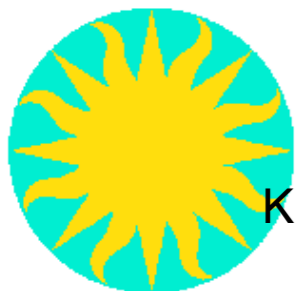


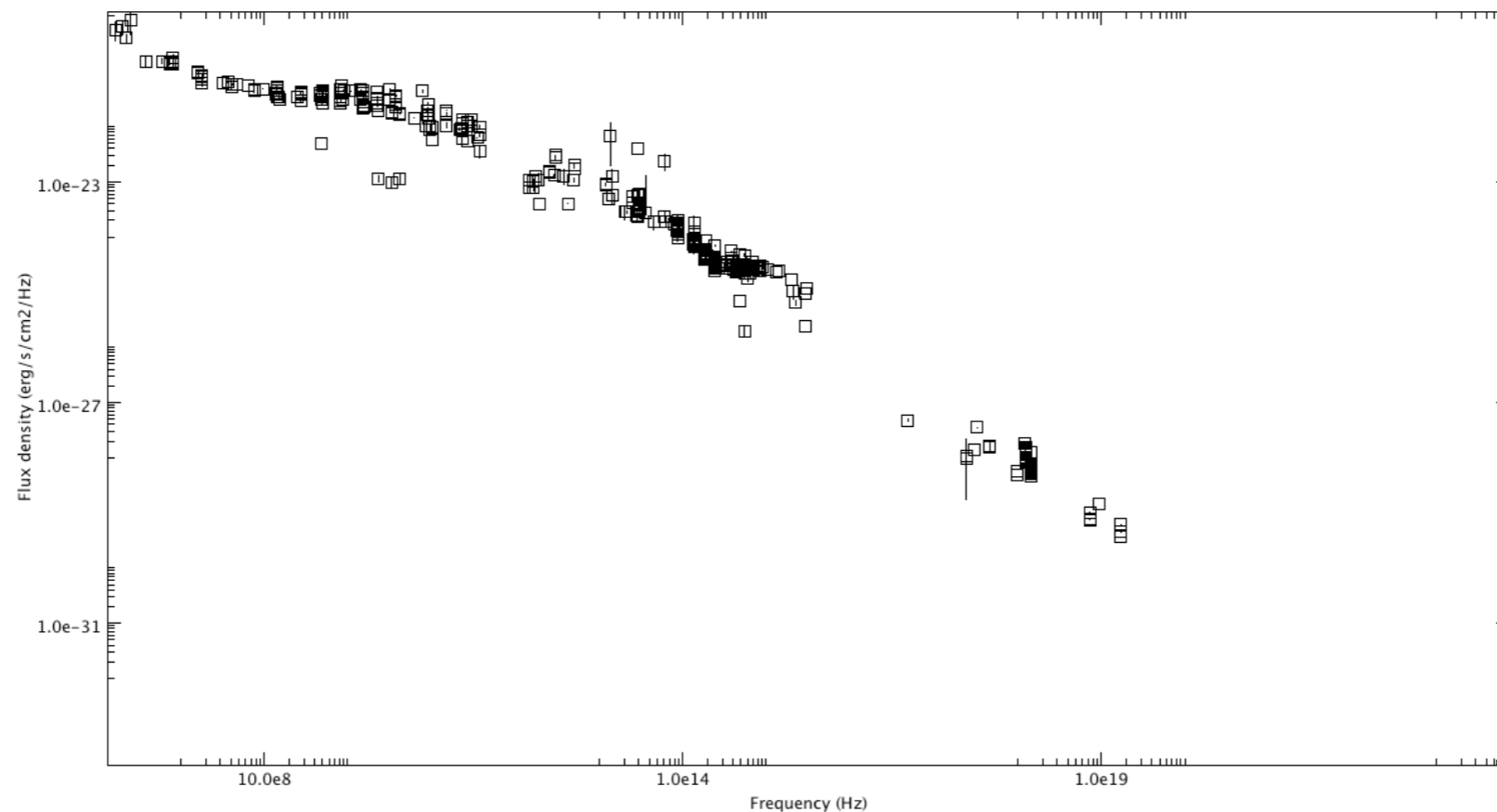
# The exploration of multi-wavelength astronomical datasets: the case of AGNs in the Chandra Source Catalog and unsupervised clustering.

**R. D'Abrusco**



# Motivations

Characterization of the distribution of AGNs in a high dimensionality parameter space obtained by combining multi-wavelength data and study of their X-rays properties.



**The primary purpose of this study is to obtain a possible census of AGN behavior in the 13-dimensional features space of X-UV-optical-IR-Radio photometry, and pick up outliers and constrain their nature.**

# The approach

Unsupervised clustering in a high dimensional features space

**AND**

use additional information (labels) to identify interesting cluster(ing)s

**IN ORDER TO**

derive new high dimensional correlations and/or expand known correlations to other bands, and/or spot unusual behaviors (**outliers**).

# Statistical issues

## Few points for a large space

> 10 dimensional features space  
 $10^2 \sim 10^3$  sources



Very low specific density

## Upper limits & Clustering

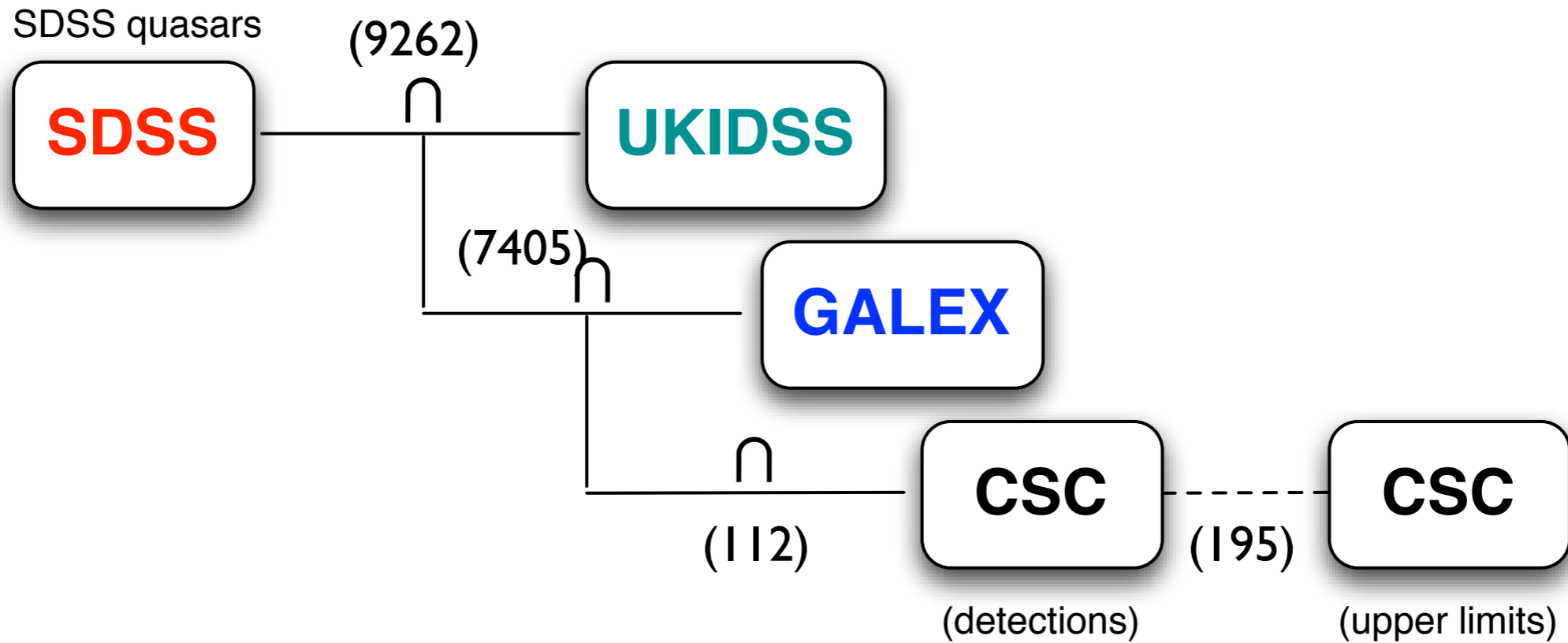
Inclusion of sources with no detections but reliable upper limits.

