



Martin Hahmann, Markus Dumat, Dirk Habich and Wolfgang Lehner Explorative Multi-View Clustering Using Frequent-Groupings

MultiClust 2012 28.04.2012

Out of many, one:

- integrate multiple solution into a single result
- ensemble clustering
- multi-source clustering

Out of one, many:

- generate different results from the same source
- alternative clustering
- subspace clustering

Our goal:

bring both sides together using frequent groupings

Database Technolog



Core Concept



> Best of both worlds



Frequent Groupings Idea (MultiClust2011)

- starting point: ensemble clustering
- goal: find cluster assignment that agrees with majority of ensemble
- find objects frequently assigned to the same cluster
- \rightarrow parallel to the concept of frequent itemsets
- mapping to ensemble clustering:



Database Technology

Frequent Grouping Graph



Database Technology Group







Scaled Up Scenario



Database Technology Group

Manual processing not feasible



- 1500 objects, 10 clusterings, 69 clusters
- apply existing frequent itemset mining algorithms
- CARPENTER with minsupp= 0.5 \rightarrow 72 frequent groupings

How to create alternative Solutions?

- assumption: every object is assigned to a cluster
- resembles exact-cover problem (np-hard)



- Algorithm X by Donald E. Knuth
- implementation: Dancing Links

570.000 alternative clusterings

9



Database Technology



Grou



Problems

- too many results to work with
- high similarity between alternatives
- smallest frequent groupings define minimum difference

Options for result limitation:

- filter small groupings
- compute similarity matrix
- top-k Dancing Links:

Top-k: 5

							_	_	
fg14	fg35	fg10	fg3	fg7	fg9	fg6	fg17	fg23	

Clusterings: {(fg14, fg 35, fg7), (fg10, fg7), (fg3, fg7)}





Second Approach: Exploring the graph

































© Martin Hahmann | UNIVERSITAT





- number of alternatives is predictable
- less or even the number of root nodes
- additional alternatives via root removal & recursion







> Result Similarity



Similarity measurement via intersection

IV			IV A3			000	
	fg71 ∩ fg48	107	ſ		fg48 ∩ fg71	107	1
	$ fg_{71} \cap fg_{55} $	151			fg48 ∩ fg65	23	
	fg71 ∩ fg47	128			$ fgss \cap fg_{71} $	151	
	fg71 ∩ fg60	0			fgss ∩ fgss	3	
	$ fg_{65} \cap fg_{48} $	23			fg ₄₇ ∩ fg ₇₁	128	
	fg ₆₅ ∩ fg ₅₅	3			fg ₄₇ ∩ fg ₆₅	0	
	fg ₆₅ ∩ fg ₄₇	0	[fg ₆₀ ∩ fg ₇₁	0	
	fg ₆₅ ∩ fg ₆₀	208	[fg ₆₀ ∩ fg ₆₅	208	

(105+208)/620 = 0,58



Future Work





Exploration so far

• top-down traversal from roots \rightarrow alternatives through splitting clusters

Other possibilities

- bottom-up traversal from leaves \rightarrow alternatives through merging clusters
- starting level between root and leaves?
- combine bottom-up and top-down







Further Graph Utilization

- examine branching for dissimilarity of alternatives
- for feedback on ensemble diversity/homogenity
- relation between ensemble and alternatives







- frequent groupings combine alternative and ensemble-clustering
- offer alternative and robust results
- automated recombination and its problems
- user-centered exploration



Questions?

