# Clustering
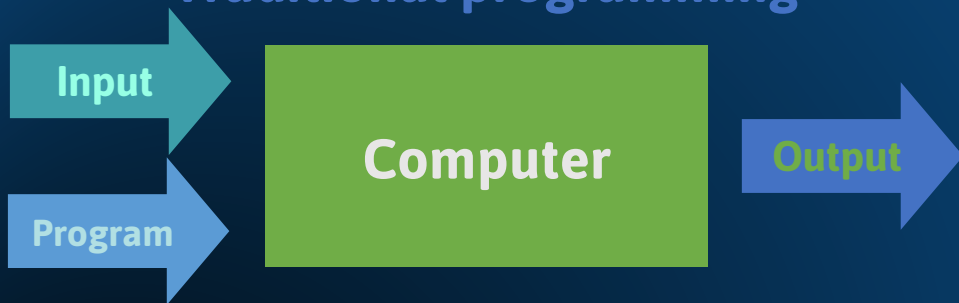
Cláudia Antunes

Instituto Superior Técnico – Universidade de Lisboa

# Unsupervised Learning

**Traditional programming**

Input

Computer
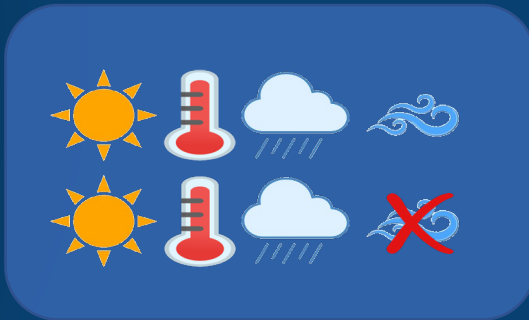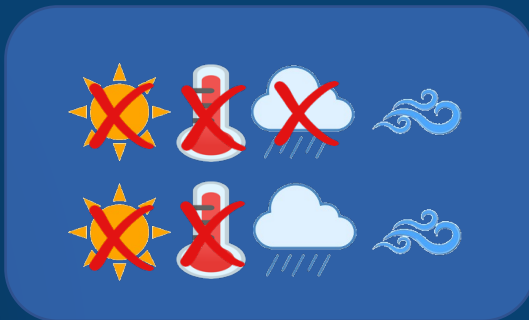
Output

Program

**Machine Learning**

Input

Computer

Output

Information

# CLUSTERING



**Dataset**

No Target Variable

# Assessment

Cohesion

Separation

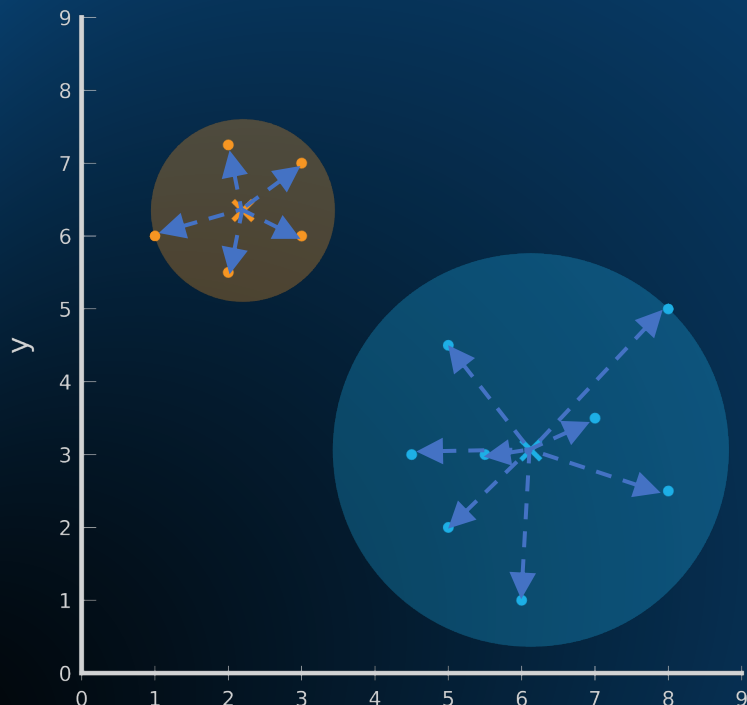$$radius(C_i) =$$
$$max_x\{d(x, \mu_i): x \in C_i\}$$

