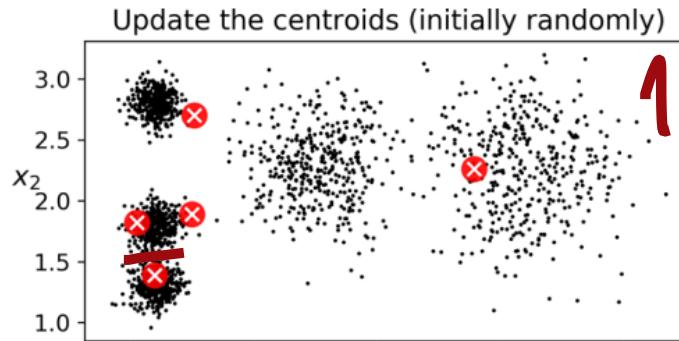
3. K-Means



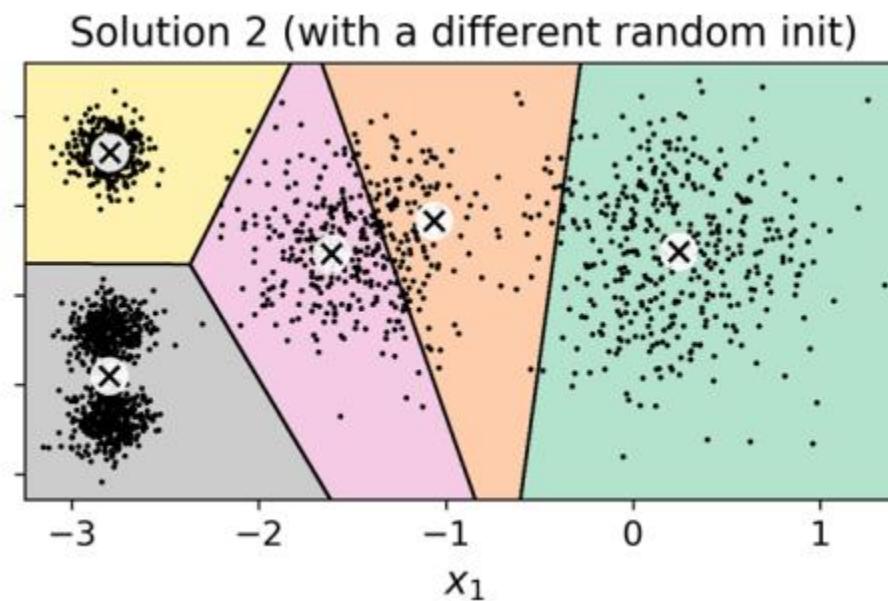
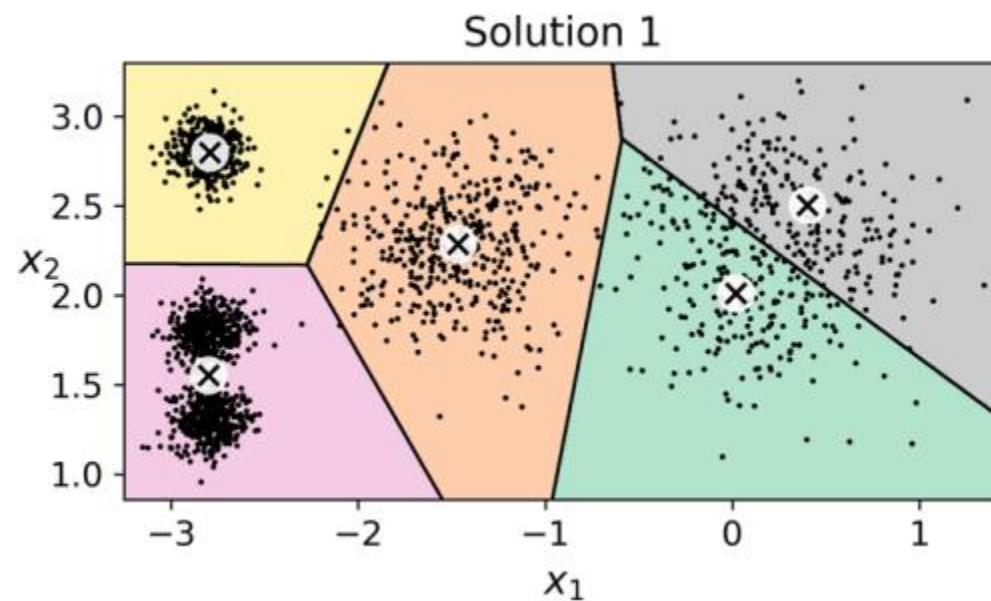
3. K-Means