## Outliers vs Clusters

Most clustering methods tends to prefer the selection of either well populated homogeneous clusters or to sparse clusters/singletons (outliers).

# The datasets

## 1) “Large area surveys” sample

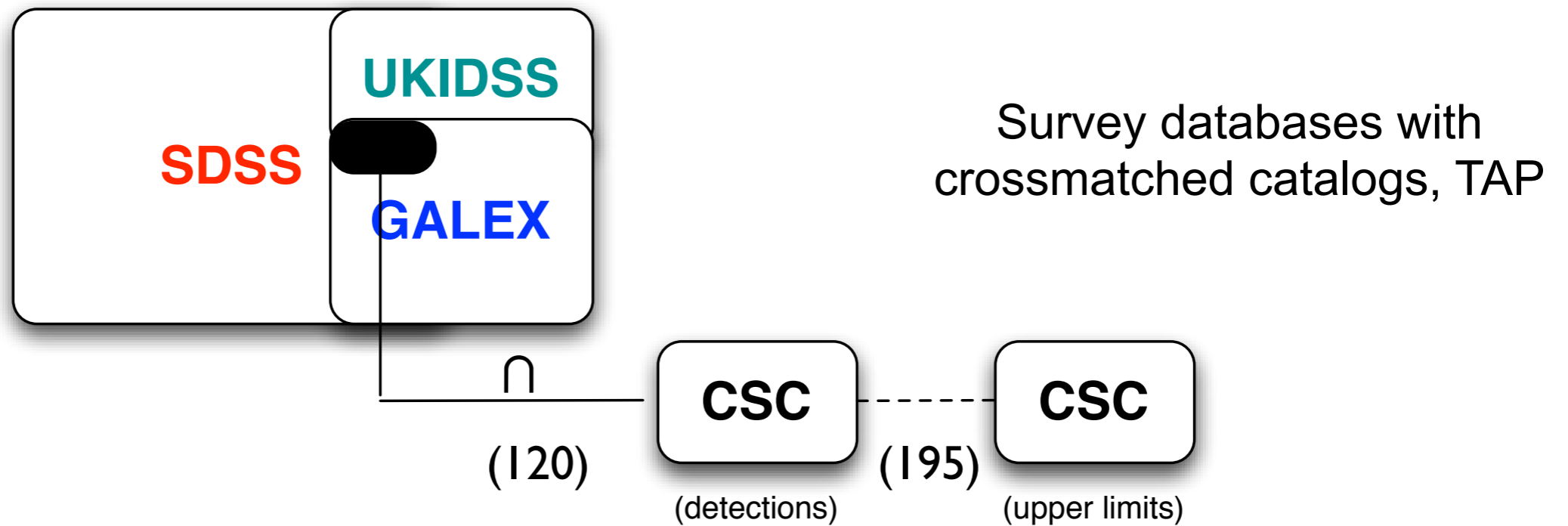


## 2) Chandra COSMOS X-ray survey

## 3) SWIRE

# The datasets

## 1) “Large area surveys” sample



## 2) Chandra COSMOS X-ray survey

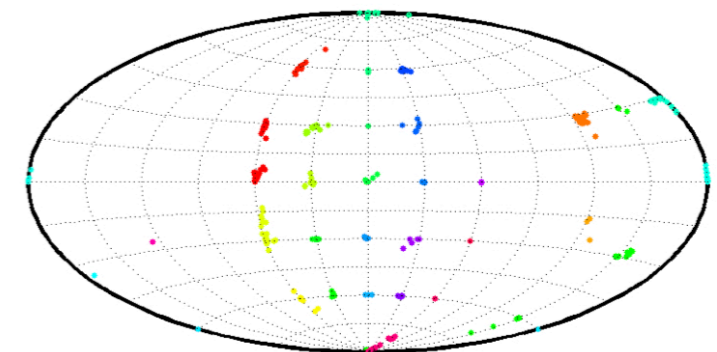
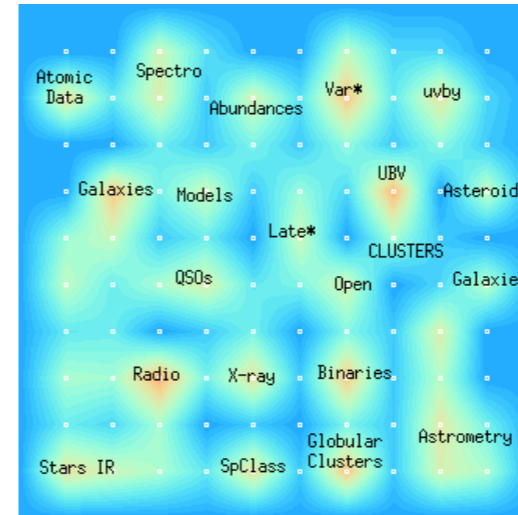
## 3) SWIRE

# Multiple techniques

Hierarchical Clustering (K-means)

Self-Organizing Maps (SOM)