$$max(C_i) =$$
$$max_{x,y}\{d(x, y): x, y \in C_i\}$$

$$avg\ dist(C_i\ ) = \frac{1}{|C_i|} \sum_{x \in C_i} d(x, \mu_i)$$

$$avg\ dist(C_i) =$$
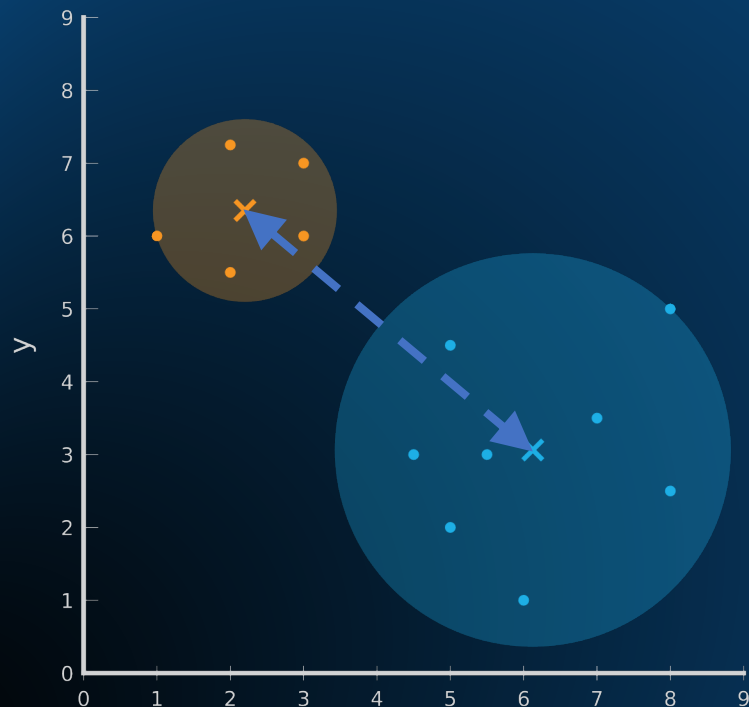$$\frac{1}{|C_i||C_i - 1|} \sum_{\substack{x,y \in C_i \\ x \neq y}} d(x, y)$$

$$MSE = \frac{1}{N} \sum_{i} \sum_{x \in C_i} d(\mu_i, x)^2$$
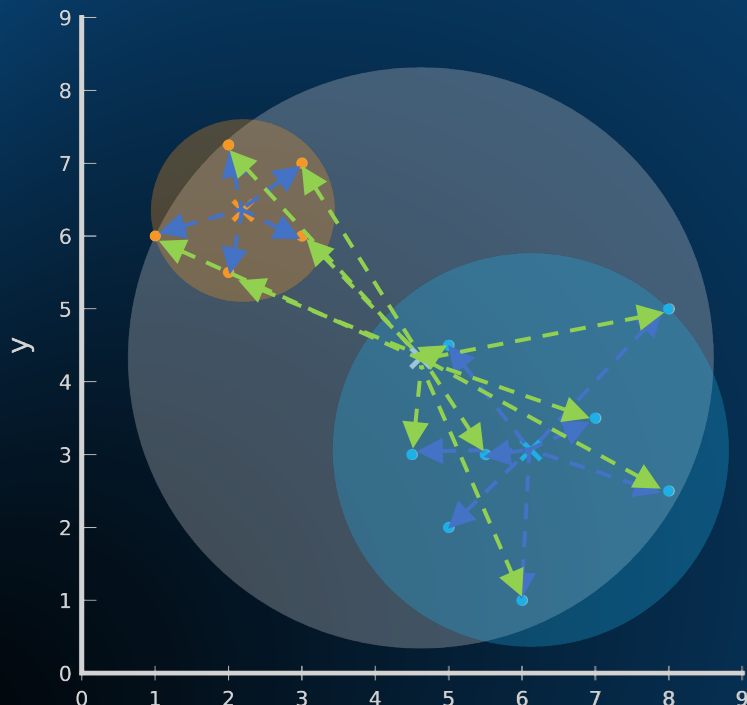
$$MAE = \frac{1}{N} \sum_{i} \sum_{x \in C_i} d(\mu_i, x)$$

