



香港城市大學
City University of Hong Kong

專業 創新 胸懷全球
Professional • Creative
For The World

Cluster Analysis: Different Evaluation Metrics

CS5483 Data Warehousing and Data Mining

Different aspects for evaluation

- **Tendency:** Do the objects actually form clusters? Statistical significance?
- **Quantity:** What is the number of clusters?
- **Quality:** Are the clusters
 - Correct if ground truth is provided? **E**_____ **measure**.
 - “Clear” if ground truth is not provided? **I**_____ **measure**.

Centroid-based methods

- The **within cluster sum of squares error** (WSS) of the clustering solution as a function of k

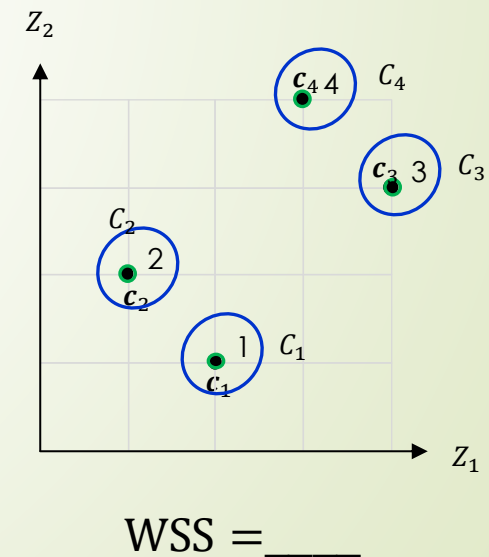
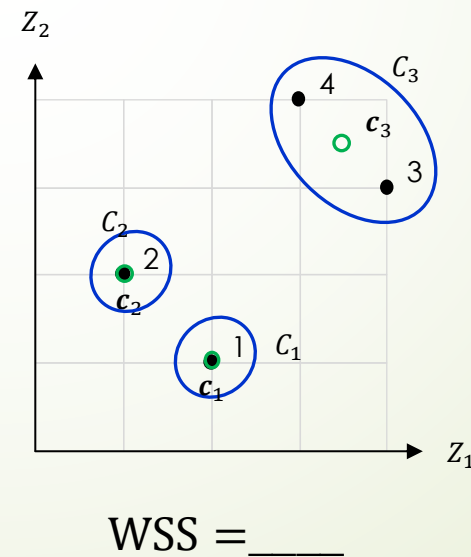
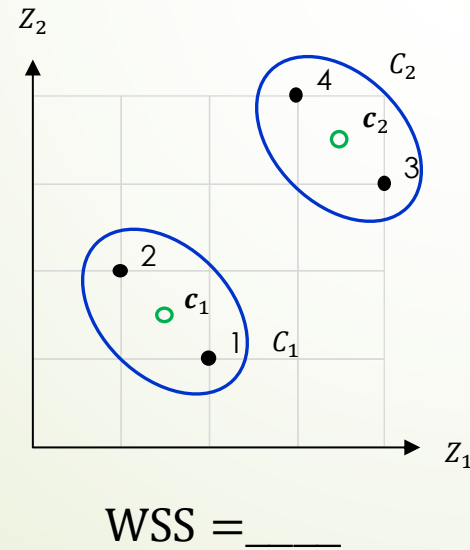
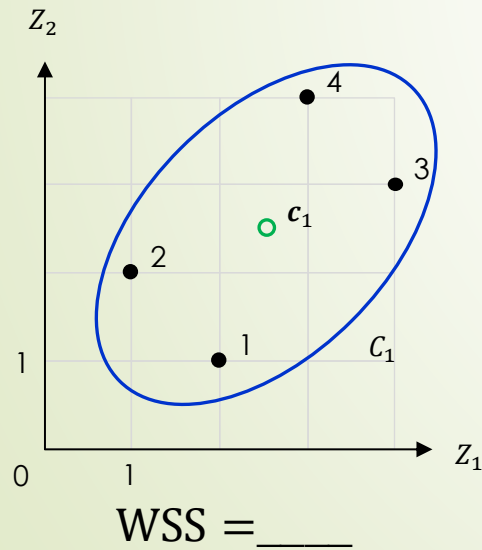
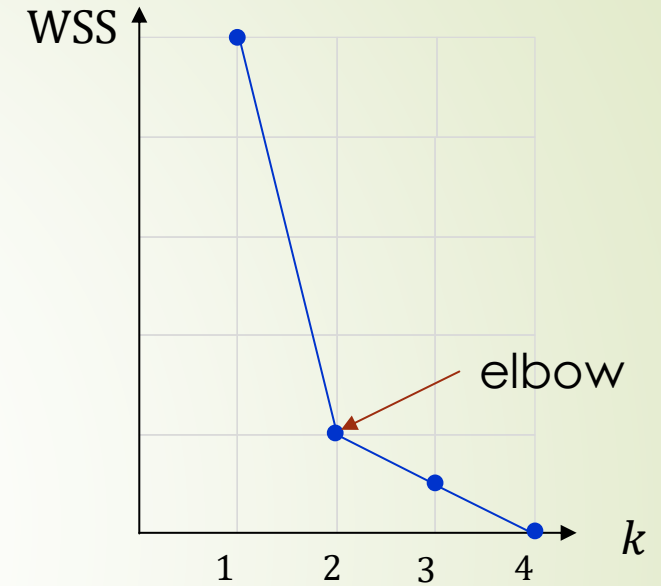
$$\text{WSS}(k) := \sum_{j=1}^k \sum_{p \in C_j} \text{dist}(\mathbf{p}, \mathbf{c}_j)^2,$$