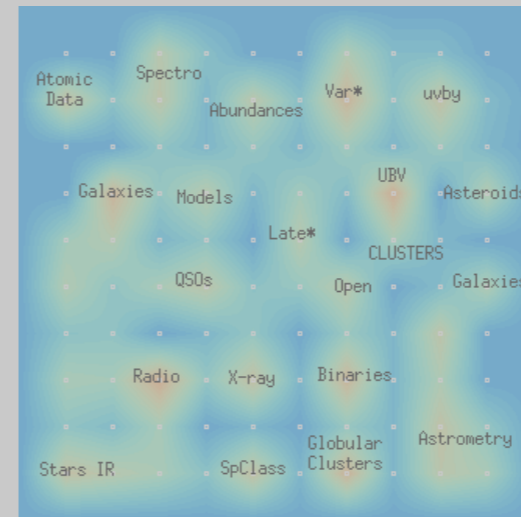
Principal Probabilistic  
Surfaces



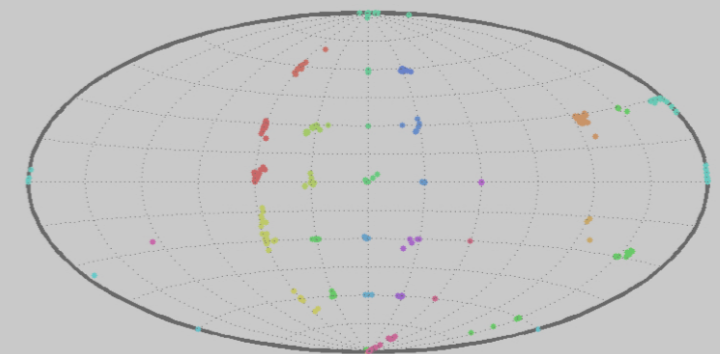
# Multiple techniques

## Hierarchical Clustering (K-means)

Self-Organizing Maps (SOM)



Principal Probabilistic  
Surfaces





# Dendrograms

## Representation of hierarchical structure - HC tree.

HC does not require a fixed number of clusters and produces all possible clustering, given a *measure of dissimilarity* (distance). Every generation of clusters maximizes the between-group dissimilarity.

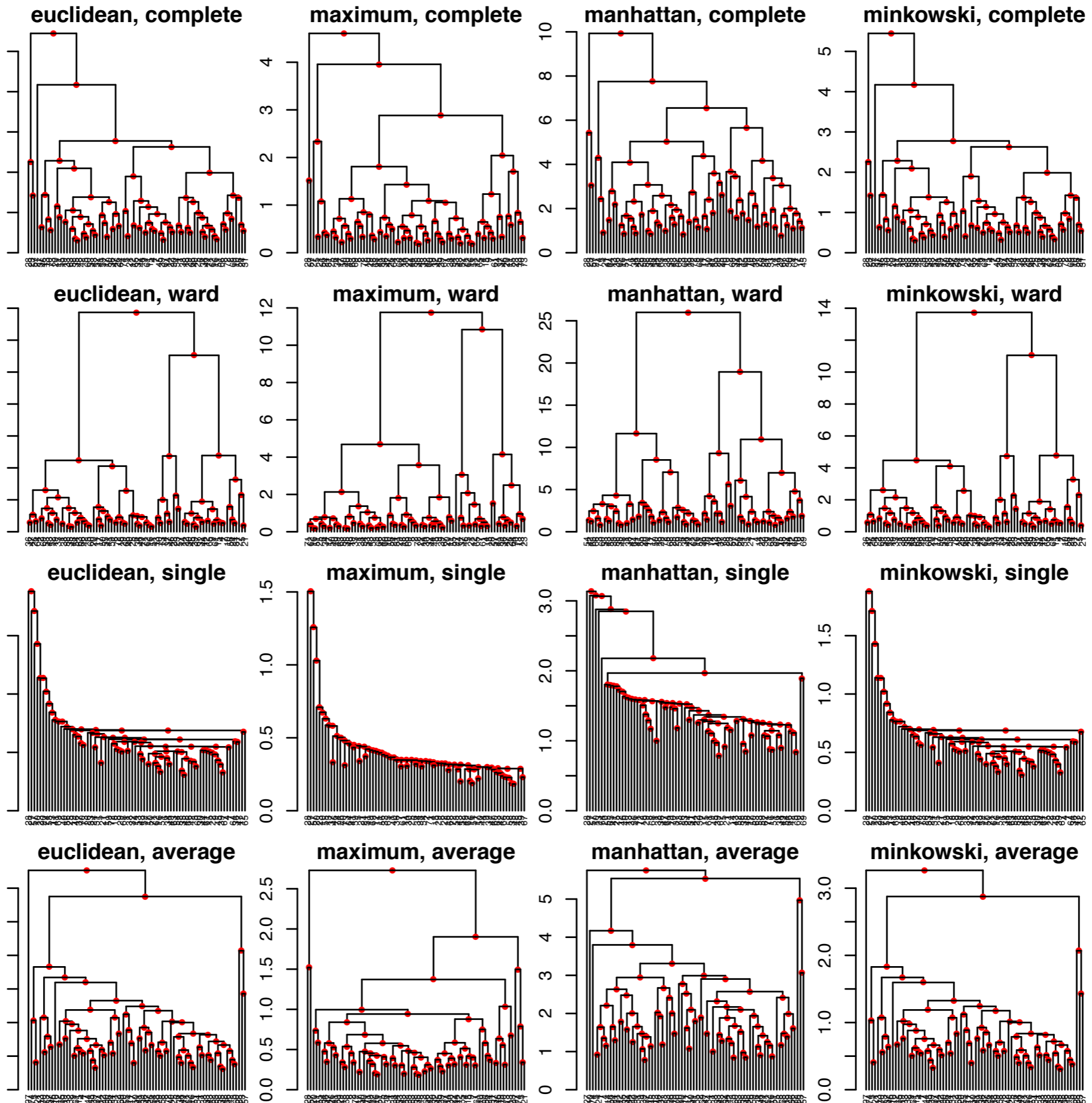
dissimilarity  $\equiv$  (metric, linkage strategy)