$$d(C_i, C_j) = d(\mu_i, \mu_j)$$

$$slink(C_i, C_j) =$$
$$min_{x,y}\{d(x, y): x \in C_i, y \in C_j\}$$

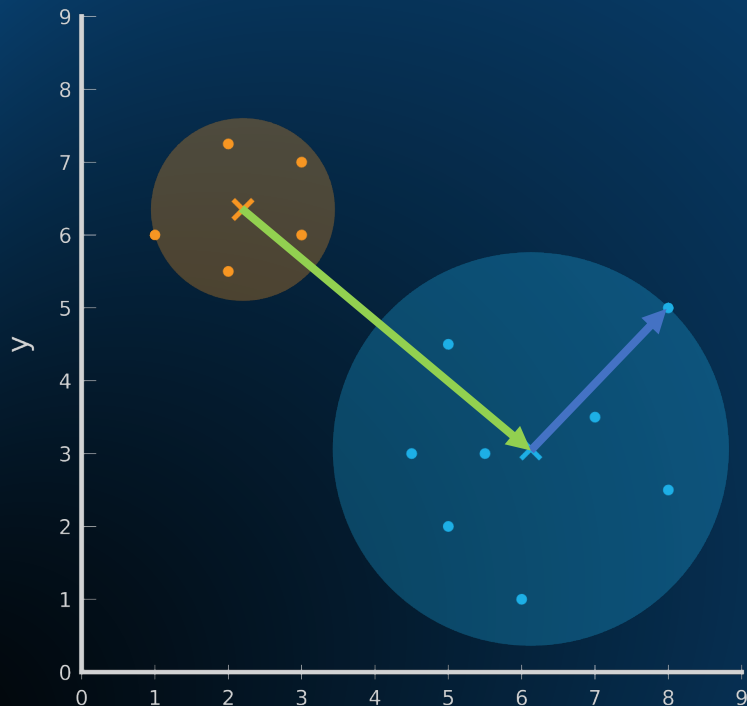$$clink(C_i, C_j) =$$
$$max_{x,y}\{d(x, y): x \in C_i, y \in C_j\}$$