also denoted as $\text{var}(k)$ in [Han11].

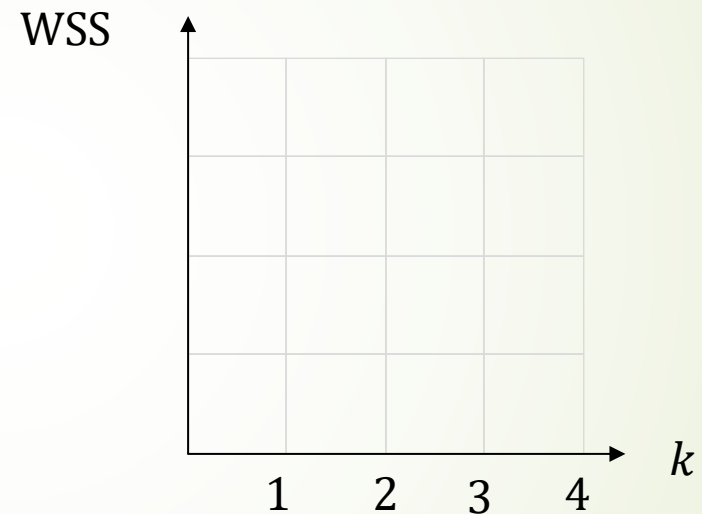
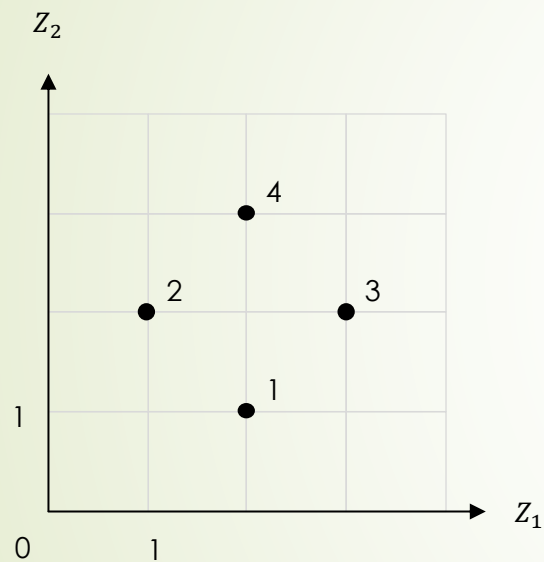
- WSS is an intrinsic/extrinsic quality measure.
 - WSS is small/large if the clusters are clear, i.e., sample points deviate less from centroids.
 - Should we choose k by minimizing WSS? Yes/No because
-

Elbow method

- Choose k that gives large m decrease in WSS from $k - 1$.
- Any problem of overfitting?



