# Cluster Analysis: Partitioning Methods

CS5483 Data Warehousing and Data Mining

# Group similar tuples together

$Z_2$

4

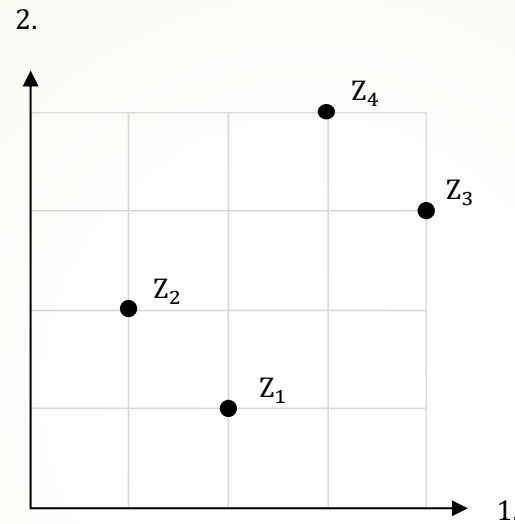|     | $Z_1$ | $Z_2$ |
|-----|-------|-------|
| 1.  | 2     | 1     |
| 2.  | 1     | 2     |
| 3.  | 3     | 4     |
| 4.  | 4     | 3     |

3

2

1

$Z_1$

Example:

- ➤ $Z_i$: intensity of $i$-th pixel in a picture
- ➤ Clustering the tuples identifies images of same/similar objects.

# Group similar features together

$$
\begin{array}{c c c c c}
 & Z_1 & Z_2 & Z_3 & Z_4 \\
1. & 2 & 1 & 3 & 4 \\
2. & 1 & 2 & 4 & 3 \\
\end{array}
$$



Example:

- $Z_i$: expression level of gene $i$
- Clustering the features identifies co-expressed genes.

# Partitioning method

- Input: A set $D := \{\boldsymbol{p}_i\}_{i=1}^{n}$ of data points

- Output: A set $\{C_j\}_{j=1}^{k}$ of non-empty disjoint clusters that partition $D$.

- Challenges:
  - there are often **too many data points**, and
  - the **d_____** can be **too high** to visualize.

- Need a mathematical criteria to automate clustering.

# Centroid-based method
## Model assumption

- Suppose:
  1. There is a typical point (**c**_____ **c**_____) in each cluster.
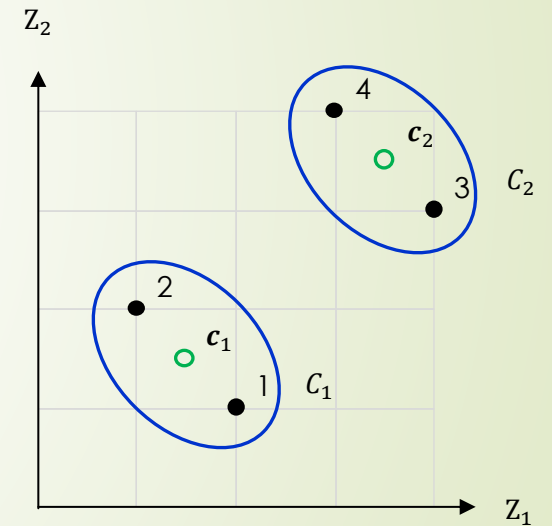  2. The **variations** of the points in the same cluster are **due to noise**.
- How to recover $c_j$ given $C_j$?