### Metrics:

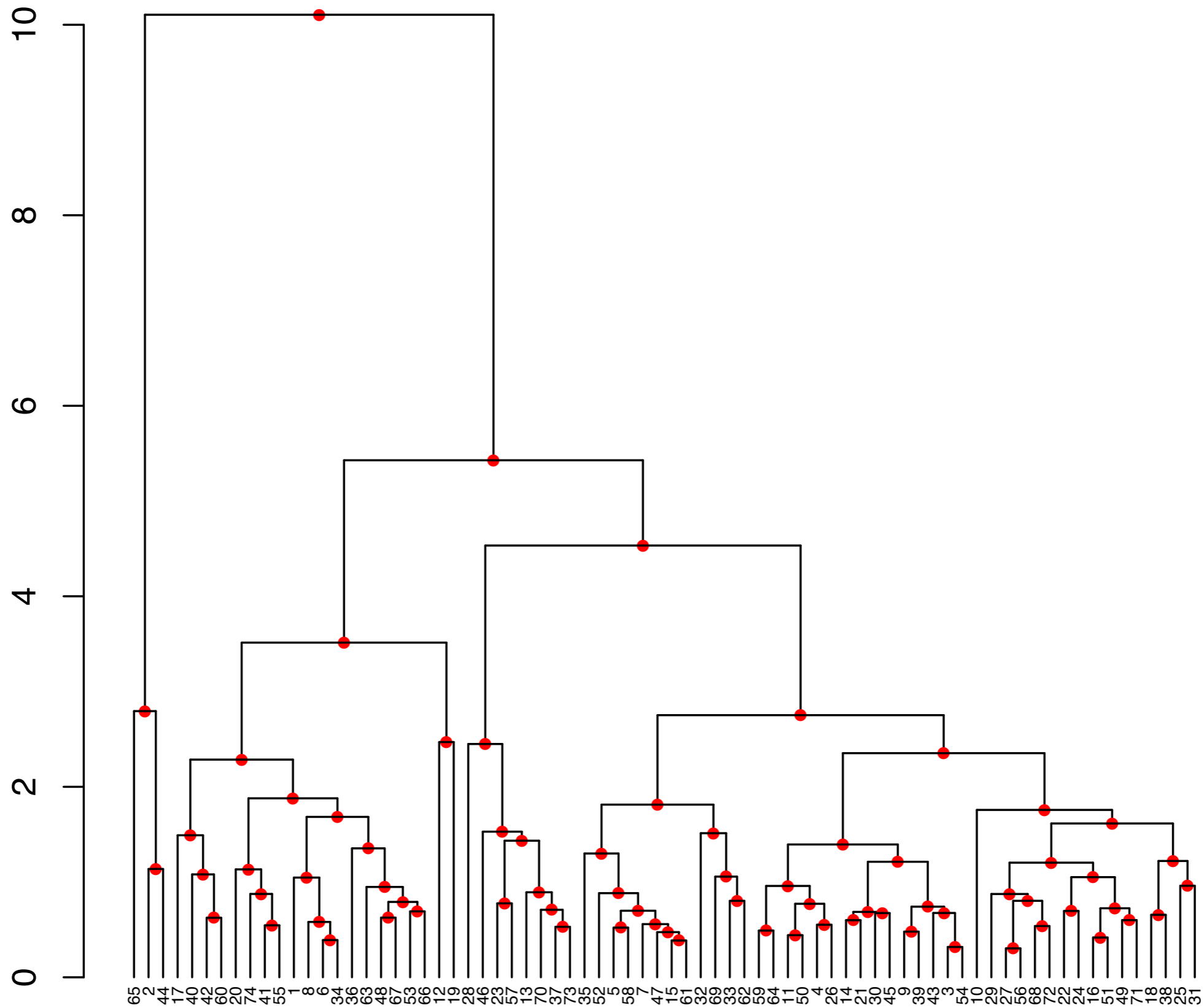
Euclidean, Manhattan,  
Mahalanobis, maximum, etc.

### Linkage strategies:

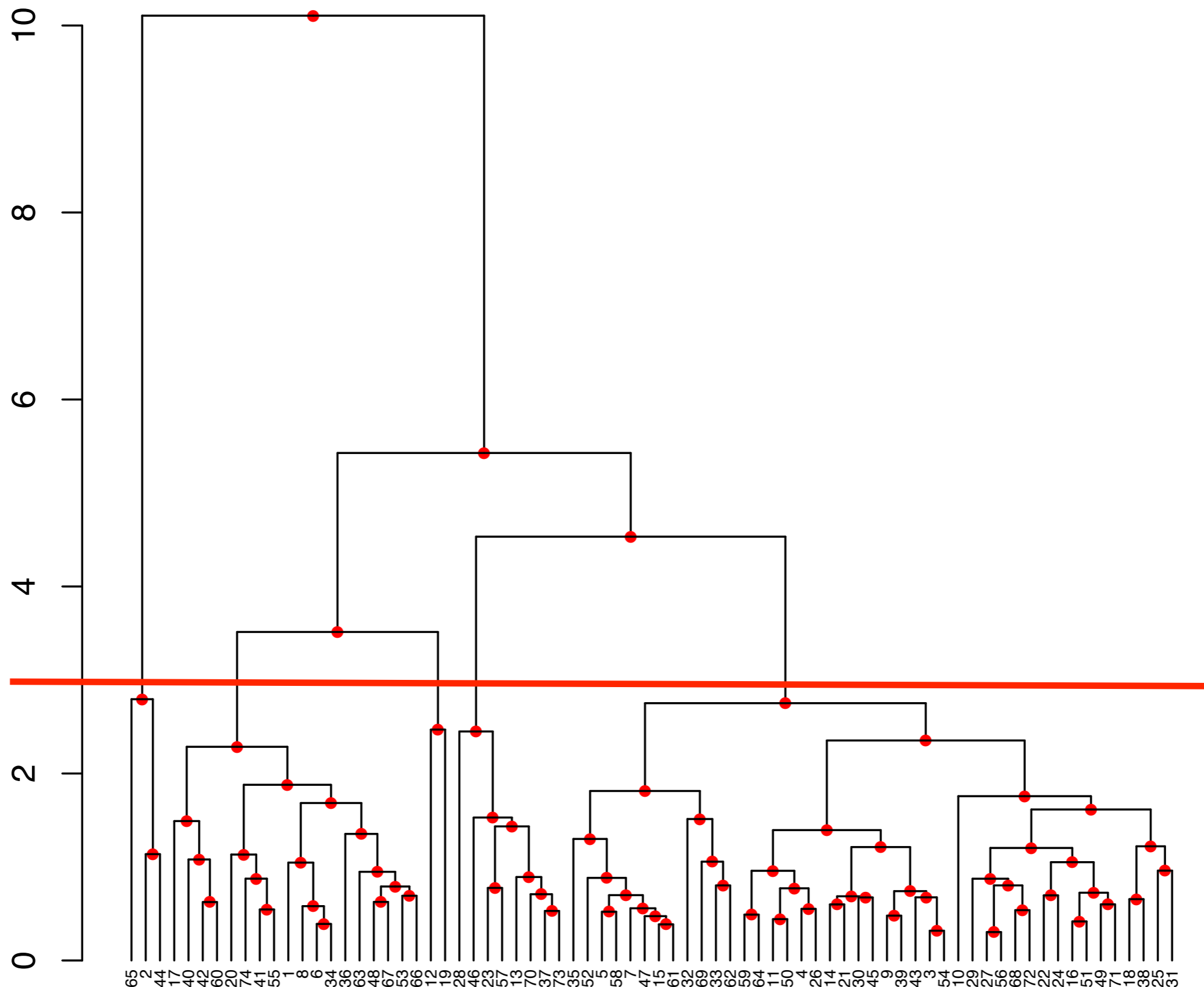
complete, single,  
average, etc.