Limitation of the elbow method



- The elbow method returns $k =$ _____.
- The elbow method can fail. It is only a **h**_____.
- Other intrinsic measures and analysis?

Silhouette coefficient

- For any point $\mathbf{p} \in D$ in the i -th cluster C_i :

$$s(\mathbf{p}) := \begin{cases} \text{undefined,} & k = 1 \\ 0, & |C_i| = 1 \\ \frac{b(\mathbf{p}) - a(\mathbf{p})}{\max\{a(\mathbf{p}), b(\mathbf{p})\}}, & |C_i| > 1 \end{cases}$$

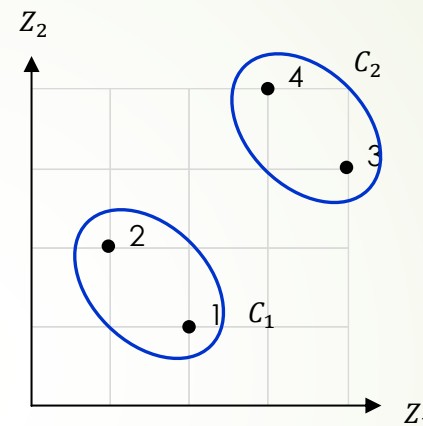
- Mean i-cluster distance:**

$$a(\mathbf{p}) := \frac{1}{|C_i| - 1} \sum_{q \in C_i: q \neq \mathbf{p}} \text{dist}(\mathbf{p}, \mathbf{q})$$

- Mean n-cluster distance:**

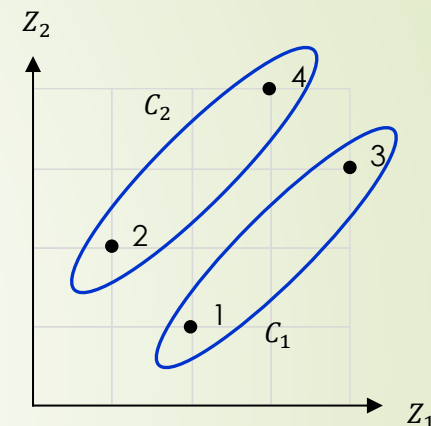
$$b(\mathbf{p}) := \min_{j: \mathbf{p} \notin C_j} \frac{1}{|C_j|} \sum_{q \in C_j} \text{dist}(\mathbf{p}, \mathbf{q})$$

- Quality of C_i : average $s(\mathbf{p})$ over $\mathbf{p} \in C_i$.
- Overall quality: average $s(\mathbf{p})$ over $\mathbf{p} \in D$.



$$\begin{aligned} a(\mathbf{p}) &= \sqrt{2} \\ b(\mathbf{p}) &= \sqrt{2} + \sqrt{5/2} \\ s(\mathbf{p}) &= \frac{\sqrt{5}}{2 + \sqrt{5}} \approx 0.53 \end{aligned}$$

Quality of C_1 : 0.53
 Quality of C_2 : 0.53
 Overall quality: 0.53



$$\begin{aligned} a(\mathbf{p}) &= 2\sqrt{2} \\ b(\mathbf{p}) &= 1/\sqrt{2} + \sqrt{5/2} \\ s(\mathbf{p}) &= \frac{\sqrt{5} - 3}{4} \approx -0.19 \end{aligned}$$

Quality of C_1 : -0.19
 Quality of C_2 : -0.19
 Overall quality: -0.19

Silhouette coefficient

- For any point $\mathbf{p} \in D$ in the i -th cluster C_i :

