# Comparing the Performance of Biomedical Clustering Methods

Barbara Ikica

Faculty of Mathematics and Physics, UL & Institute of Mathematics, Physics and Mechanics

9 January 2018

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# ClustEval - Integrative Clustering Evaluation Framework

C. Wiwie, R. Röttger, J. Baumbach, Comparing the performance of biomedical clustering methods, *Nat. Methods* **12** (2015), 1033–1038.

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# ClustEval - Integrative Clustering Evaluation Framework

C. Wiwie, R. Röttger, J. Baumbach, Comparing the performance of biomedical clustering methods, *Nat. Methods* **12** (2015), 1033–1038.

- a clustering analysis platform *ClustEval*, available at
  https://clusteval.sdu.dk/1/mains

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# ClustEval - Integrative Clustering Evaluation Framework

C. Wiwie, R. Röttger, J. Baumbach, Comparing the performance of biomedical clustering methods, *Nat. Methods* **12** (2015), 1033–1038.

- a clustering analysis platform *ClustEval*, available at https://clusteval.sdu.dk/1/mains
- 18 different clustering methods, 24 datasets, 18 common cluster validity indices

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# ClustEval - Integrative Clustering Evaluation Framework

C. Wiwie, R. Röttger, J. Baumbach, Comparing the performance of biomedical clustering methods, *Nat. Methods* **12** (2015), 1033–1038.

- a clustering analysis platform *ClustEval*, available at
  https://clusteval.sdu.dk/1/mains
- 18 different clustering methods, 24 datasets, 18 common cluster validity indices
- robustness analysis (density reduction, noise addition)

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# ClustEval - Integrative Clustering Evaluation Framework

C. Wiwie, R. Röttger, J. Baumbach, Comparing the performance of biomedical clustering methods, *Nat. Methods* **12** (2015), 1033–1038.

- a clustering analysis platform *ClustEval*, available at https://clusteval.sdu.dk/1/mains
- 18 different clustering methods, 24 datasets, 18 common cluster validity indices
- robustness analysis (density reduction, noise addition)
- parameter optimisation (Naïve Grid Optimisation, Adaptive Grid Optimisation)

# Clustering in Biomedicine

Applicable to:

- cancer subtyping on the basis of gene expression levels,

# Clustering in Biomedicine

Applicable to:

- cancer subtyping on the basis of gene expression levels,
- protein homology detection from amino acid sequences of structures,

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Clustering in Biomedicine

Applicable to:

- cancer subtyping on the basis of gene expression levels,
- protein homology detection from amino acid sequences of structures,
- the identification of protein complexes using protein-protein interactions

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# Key Problems

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# Key Problems

1. **Tool picking**

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval

Motivation

Key Problems

Clustering
Methods

Partitioning, hierarchical

Density-based, graph-based

Datasets

Clustering Validity
Indices

Internal, External

Correlation

Guidelines

# Key Problems

1. **Tool picking** (partitioning ($k$-means), hierarchical, density-based and graph-based approaches)

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval

Motivation

Key Problems

Clustering
Methods

Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices

Internal, External
Correlation

Guidelines

# Key Problems

1. **Tool picking** (partitioning ($k$-means), hierarchical, density-based and graph-based approaches)
2. **Parameter optimisation**

# Key Problems

1. **Tool picking** (partitioning ($k$-means), hierarchical, density-based and graph-based approaches)
2. **Parameter optimisation** (parameters influence the number and size of resulting clusters)

# Key Problems

1. **Tool picking** (partitioning ($k$-means), hierarchical, density-based and graph-based approaches)
2. **Parameter optimisation** (parameters influence the number and size of resulting clusters)
3. **Quality measures/Cluster validity indices**

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval

Motivation

Key Problems

Clustering
Methods

Partitioning, hierarchical

Density-based, graph-based

Datasets

Clustering Validity
Indices

Internal, External

Correlation

Guidelines

# Key Problems

1. **Tool picking** (partitioning ($k$-means), hierarchical, density-based and graph-based approaches)
2. **Parameter optimisation** (parameters influence the number and size of resulting clusters)
3. **Quality measures/Cluster validity indices** (internal, external)

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Key Problems

1. **Tool picking** (partitioning ($k$-means), hierarchical, density-based and graph-based approaches)
2. **Parameter optimisation** (parameters influence the number and size of resulting clusters)
3. **Quality measures/Cluster validity indices** (internal, external)
4. **Standardised evaluation**