## Ward's distance
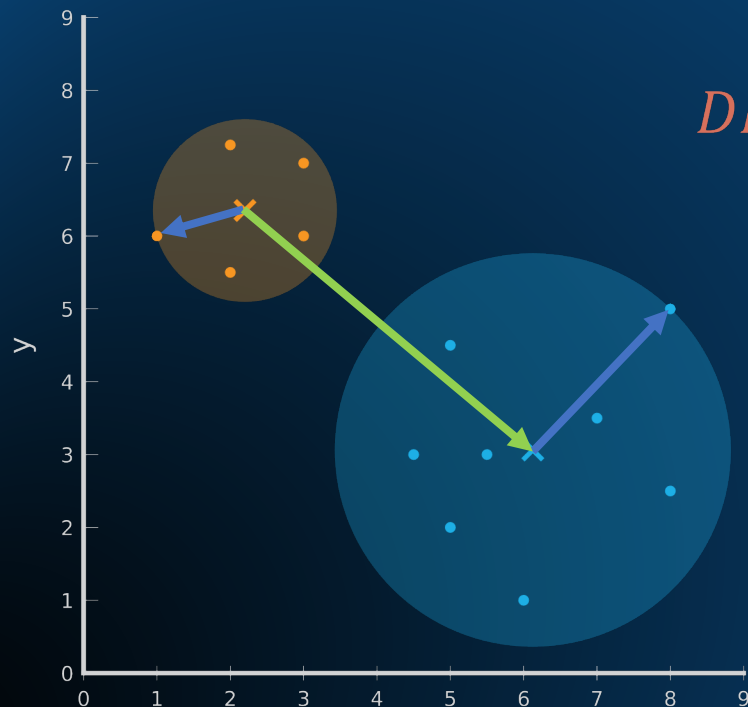
$$d(C_i, C_j) =$$

$$\sum_{x \in C_i} d(x, \mu_i)^2$$

$$+ \sum_{x \in C_j} d(x, \mu)^2$$
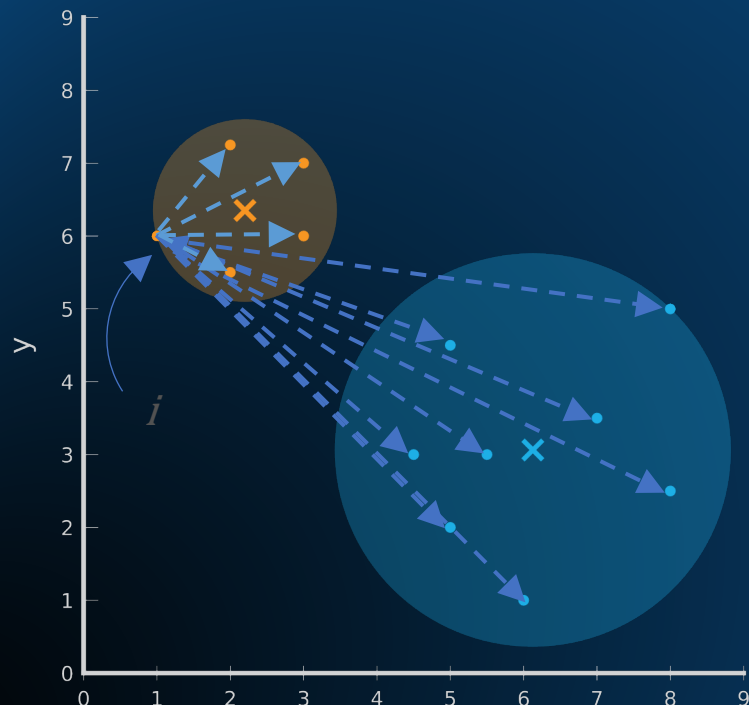
$$- \sum_{x \in C_{ij}} d(x, \mu_{ij})^2$$

$$DI_C =$$
$$\frac{min_{i,j}\{d(C_i, C_j): 1 \leq i, j \leq k\}}{max_i\{diam(C_i): 1 \leq i \leq k\}}$$
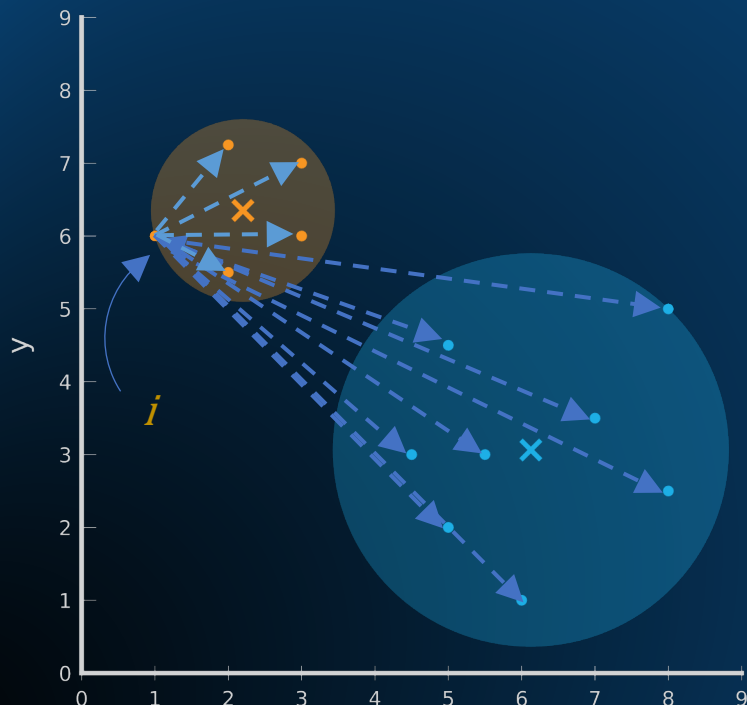
$$\mu_{in}(x_i) =$$

$$\frac{1}{|C_{x_i}| - 1} \sum_{x_j \in C_{x_i}, i \neq j} d(x_i, x_j)$$

$$\mu_{out}(x_i) =$$

$$min_k \frac{1}{|C_k|} \sum_{\substack{x_j \in C_k \\ C_k \neq C_{x_i}}} d(x_i, x_j)$$

