# Lecture 17: Clustering and density estimation

## Professor Ilias Bilionis

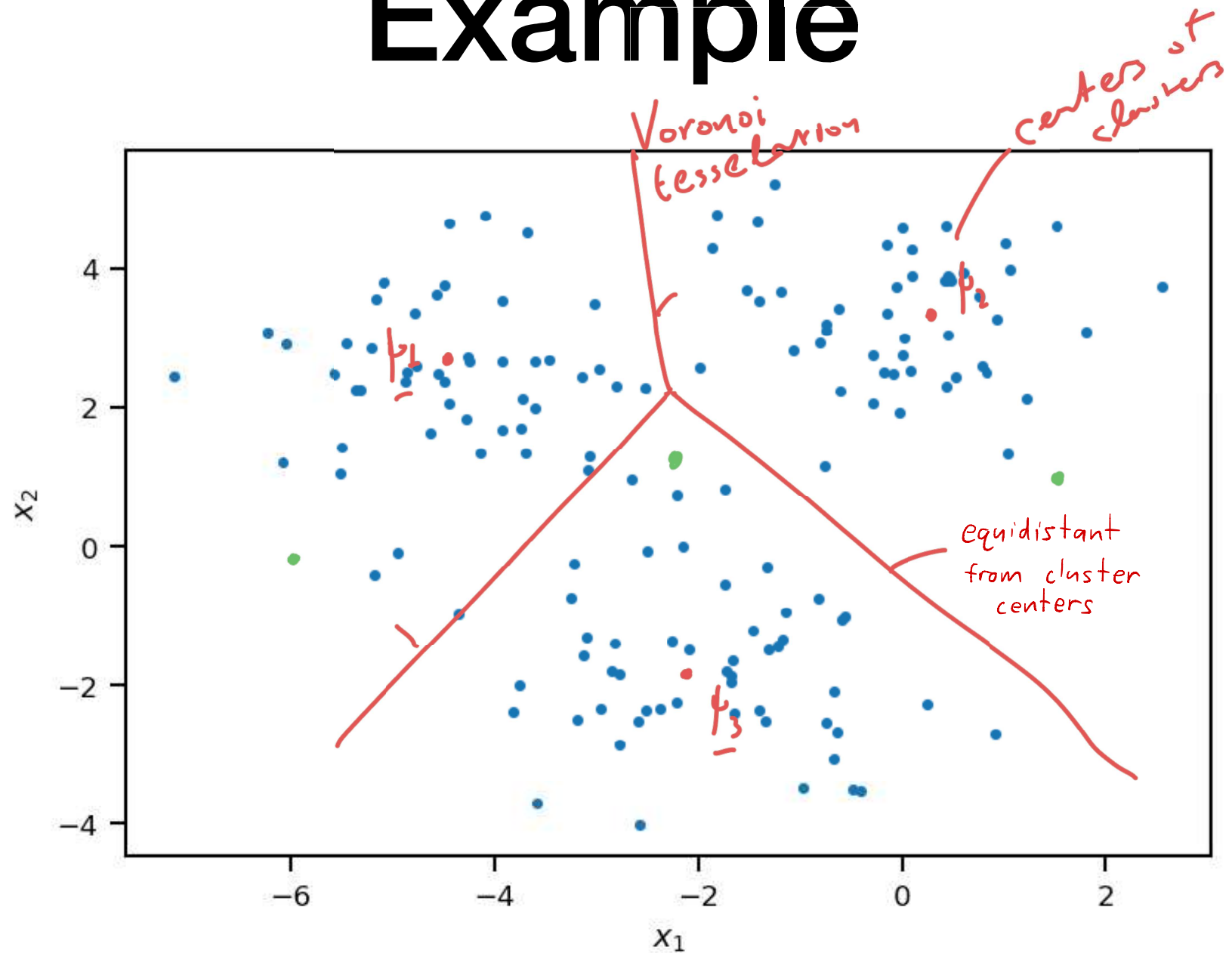## Clustering using k-means

# Clustering

Your are given n observations:

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \ldots, \mathbf{x}_n\}$$

(inputs, features, …)

**Problem**: Separate the data into K groups? How many such groups exist?

**PREDICTIVE SCIENCE LABORATORY**
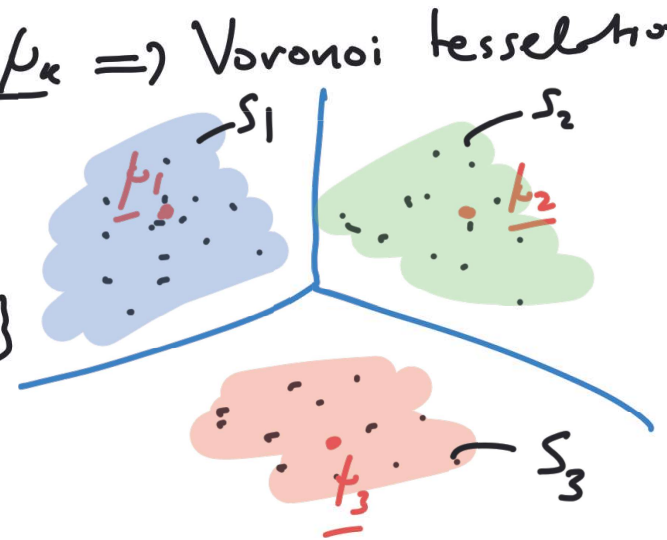
# Example

# K-means objective

$K$ clusters w/ centers $\underline{\mu_1}, \ldots, \underline{\mu_k} \Rightarrow$ Voronoi tesselation

$\underline{x}_{1:n} = (\underline{x_1}, \underline{x_2}, \ldots, \underline{x_n})$

$S_1 \subset \{\underline{x_1}, \ldots, \underline{x_n}\}, \quad S_2 \subset \{\underline{x_1}, \ldots, \underline{x_n}\}$

$S_3, \ldots, S_k$

$$\min_{\underline{\mu_1}, \ldots, \underline{\mu_k}} \sum_{i=1}^{K} \sum_{\underline{x} \in S_i} \|\underline{\mu_i} - \underline{x}\|^2$$

$S_1$  $\mu_1$

$S_2$  $\mu_2$

$\mu_3$  $S_3$

# Standard k-means algorithm

$$\min_{\underline{\mu}_1,\dots,\underline{\mu}_k} \sum_{i=1}^{K} \sum_{\underline{x} \in S_i} \|\underline{\mu}_i - \underline{x}\|^2$$
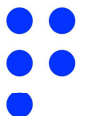
1. Start by randomly choosing $\underline{\mu}_1^{(1)}, \dots, \underline{\mu}_k^{(1)}$. $t \leftarrow 1$.

2. <u>Assignment step</u>: $S_1^{(t)}, \dots, S_k^{(t)}$

   $$S_i^{(t)} = \{ \underline{x}_j \in \{\underline{x}_1, \dots, \underline{x}_n\} : \|\underline{x}_j - \underline{\mu}_i^{(t)}\| \leq \|\underline{x}_j - \underline{\mu}_r^{(t)}\| \quad r \neq i \}$$

3. <u>Update step</u>:

   $$\underline{\mu}_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \cdot \sum_{\underline{x} \in S_i^{(t)}} \underline{x}$$

   <span style="color:red"># points in cluster</span>     <span style="color:red">} empirical mean</span>

4. <u>Convergence test</u>: If $\|\underline{\mu}_i^{(t+1)} - \underline{\mu}_i^{(t)}\| < \varepsilon \quad \forall i = 1, \dots, k$, then STOP. Otherwise GOTO 2.