3. K-Means

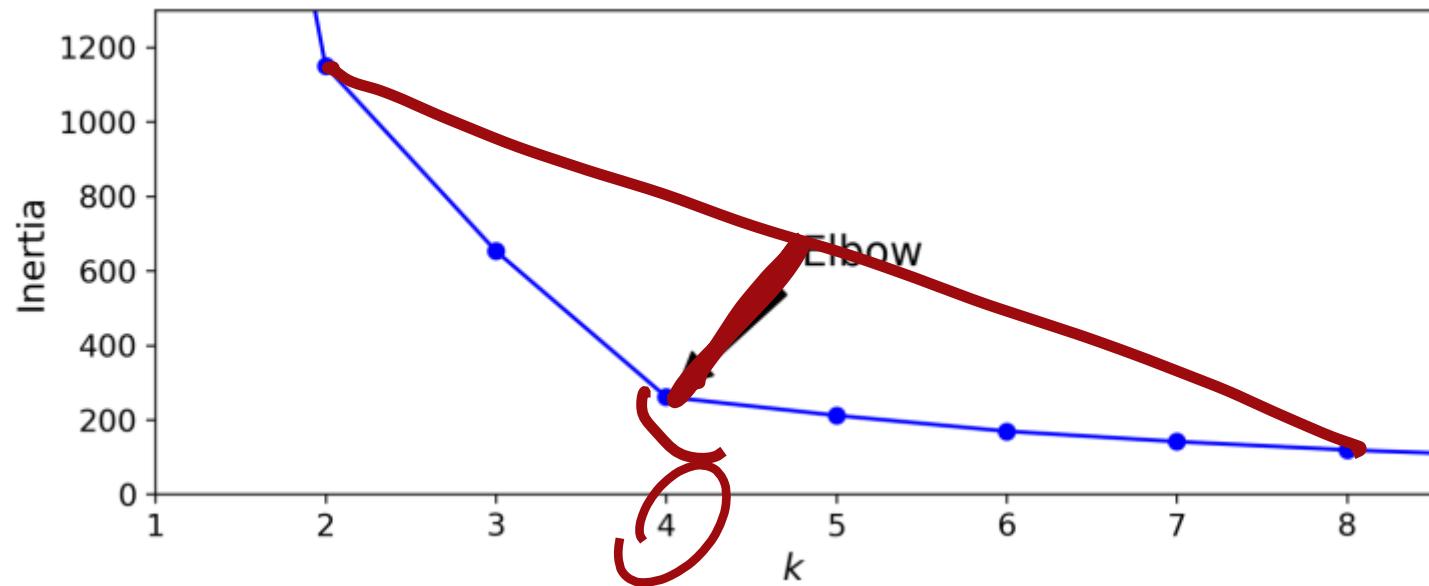


3. K-Means



3. K-Means

Elbow Method



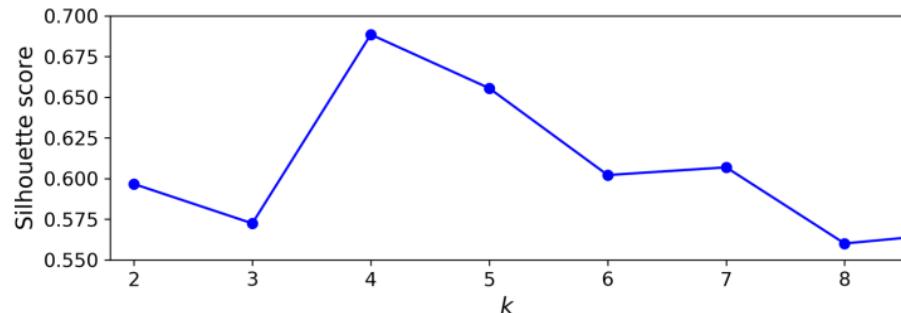
How many clusters

$$J = \sum_{k=1}^n \sum_{i \in C_k} (x_i - \bar{\mu}_k)^2 \rightarrow \text{sum-of-squared distances}$$

of the data set

3. K-Means

Silhouette Method
 $s = \max \tilde{s}$, $\tilde{s} = \frac{1}{n} \sum s_u(i)$



$$a_u = \frac{1}{|N_u|-1} \sum_{j \in C_u} d^2(x_i, x_j)$$

$$b_u = \min_{j \neq u} \frac{1}{|C_p|} \sum_{j \in C_p} d^2(x_i, x_j)$$

$$s_u(i) = \frac{b_u(i) - a_u(i)}{\max\{a_u(i), b_u(i)\}}$$