$$s(x_i) = \frac{\mu_{out}(x_i) - \mu_{min}(x_i)}{\max\{\mu_{out}(x_i), \mu_{in}(x_i)\}}$$
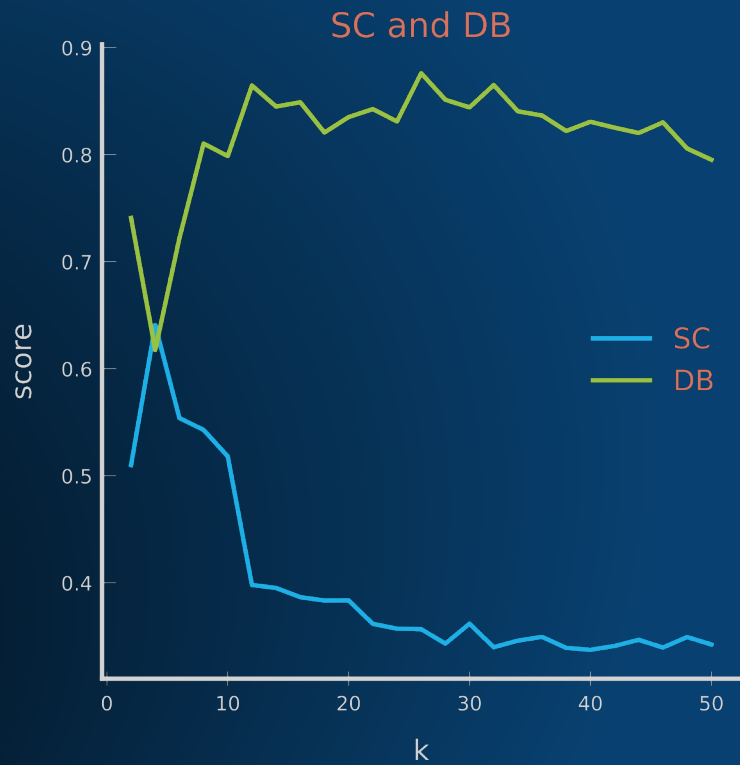
$$s(x_i) = \frac{\mu_{out}(x_i) - \mu_{min}(x_i)}{\max\{\mu_{out}(x_i), \mu_{in}(x_i)\}}$$
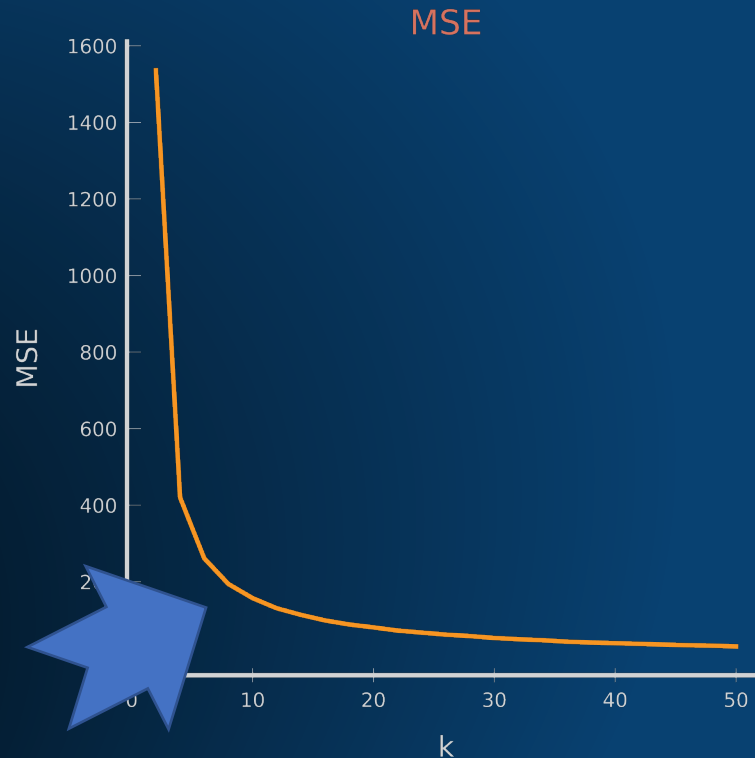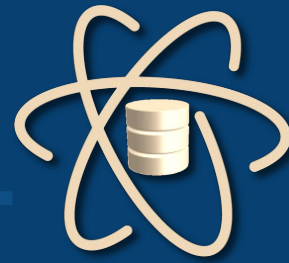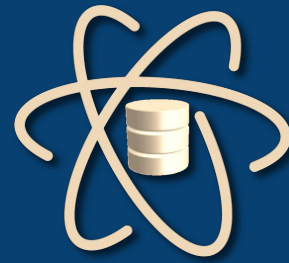
$$SC = \frac{1}{|D|} \sum_{x_i \in D} s(x_i)$$

$0.75 \leq SC \leq 1.00 \rightarrow$ excellent
$0.50 \leq SC < 0.75 \rightarrow$ good
$0.25 \leq SC < 0.50 \rightarrow$ weak
$SC < 0.25 \qquad \rightarrow$ no structure

# Choosing K



SC and DB

# CHOOSING K – THE ELBOW METHOD

# Thank you!