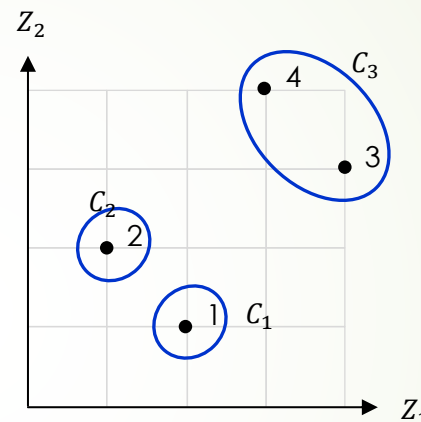
$$s(\mathbf{p}) := \begin{cases} \text{undefined,} & k = 1 \\ 0, & |C_i| = 1 \end{cases} \text{ Why?} \\ \frac{b(\mathbf{p}) - a(\mathbf{p})}{\max\{a(\mathbf{p}), b(\mathbf{p})\}}, |C_i| > 1$$

- $s(\mathbf{p}) \in [-1, 1]$: why?

- $b(\mathbf{p}) > a(\mathbf{p})$: $s(\mathbf{p}) =$ _____

- $b(\mathbf{p}) < a(\mathbf{p})$: $s(\mathbf{p}) =$ _____

- Sometimes we need a more detailed analysis by a [silhouette plot](#).
- The method can fail on non-s_____ clusters.



$$s(\mathbf{p}_1) = s(\mathbf{p}_2) = \underline{\hspace{2cm}}$$

$$a(\mathbf{p}_3) = a(\mathbf{p}_4) = \underline{\hspace{2cm}}$$

$$b(\mathbf{p}_3) = b(\mathbf{p}_4) = \underline{\hspace{2cm}}$$

$$s(\mathbf{p}_3) = s(\mathbf{p}_4) = \underline{\hspace{2cm}}$$

Quality of C_1 : _____

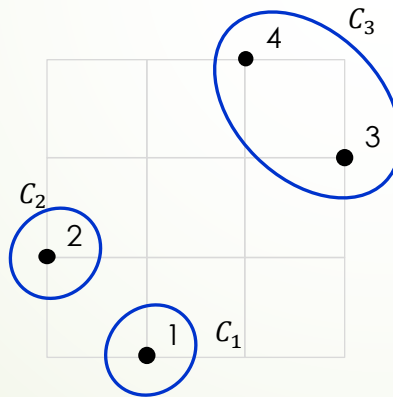
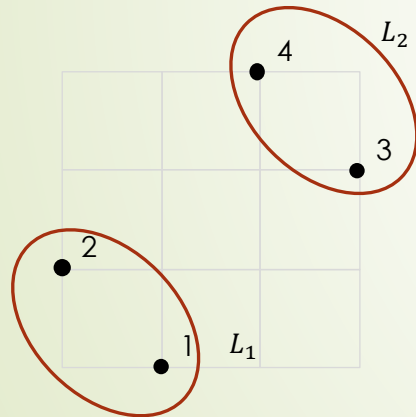
Quality of C_2 : _____

Quality of C_3 : _____

Quality of the clustering: _____

Extrinsic cluster quality measures

- Setting: Ground truth is available
 - $L(\mathbf{p})$: \mathbf{C} _____ (class) of \mathbf{p} . (Ground truth.)
 - $C(\mathbf{p})$: Cluster index of \mathbf{p} .
- How to compare the clustering solution to the ground truth?



$$L(\mathbf{p}_1) = L(\mathbf{p}_2) = 1$$

$$L(\mathbf{p}_3) = L(\mathbf{p}_4) = 2$$

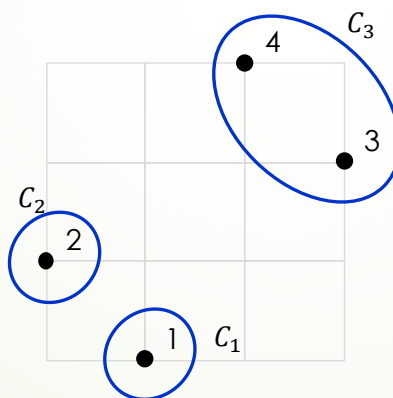
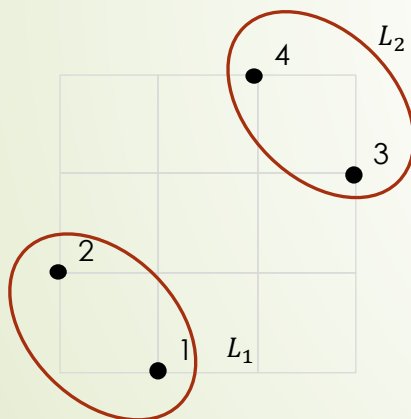
$$C(\mathbf{p}_1) = 1$$