   <span style="color:red">threshold</span>
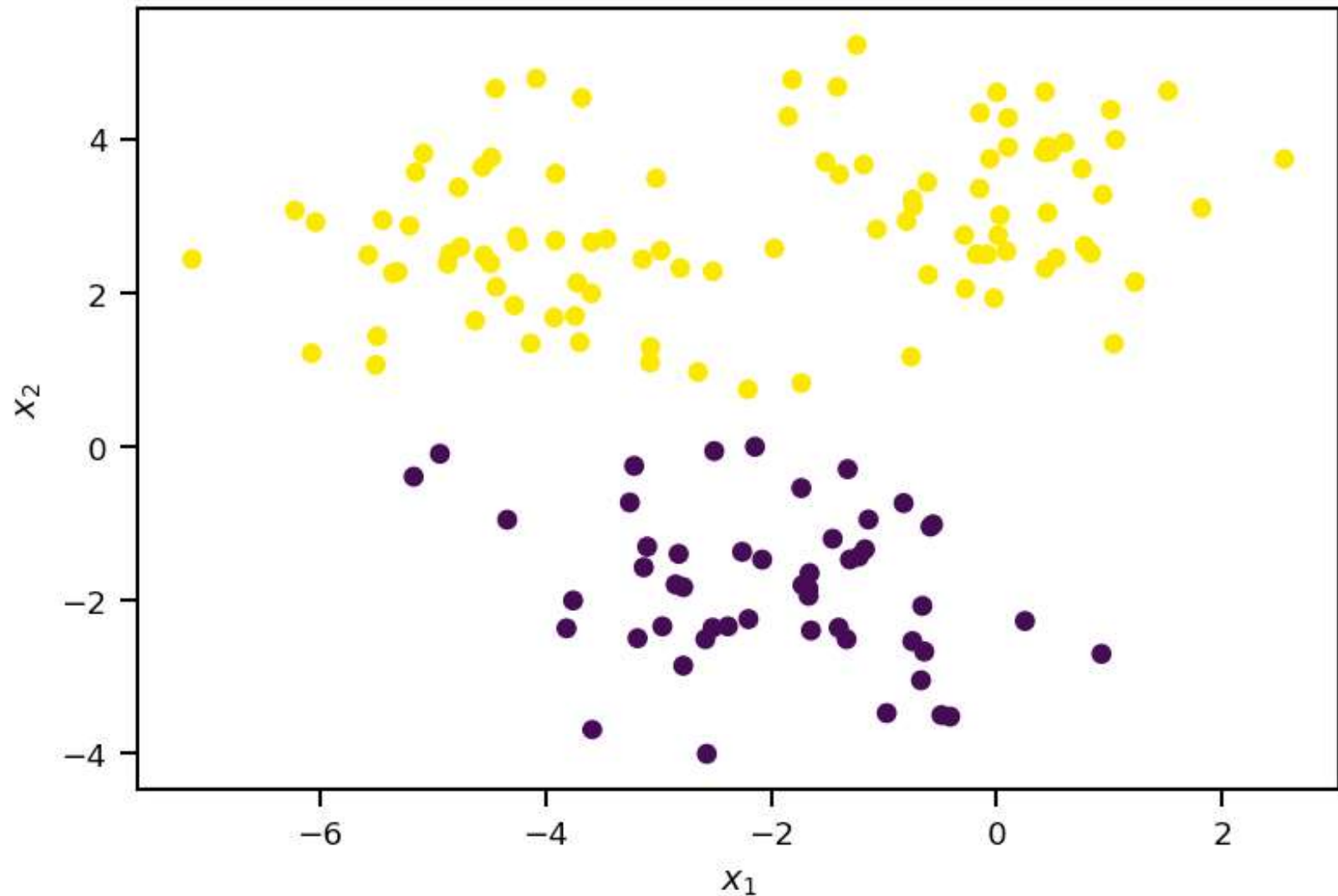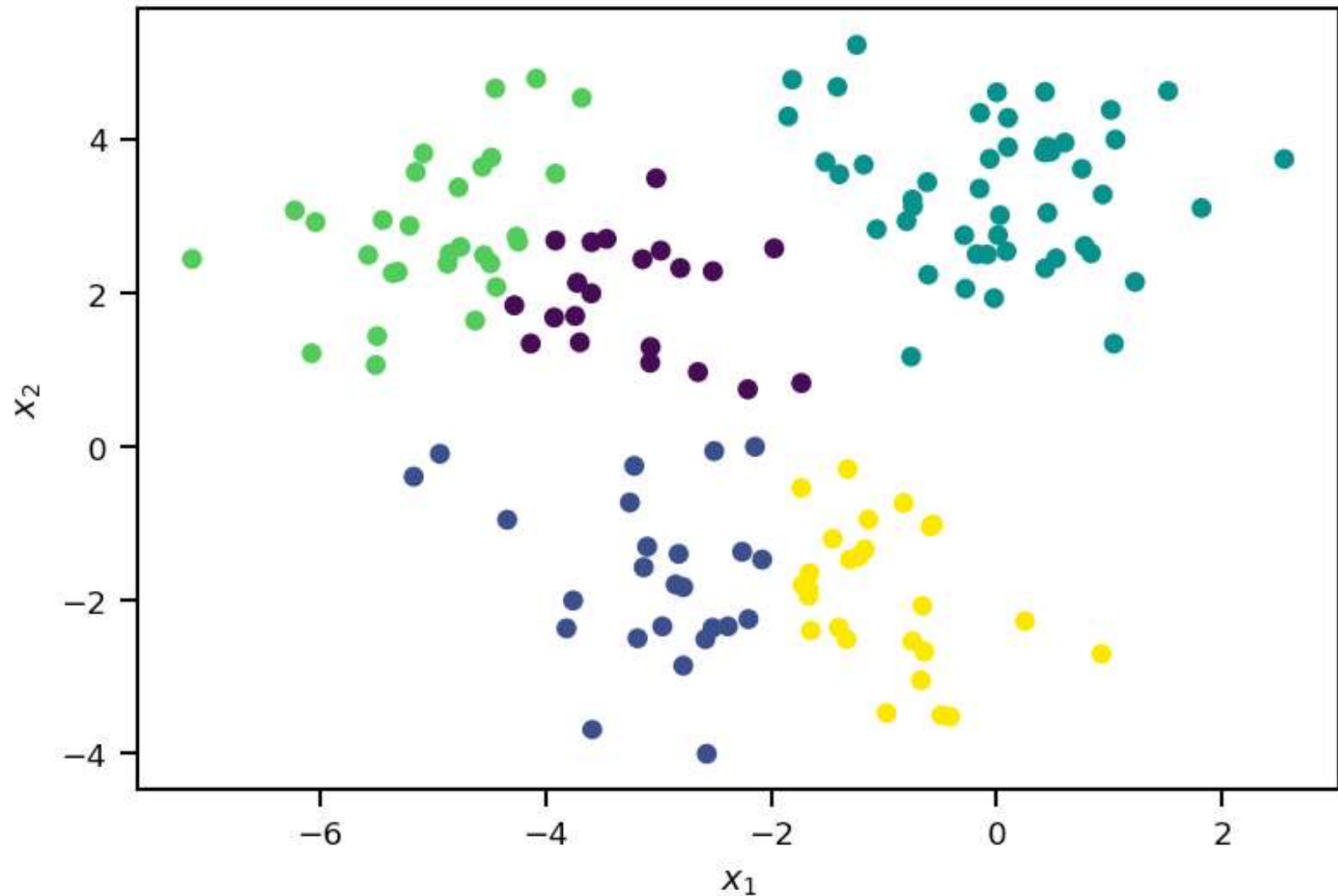
# Example

# Example

PREDICTIVE
SCIENCE LABORATORY

# What if I used two clusters?

# What if I used five clusters?

# Limitations of k-means

- How many clusters?

- Assumes spherical clusters.

- Cannot be applied to high-dimensional datasets, e.g., images.

**PREDICTIVE
SCIENCE LABORATORY**

# Beyond k-means

- Clustering is related to density estimation.

- Idea:

  - Make hypothesis about how data are generated.

  - Train your model.

  - Let the structure arise naturally.