# K-Means

$$C_1 \cup C_2 \cup \ldots \cup C_K = \{1, \ldots, n\}$$
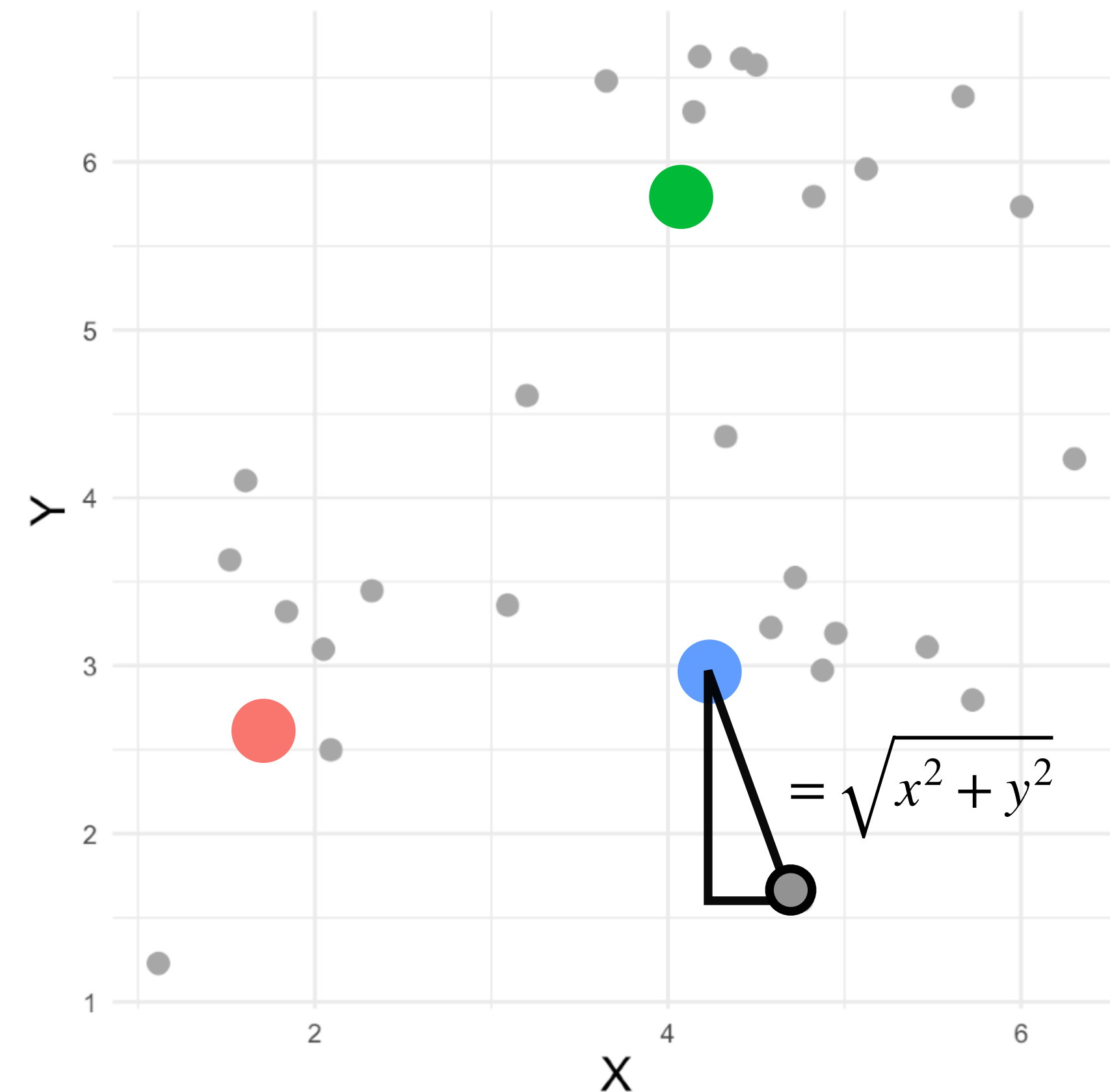
$$C_k \cap C_{k'} = \varnothing \text{ for all } k \neq k'$$

$$\min_{C_1, \ldots, C_K} \left\{ \sum_{k=1}^{K} W(C_k) \right\}$$

within cluster variance

$$\text{where } W(C_k) = \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^{p} (x_{ij} - x_{i'j})^2$$

squared Euclidean distance



$= \sqrt{x^2 + y^2}$

# K-Means: Algorithm

1. Choose **k** random points as cluster centers

2. For each data point, assign it the cluster whose centroid is the closest

3. Using these assignments, recalculate the centers

4. Reiterate from step (2) until **convergence**:
   - cluster membership does not change
   - center only changes very very little



**Data**

**Step 1**

**Iteration 1, Step 2a**

**Iteration 1, Step 2b**

**Iteration 2, Step 2a**

**Final Results**