$$C(\mathbf{p}_2) = 2$$

$$C(\mathbf{p}_3) = C(\mathbf{p}_4) = 3$$

Extrinsic cluster quality measures

- For two points $p, q \in D$ to be clustered correctly:
 - p, q should belong to the same cluster if they are in the same category.
 - p, q should NOT belong to the same cluster if they are NOT in the same category.
 - Success indicator: $\text{correctness}(p, q) = \mathbb{1}(L(p) = L(q) \leftrightarrow C(p) = C(q))$



$\text{correctness}(p_i, p_j)$

$i \backslash j$	1	2	3	4
1				
2				
3				
4				

Accuracy: _____

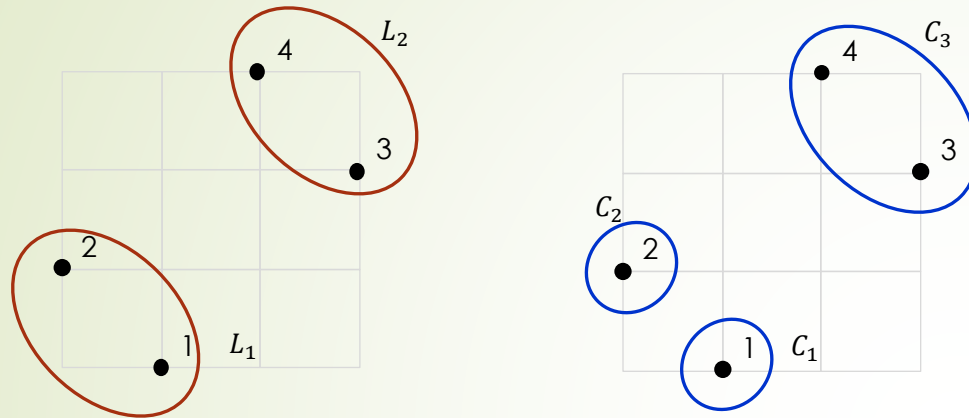
- Is accuracy a good performance metric? Yes/No because _____

B-Cubed precision and recall

- For $p, q \in D$ (allowing $p = q$), compute the confusion matrix

	$C(p) = C(q)$	$C(p) \neq C(q)$
$L(p) = L(q)$	TP	FN
$L(p) \neq L(q)$	FP	TN

- Sample p and q from D with replacement.
 - If p, q are in the same cluster, what is the chance they are in the same category?
 - *recall/precision := —
 - If p, q are in the same category, what is the chance they are in the same cluster?
 - *recall/precision := —



precision = _____

recall = _____

	$C(p) = C(q)$	$C(p) \neq C(q)$
$L(p) = L(q)$	TP=6: (1,1),(2,2),(3,3),(4,4),(3,4),(4,3)	FN=____
$L(p) \neq L(q)$	FP=____	TN=8: (1,3),(3,1),(2,3),(3,2),(1,4),(4,1),(2,4),(4,1)

Sanity check: total number of counts should be _____.