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval

Motivation

Key Problems

Clustering
Methods

Partitioning, hierarchical

Density-based, graph-based

Datasets

Clustering Validity
Indices

Internal, External

Correlation

Guidelines

# Key Problems

1. **Tool picking** (partitioning ($k$-means), hierarchical, density-based and graph-based approaches)
2. **Parameter optimisation** (parameters influence the number and size of resulting clusters)
3. **Quality measures/Cluster validity indices** (internal, external)
4. **Standardised evaluation** (poor comparability of clustering results)

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# Clustering Methods

# Clustering Methods

## Partitioning

The objects are assigned to clusters and iteratively change clusters based on their dissimilarity in order to optimise a given target function.

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Clustering Methods

## Partitioning

The objects are assigned to clusters and iteratively change clusters based on their dissimilarity in order to optimise a given target function.

*Fanny, $k$-means, Partitioning Around Medoids*

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# Clustering Methods

## Partitioning

The objects are assigned to clusters and iteratively change clusters based on their dissimilarity in order to optimise a given target function.

*Fanny, $k$-means, Partitioning Around Medoids*

## Hierarchical

The dataset is transformed into a tree-like structure where the leaves represent the objects and the inner nodes the hierarchical relationship between them.

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Clustering Methods

## Partitioning

The objects are assigned to clusters and iteratively change clusters based on their dissimilarity in order to optimise a given target function.

*Fanny, $k$-means, Partitioning Around Medoids*

## Hierarchical

The dataset is transformed into a tree-like structure where the leaves represent the objects and the inner nodes the hierarchical relationship between them.

*Hierarchical Clustering, Spectral Clustering*

# Clustering Methods

## Density-based

Identifying regions with a locally similar object density.

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# Clustering Methods

## Density-based

Identifying regions with a locally similar object density.

*Clusterdp, DBSCAN*

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Clustering Methods

## Density-based

Identifying regions with a locally similar object density.

*Clusterdp, DBSCAN*

## Graph-based

The input is considered as a graph with the objects being the nodes connected by weighted or unweighted edges. The clustering problem is solved by solving an analogous graph theoretical problem (e.g. clique finding, simulating random walks).

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Clustering Methods

## Density-based

Identifying regions with a locally similar object density.

*Clusterdp, DBSCAN*

## Graph-based

The input is considered as a graph with the objects being the nodes connected by weighted or unweighted edges. The clustering problem is solved by solving an analogous graph theoretical problem (e.g. clique finding, simulating random walks).

*Affinity Propagation, clusterONE, Markov Clustering, MCODE, Transitivity Clustering*

# Datasets

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Datasets

- Gene expression levels,

# Datasets

- Gene expression levels,
- Protein-protein interaction,

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Datasets

- Gene expression levels,
- Protein-protein interaction,
- Protein-sequence similarity,

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Datasets

- Gene expression levels,
- Protein-protein interaction,
- Protein-sequence similarity,
- Social network,

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Datasets

- Gene expression levels,
- Protein-protein interaction,
- Protein-sequence similarity,
- Social network,
- Synthetic (easy/medium/hard),

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Datasets

- Gene expression levels,
- Protein-protein interaction,
- Protein-sequence similarity,
- Social network,
- Synthetic (easy/medium/hard),
- Language-processing datasets used for word-sense disambiguation

# Clustering Validity Indices

## Internal Indices

The clustering is judged on the basis of certain intrinsic statistical properties of the clustering itself.

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
  Motivation
  Key Problems

Clustering
Methods
  Partitioning, hierarchical
  Density-based, graph-based

Datasets

Clustering Validity
Indices
  Internal, External
  Correlation

Guidelines

# Clustering Validity Indices

## Internal Indices

The clustering is judged on the basis of certain intrinsic statistical properties of the clustering itself.

*Davies-Bouldin Index, Dunn Index, Silhouette Value*

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Clustering Validity Indices

## Internal Indices

The clustering is judged on the basis of certain intrinsic statistical properties of the clustering itself.

*Davies-Bouldin Index, Dunn Index, Silhouette Value*

## External Indices

The clustering is compared to a user-given gold-standard clustering (using a pairwise/mapping approach).

# Clustering Validity Indices

## Internal Indices

The clustering is judged on the basis of certain intrinsic statistical properties of the clustering itself.