$$\min_{c_j} \sum_{p \in C_j} \text{dist}(p, c_j)^2$$

- For Euclidean distance, the solution is the **c**_____

$$c_j = \frac{1}{|C_j|} \sum_{p \in C_j} p$$
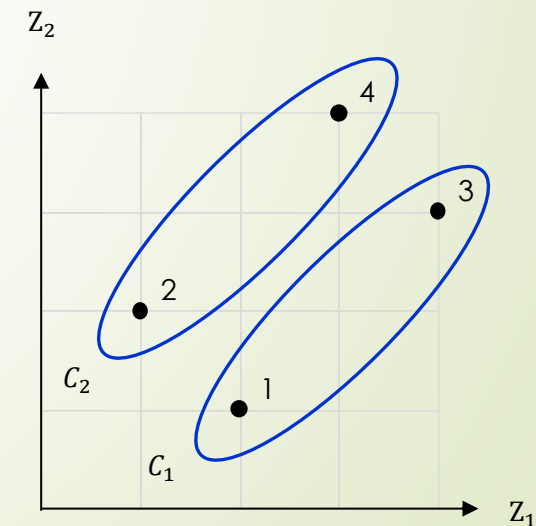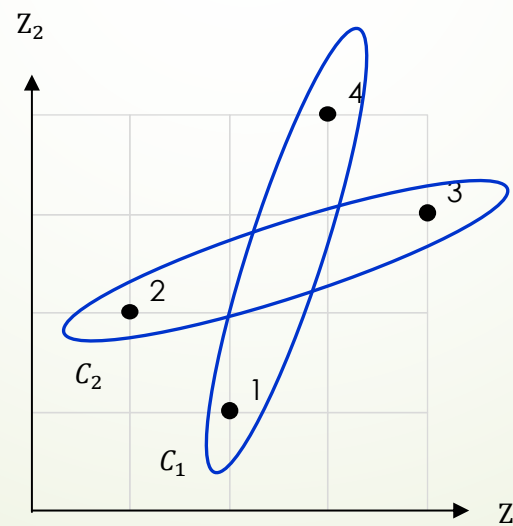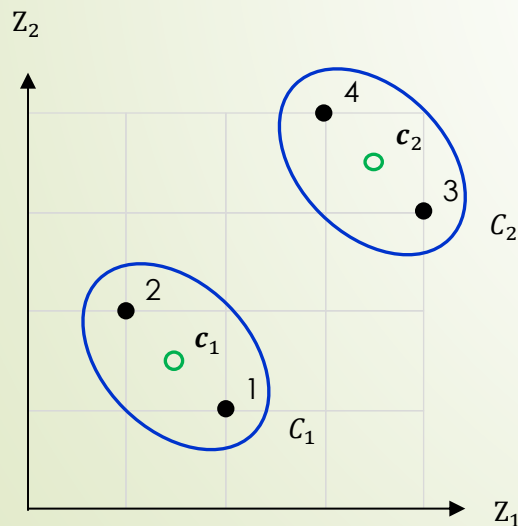
- How to find $C_j$'s?

# Centroid-based method
## Mathematical criteria

- Given the number $k$ of clusters, solve

$$\min_{\{C_j\}_{j=1}^{k}} \sum_{j=1}^{k} \min_{c_j} \sum_{p \in C_j} \mathrm{dist}(p, c_j)^2$$

- Example: Left/Middle/Right is the optimal clustering solution.

# Centroid-based method
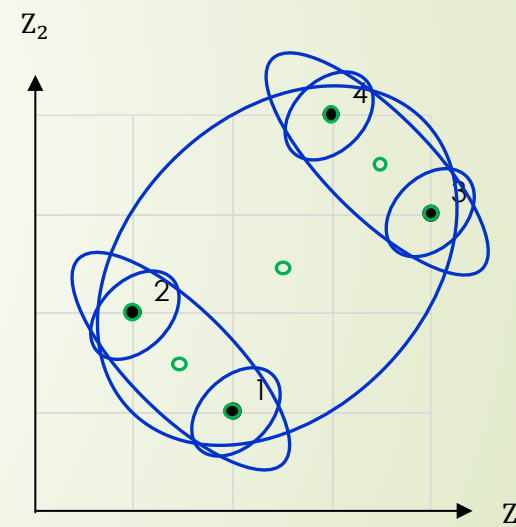## Challenges

- Given the number $k$ of clusters, solve

$$\min_{\{C_j\}_{j=1}^{k}} \sum_{j=1}^{k} \min_{c_j} \sum_{p \in C_j} \text{dist}(p, c_j)^2$$

- What if we further minimize over $k$?

  - $k = \_\_$, $c_j = \_\_$, $C_j = \_\_$ (good? Why or why not?)

  - Not the right objective to find $k$.

  - Remedy? Assume $k$ is given for now.

- Another Issue: Minimization is _____.