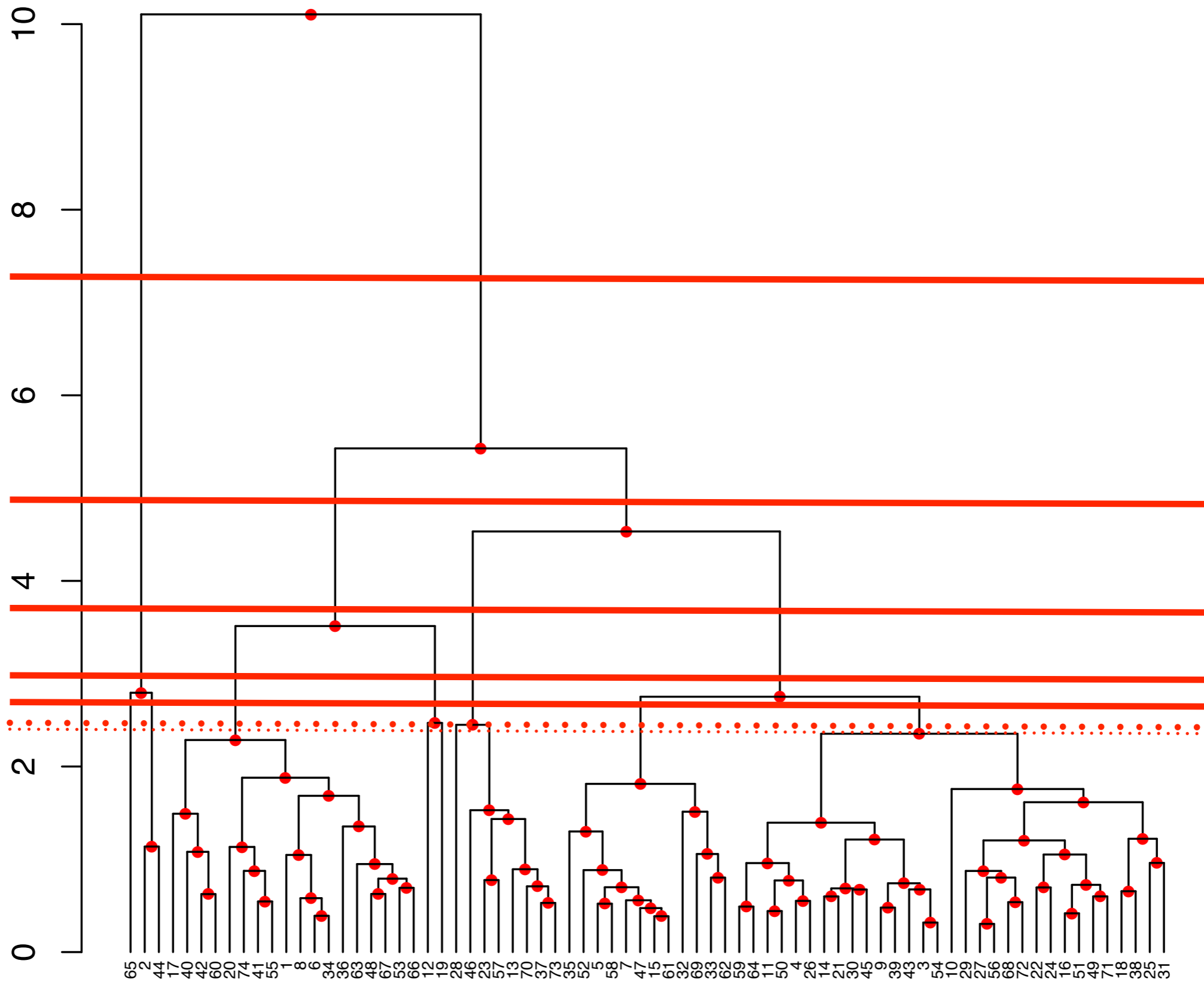
**Distance: euclidean; linking strategy: complete**



# Distance: euclidean; linking strategy: complete



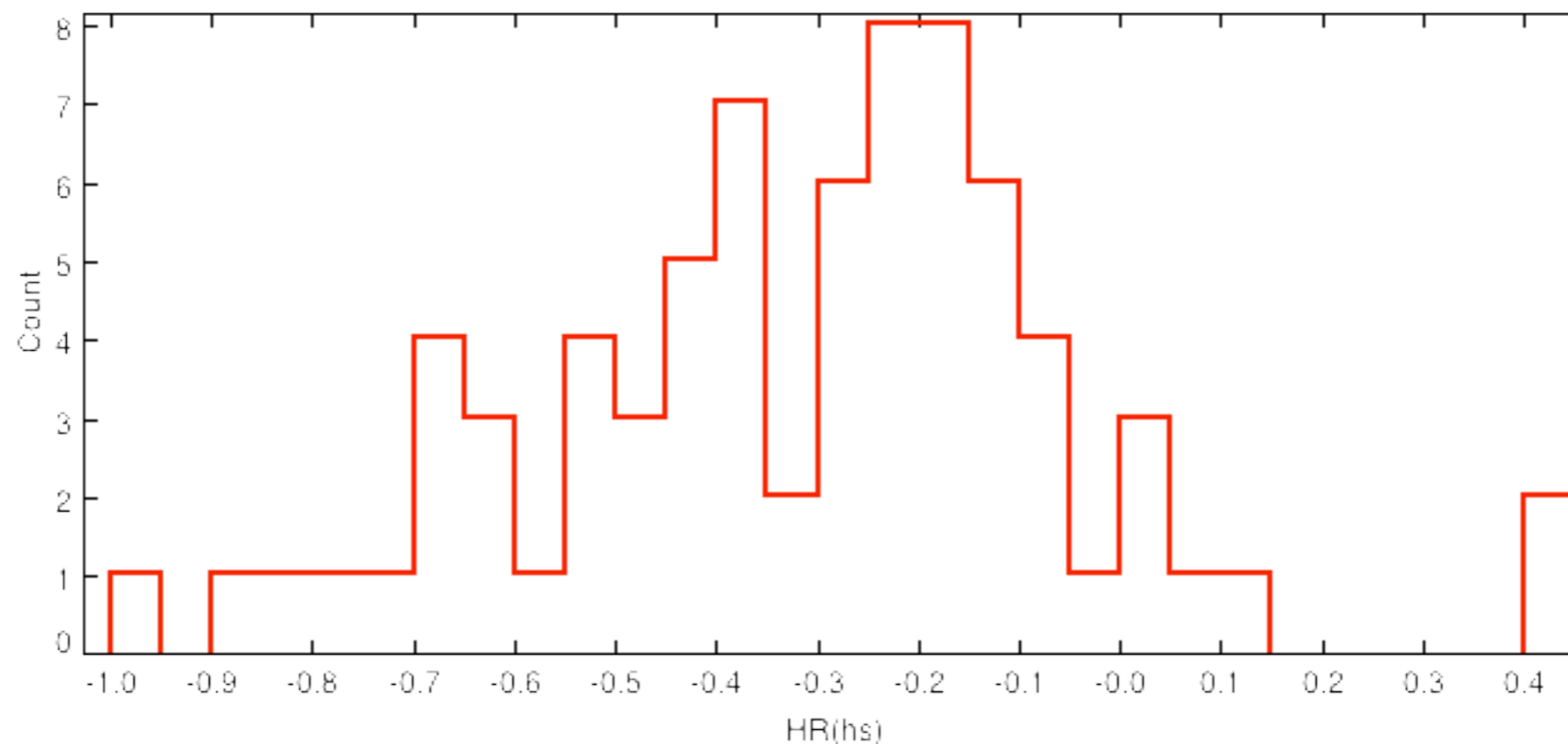
# Distance: euclidean; linking strategy: complete



# Labels to pick clusters

Other measured quantities, called “labels” (continuous or discrete) are separated into bins (if continuous) and used to pick those cluster(ing)s which seem interesting.