*Davies-Bouldin Index, Dunn Index, Silhouette Value*

## External Indices

The clustering is compared to a user-given gold-standard clustering (using a pairwise/mapping approach).

*$F_\beta$ score, False Discovery rate, False Positive Rate, Fowles-Mallows Index, Jaccard Index, Rand Index, Sensitivity (Recall), Specificity, V-measure*

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# Clustering Validity Indices

**Davies-Bouldin Index**

$$DB = \frac{1}{n} \cdot \sum_{c_i \in C} \max_{c_i \neq c_j \in C} \left( \frac{\sigma_i + \sigma_j}{\|\overline{c_i} - \overline{c_j}\|} \right),$$

$$\sigma_i = \sqrt{\frac{1}{|c_i|} \sum_{x_i \in c_i} \|x_i - \overline{c_i}\|}$$

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Clustering Validity Indices

**Davies-Bouldin Index**

$$DB = \frac{1}{n} \cdot \sum_{c_i \in C} \max_{c_i \neq c_j \in C} \left( \frac{\sigma_i + \sigma_j}{\|\overline{c_i} - \overline{c_j}\|} \right),$$

$$\sigma_i = \sqrt{\frac{1}{|c_i|} \sum_{x_i \in c_i} \|x_i - \overline{c_i}\|}$$

**Dunn Index**

$$D = \frac{\min\limits_{c_i \neq c_j \in C} \left( \min\limits_{x_i \in c_i, x_j \in c_j} d(x_i, x_j) \right)}{\max\limits_{c_k \in C} \left( \max\limits_{x_i, x_j \in c_k} d(x_i, x_j) \right)}$$

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval

Motivation
Key Problems

Clustering
Methods

Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices

Internal, External

Correlation

Guidelines

# Clustering Validity Indices

**Silhouette Value**

$c(x_i)$ – the cluster of object $x_i$

$o(x_i)$ – the closest cluster to $x_i$

$$S = \frac{1}{n} \sum_i sv_i$$

$$sv_i = \frac{1}{n} \frac{b(x_i) - a(x_i)}{\max\{a(x_i), b(x_i)\}}$$

$$a(x_i) = \frac{1}{|c(x_i)|} \sum_{x_j \in c(x_i)} d(x_i, x_j)$$

$$b(x_i) = \frac{1}{|o(x_i)|} \sum_{x_j \in o(x_i)} d(x_i, x_j)$$

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# Clustering Validity Indices

**Silhouette Value**

$c(x_i)$ – the cluster of object $x_i$
$o(x_i)$ – the closest cluster to $x_i$

$$S = \frac{1}{n} \sum_i sv_i$$

## Correlation among internal and external indices

The silhouette value correlated best (0.71) with the F1
score, F2 score, FM index, Jaccard index and V-measure.

# Guidelines (working with a new *biomedical* set without a gold standard)

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

# Guidelines (working with a new *biomedical* set without a gold standard)

1. Use *Transitivity Clustering*, *Hierarchical Clustering* or *Partitioning Around Medoids*.

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Guidelines (working with a new *biomedical* set without a gold standard)

1. Use *Transitivity Clustering*, *Hierarchical Clustering* or *Partitioning Around Medoids*.

2. Compute the silhouette values for clustering results using a broad range of parameter-set variations.

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Guidelines (working with a new *biomedical* set without a gold standard)

1. Use *Transitivity Clustering*, *Hierarchical Clustering* or *Partitioning Around Medoids*.
2. Compute the silhouette values for clustering results using a broad range of parameter-set variations.
3. Pick the result for the parameter set yielding the highest silhouette value.

Comparing the
Performance of
Biomedical
Clustering
Methods

Barbara Ikica

ClustEval
Motivation
Key Problems

Clustering
Methods
Partitioning, hierarchical
Density-based, graph-based

Datasets

Clustering Validity
Indices
Internal, External
Correlation

Guidelines

# Guidelines (working with a new *biomedical* set without a gold standard)

1. Use *Transitivity Clustering*, *Hierarchical Clustering* or *Partitioning Around Medoids*.
2. Compute the silhouette values for clustering results using a broad range of parameter-set variations.
3. Pick the result for the parameter set yielding the highest silhouette value.

**Remark.** The silhouette value is a particularly poor measure for entangled and highly overlapping datasets.