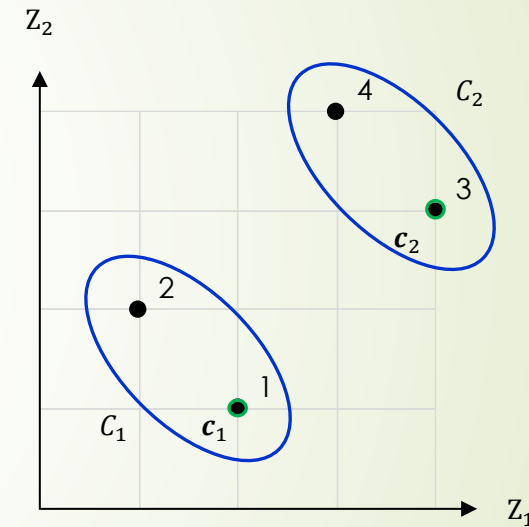  - Best bound: $O(n^{d(k+1)} \log n)$, exponential in $k$ and the dimension $d$.

# $k$-means clustering
## Greedy algorithm

1. Select $k$ tuples randomly as cluster centers initially

2. Calculate the clusters given the cluster centers
   for each $p \in D$
         assign $p$ to $C_j$ where $j$ minimizes $\text{dist}(p, c_j)$

3. Calculate the cluster centers given the clusters
   for each $j$ from 1 to $k$
   $$c_j \leftarrow \frac{1}{|C_j|}\sum_{p \in C_j} p$$
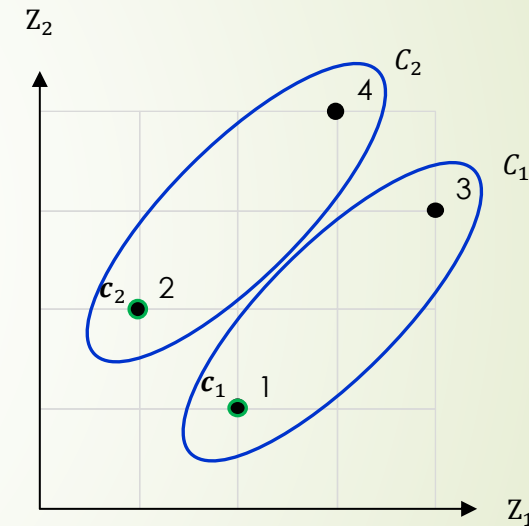
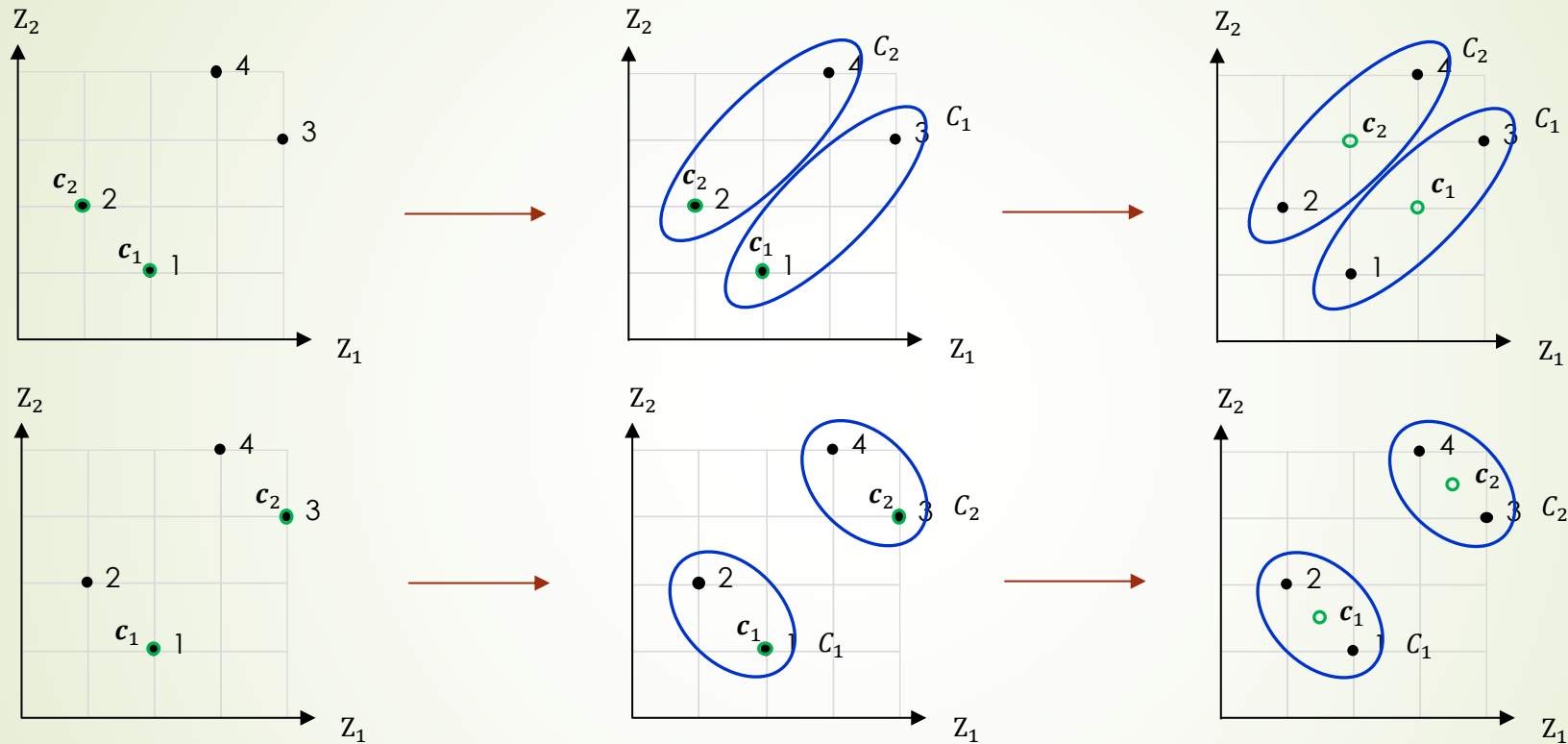4. Repeat 2 to 3 until no change in clusters.