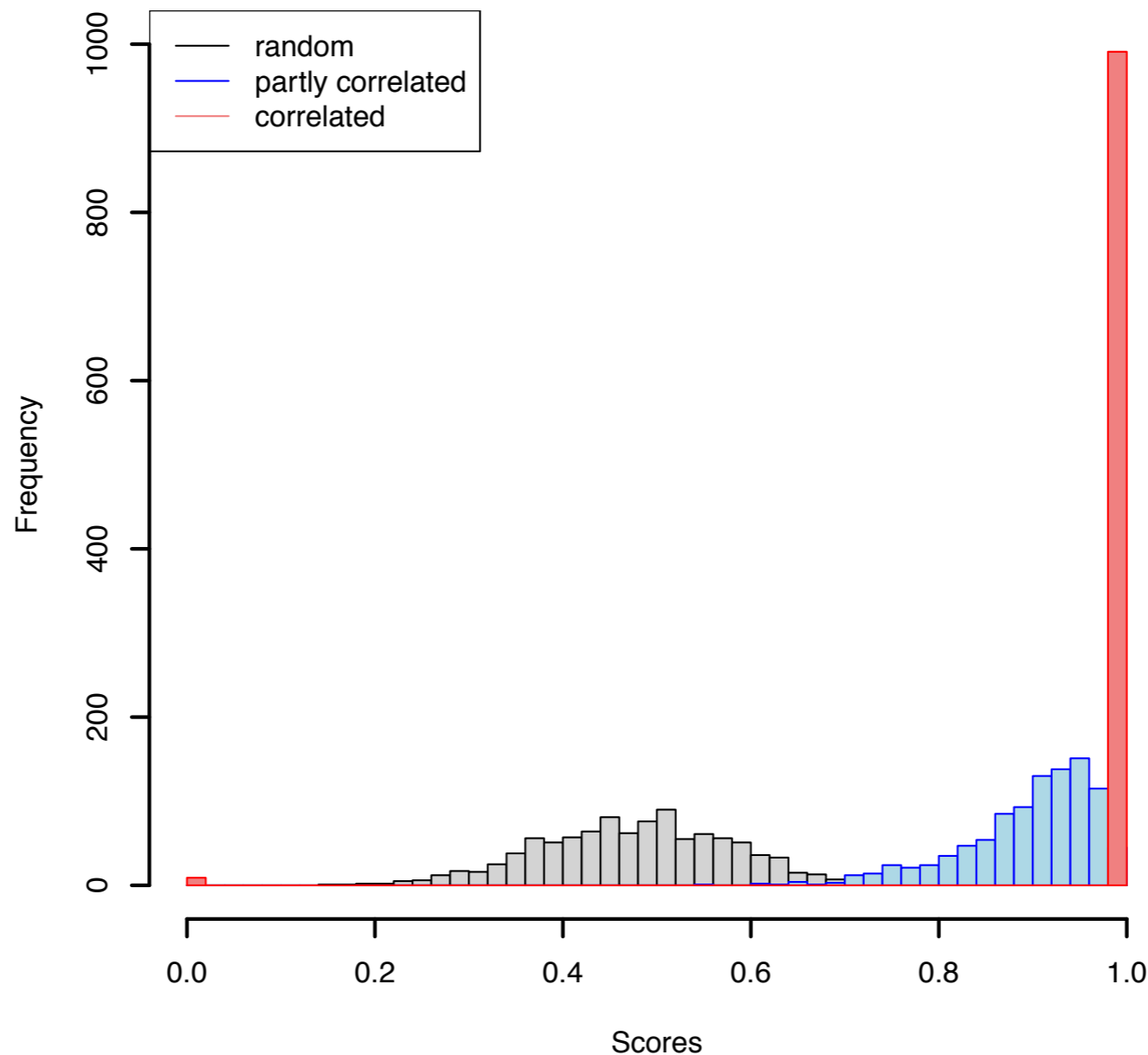
**$L_x$ , HR,  $\Gamma$ ,  $n_h$ , X-ray/optical ratio, X-ray time variability, radio morphology (if available)**



# Discriminating clusterings

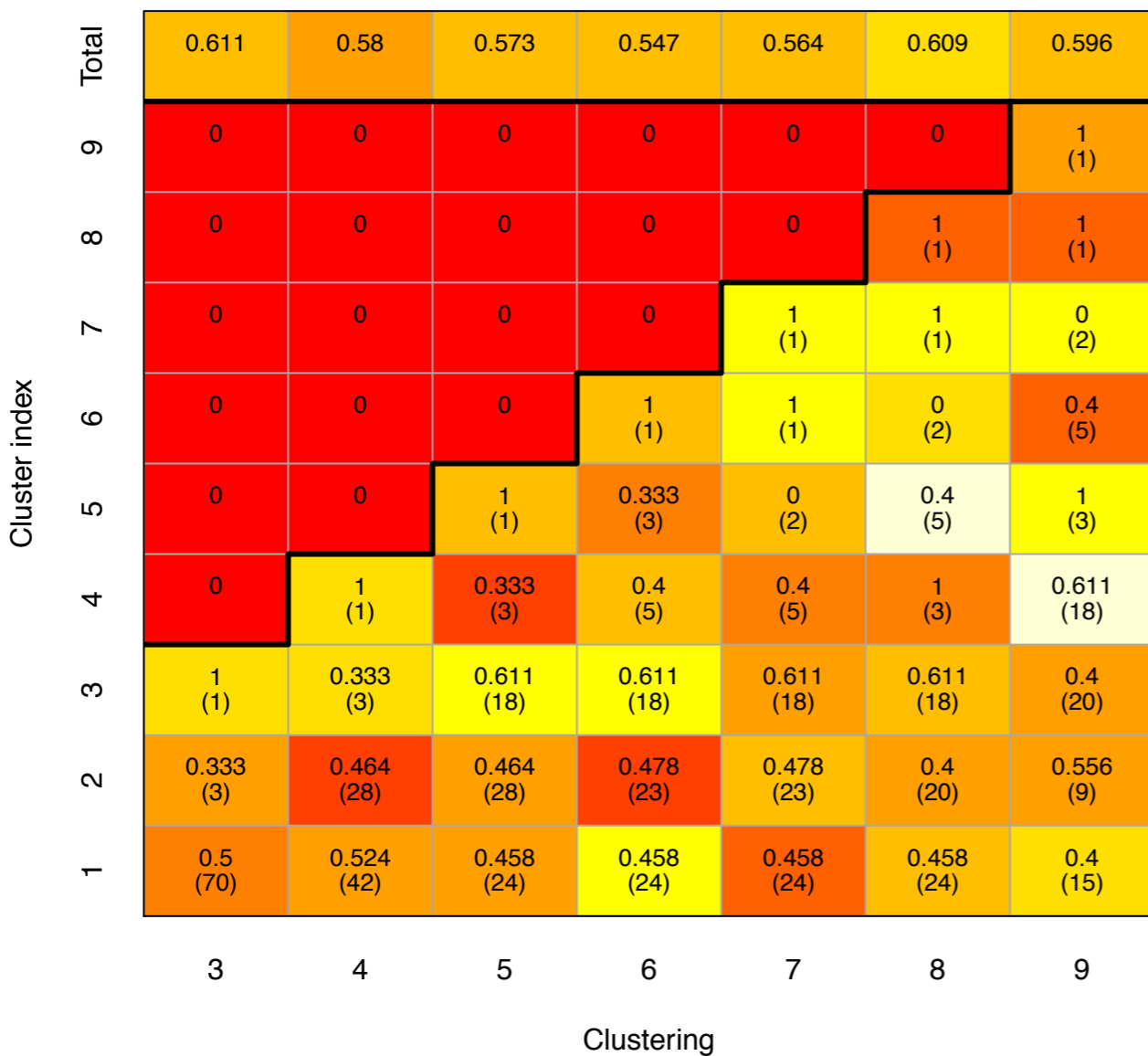
The score, a modified version of the “total variation” of a feature vector:

$$S_{TOT} = \frac{\sum_{i=1}^{N_{Cl}} S_i}{N_{Cl}} = \frac{\sum_{i=1}^{N_{Cl}} \left( \sum_{j=1}^{K-1} \|f_{ij} - f_{i(j+1)}\| \right)}{N_{Cl}}$$

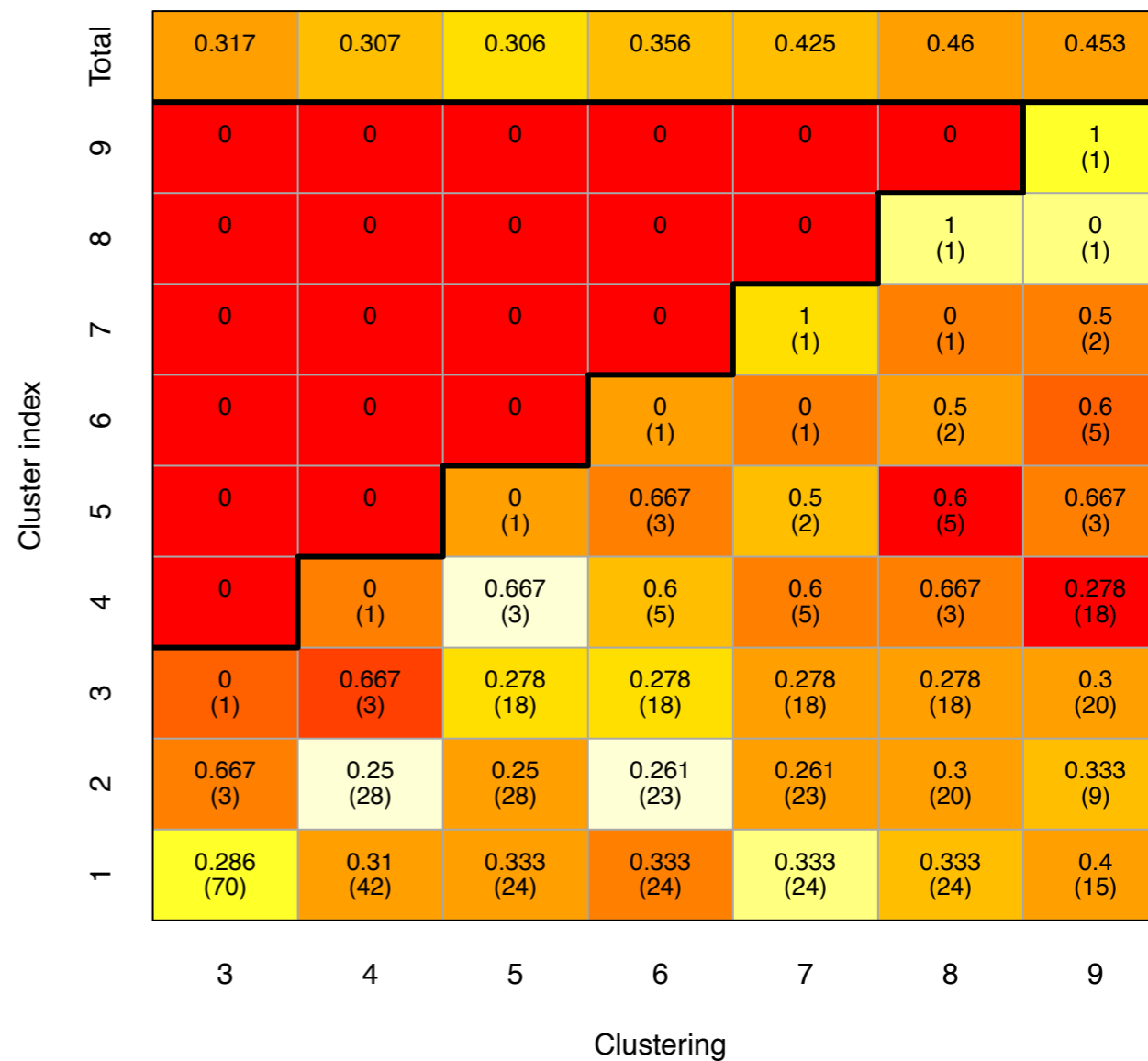


# Scores for single labels and clusters

Score