Clusterability in Model Selection

Johannes Kiesel

Bauhaus-Universität Weimar

28th May, 2014

Cluster Analysis: Motivation



Given data (a set of comparable *entities* or *objects*) Find a categorization of it

Cluster Analysis: Motivation



Given data (a set of comparable *entities* or *objects*) Find a categorization of it (without labels)

Cluster Analysis: Motivation



Given data (a set of comparable *entities* or *objects*) Find a categorization of it (without labels)

Cluster Analysis: In the Beginning was the Data



Data





































Cluster Analysis: Model Evaluation



Cluster Analysis: Overview





- Task: calculate a score for a model
- Has to be comparable at least among similar models (same number of objects)



 A clusterable model (high score) has a dominant structure of mutually separated parts that are cohesive groups of objects.

Clusterability I: Salient Clustering

Idea Model selection by cluster evaluation ("one-step")

- Cluster the model with different algorithms and/or parameter settings
- Evaluate all clusterings
- Choose best combination of model & clustering



Clusterability I: Dunn Index



Minimum spanning tree Dunn index (Dunn MST)

1/● Largest edge length in the minimum spanning tree of the cluster

⇔ Smallest dissimilarity of objects from different clusters

Optimum clustering is feasibly computable (no other clustering algorithm necessary)

Clusterability I: Salient Clustering

+

+ Needs no additional clusterability index
+ Evaluation indices are better understood

- Most evaluation indices require local optimization
 Not all evaluation indices
- can compare clusterings of different models



Clusterability II: Statistical Tests on Structure

Idea Use a statistical test for unstructured models

 Null hypothesis: model generated from a model distribution that generates non-clusterable models (e.g., uniform distribution)



- Calculate a test statistic with known distribution under the null hypothesis
- Use the probability that a similar large value occurs under the null hypothesis for the clusterability assessment

Clusterability II: Hopkins and Skellam Statistic



+ **X** • **X**⁰



Clusterability II: Hopkins and Skellam Statistic



$\psi_{nn}(\mathbf{x})$ Dissimilarity of \mathbf{x} to its nearest neighbor

Clusterability II: Hopkins and Skellam Statistic



Compare distribution of original objects (\mathbf{x}) and *r* uniformly sampled \mathbf{x}^0 (null hypothesis)

$$\mathbf{H}_{r} = \frac{\sum_{i=1}^{r} (\psi_{nn}(\mathbf{x}_{i}^{0}))^{m}}{\sum_{i=1}^{r} (\psi_{nn}(\mathbf{x}_{i}^{0}))^{m} + \sum_{i=1}^{r} (\psi_{nn}(\mathbf{x}_{\pi(i)}))^{m}}$$

$\psi_{nn}(\mathbf{x})$ Dissimilarity of \mathbf{x} to its nearest neighbor *m* Number of dimensions

Clusterability II: Statistical Tests on Structure

+

+ The distribution under the null hypothesis allows for an interpretation of the score

+ Often requires only a sample



- Adjustment of statistics is not trivial



Clusterability III: Concentration of Dissimilarities

Idea In a clusterable model most object pairs should be either very dissimilar (different clusters) or very similar (same clusters)



Test if relatively few dissimilarities are of average size









[Dash et al. Dimensionality Reduction for Unsupervised Data. 1997]



[Dash et al. Dimensionality Reduction for Unsupervised Data. 1997]



Clusterability III: Concentration of Dissimilarities

+ Very general idea
+ Related to the concept of intrinsic dimensionality

- Not clear when the used heuristic (see right figure) applies
- Lacks the interpretability of statistical tests



Clusterability: Overview

A clusterable model has a dominant structure of mutually separated parts that are cohesive groups of objects.



- Clusterability is related to various other topics in data analysis
 - Evaluation indices (Dunn)
 - Tests on model distributions (Hopkins and Skellam)
 - Methods of unsupervised feature selection (Dash et al.)
 - Estimators of intrinsic dimensionality
 - ▶ ...?

Can the clusterability indices identify clusterable models?

Experiment setup:

10 model distributions of varying intuitive clusterability



1 model from the uniform distribution

- 1 000 models per distribution (results are means)
- 180 2-dimensional objects per model

s = 0 s = 0.1 s = 0.2 s = 0.3



<i>s</i> = 0	<i>s</i> = 0.1	<i>s</i> = 0.2	s = 0.3	symbol



[1] Limited to clusterings with 13 or less clusters[2] Mean of 1 000 applications per model







[1] Limited to clusterings with 13 or less clusters[2] Mean of 1 000 applications per model



Contributions

A clusterable model has a dominant structure of mutually **separated** parts that are **cohesive** groups of objects.

- Clusterability indices can be used for model selection
- The indices differ, among others, with respect to their preference for fine or coarse structure
- If models are (somewhat) meaningful for a dataset, the more clusterable models are assumed to be also the more meaningful
- Clusterability can incorporate ideas from various related topics (especially clustering evaluation)
- Formal properties of clustering evaluation indices can be converted to properties of clusterability indices



- Further formalization of clusterability indices
- Application to large datasets
- Application to high-dimensional problems
- Relation to cluster stability
- Incorporation of additional knowledge (constraint clustering)

Thank you for your attention.