# Complexity

- $O(\underline{\quad})$ where $t$ is the number of iterations.

- Efficient when $k, t \ll n$.

- Does the algorithm always converge to an optimal solution?

# Run again

1. Select $k$ tuples randomly as cluster centers initially

2. Calculate the clusters given the cluster centers
   for each $\boldsymbol{p} \in D$
   
   assign $\boldsymbol{p}$ to $C_j$ where $j$ minimizes $\text{dist}(\boldsymbol{p}, \boldsymbol{c}_j)$

3. Calculate the cluster centers given the clusters
   for each $j$ from 1 to $k$
   
   $$\boldsymbol{c}_j \leftarrow \frac{1}{|C_j|}\sum_{\boldsymbol{p} \in C_j} \boldsymbol{p}$$

4. Repeat 2 to 3 until no change in clusters.
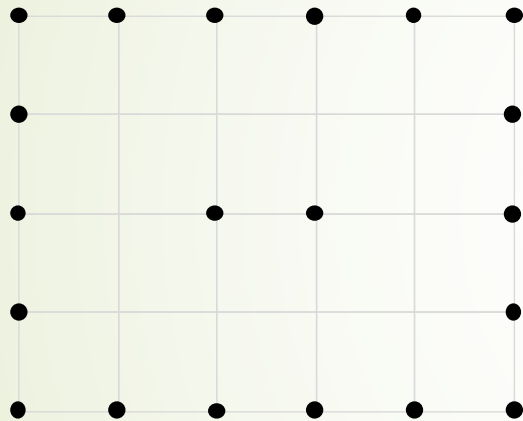
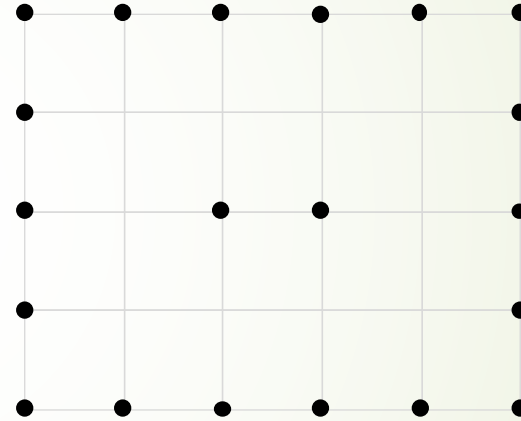# Can fail to converge to the optimum



- ➡ What went wrong? _____
- ➡ What is the chance of failure? _____

# Limitation of centroid-based methods

Desired

Centroid-based

- Fails because the two clusters have the same **c_____**.
- Fail more generally when the cluster shape is **non-c_____/non-s_____**.
- How to resolve?

# References

- 10.1 Cluster Analysis
- 10.2 Partitioning Methods