Alternative formulae

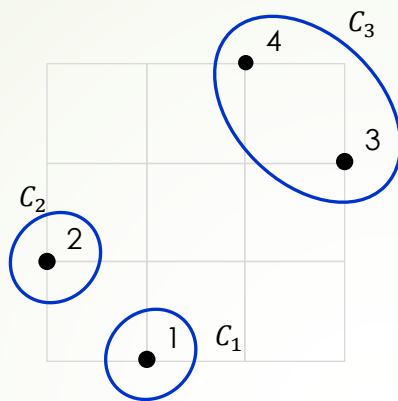
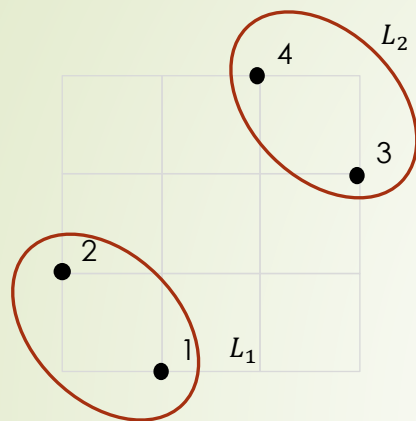
- B-Cubed precision is the average over $\mathbf{p} \in D$ of

$$\text{precision}(\mathbf{p}) := \frac{|\{\mathbf{q} \in D \mid \mathcal{C}(\mathbf{p}) = \mathcal{C}(\mathbf{q}), L(\mathbf{p}) = L(\mathbf{q})\}|}{|\{\mathbf{q} \in D \mid \mathcal{C}(\mathbf{p}) = \mathcal{C}(\mathbf{q})\}|}$$

- B-Cubed recall is the average over $\mathbf{p} \in D$ of

$$\text{recall}(\mathbf{p}) := \frac{|\{\mathbf{q} \in D \mid \mathcal{C}(\mathbf{p}) = \mathcal{C}(\mathbf{q}), L(\mathbf{p}) = L(\mathbf{q})\}|}{|\{\mathbf{q} \in D \mid L(\mathbf{p}) = L(\mathbf{q})\}|}$$

- Advantage: Breaks down the quality measure and its computation to individual point.



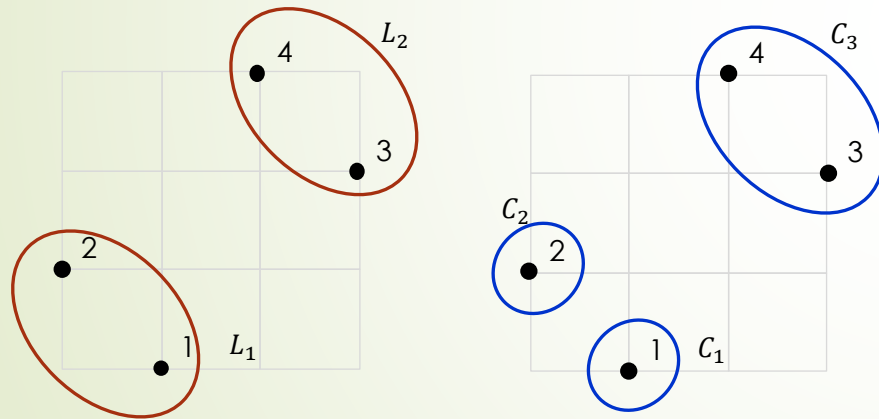
precision = _____

recall = _____

i	1	2	3	4
precision(p_i)			100%	100%
recall(p_i)			100%	100%

Classes to clusters evaluation (WEKA)

- Match class to cluster labels to maximize correctly classified tuples.
- Use the “classification” error rate as the performance measure.



Number of instances

Class assignment

$L \backslash C$	1	2	3
1	1	1	0
2	0	0	2

Cluster \leftarrow Class

1 \leftarrow _____

2 \leftarrow _____

3 \leftarrow _____

Accuracy = _____%

- (Optional) There are other extrinsic measures such as the adjusted rand index.

References

- 10.4.2 Ordering Points to Identify the Clustering Structure
- 10.6 Evaluation of Clustering
- Error in Han11:
 - Should remove the constraint $i \neq j$ from (10.29) and (10.30).
- Supplementary readings:
 - Amigó, Enrique, et al. "A comparison of extrinsic clustering evaluation metrics based on formal constraints." Information retrieval 12.4 (2009): 461-486.
 - <https://scikit-learn.org/stable/modules/clustering.html#clustering-evaluation>