distribution for HC clustering, relative to HR.ms. label

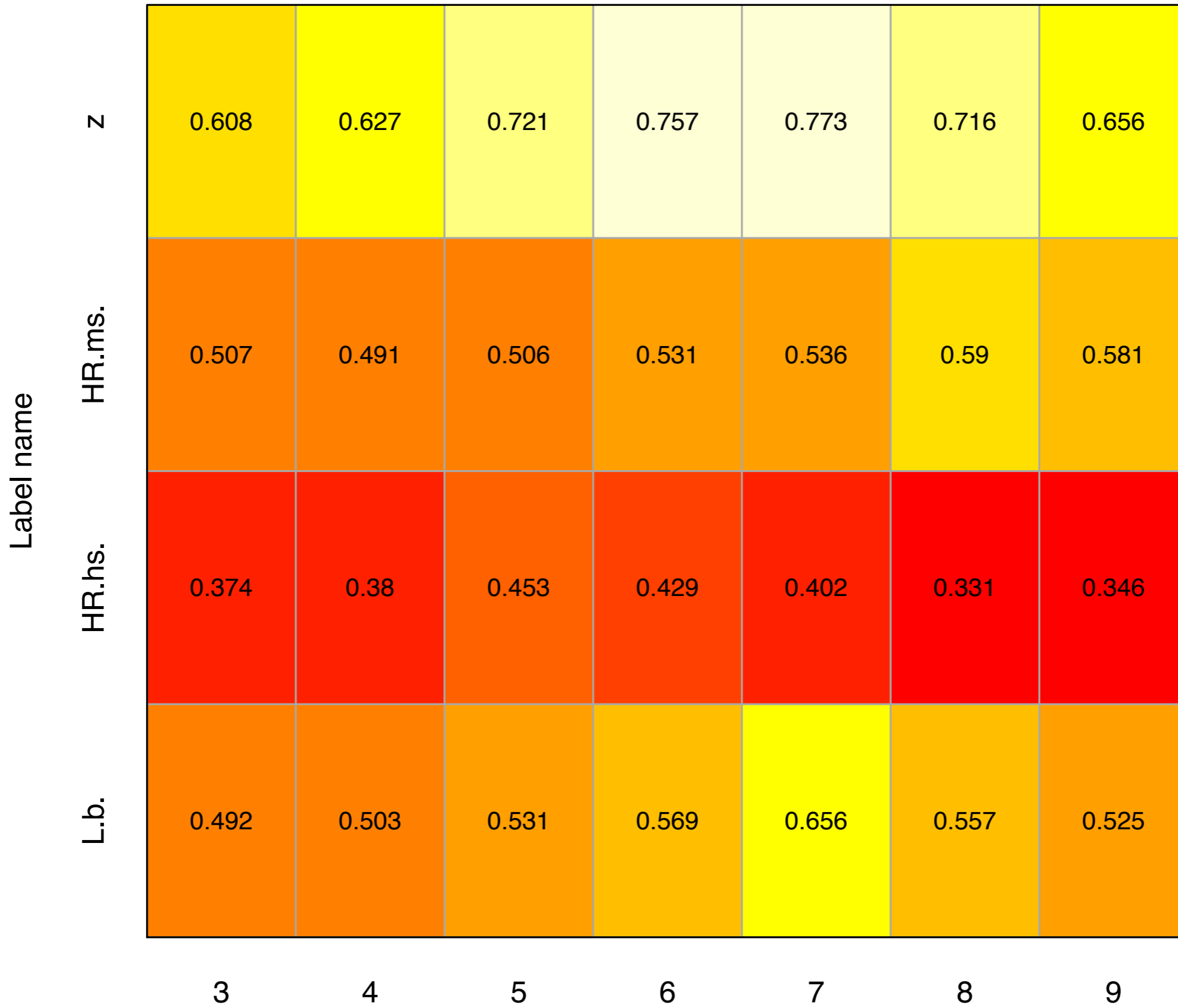


Score distribution for HC clustering, relative to HR.hs. label

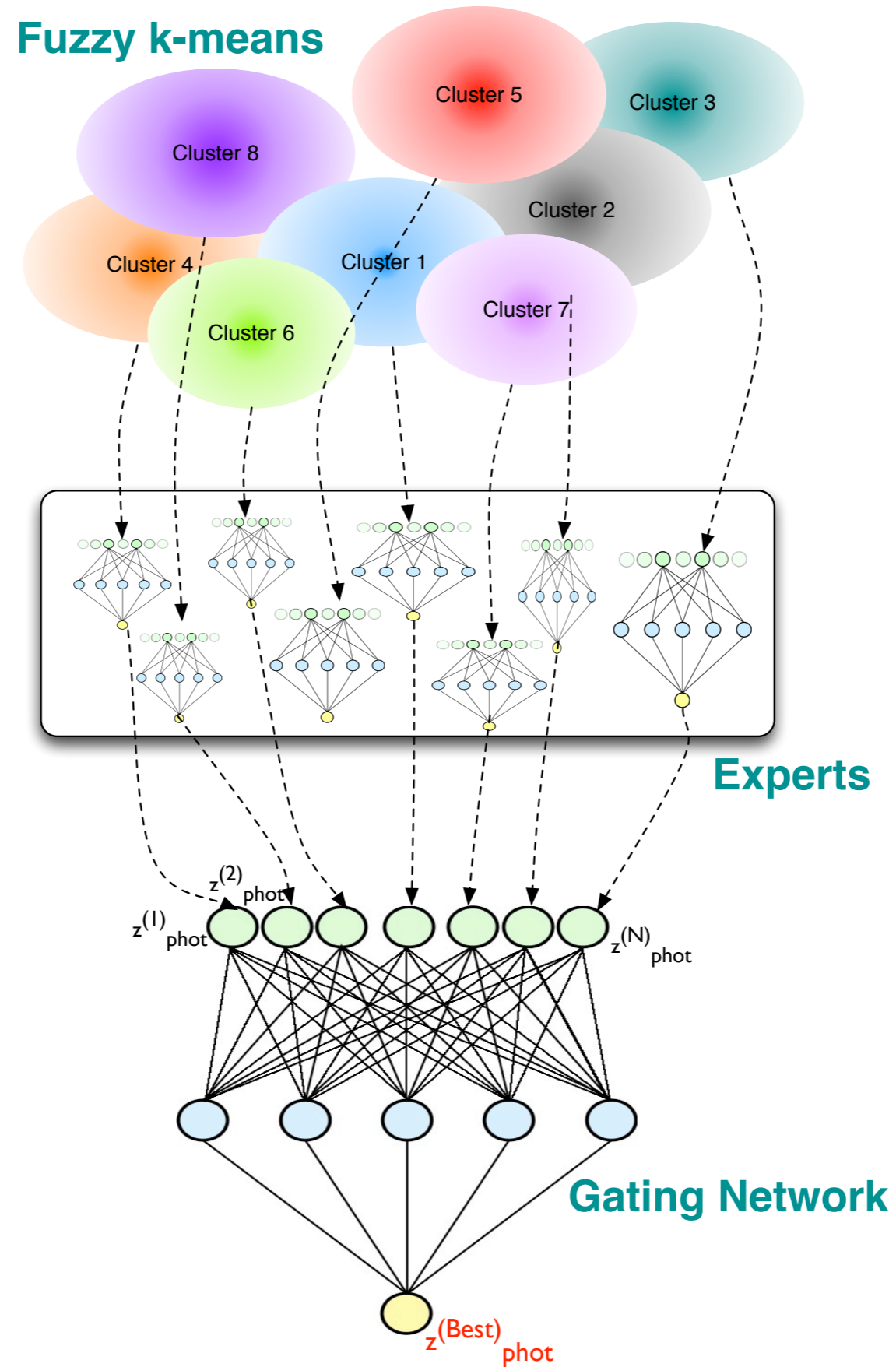




# Scores for clusterings

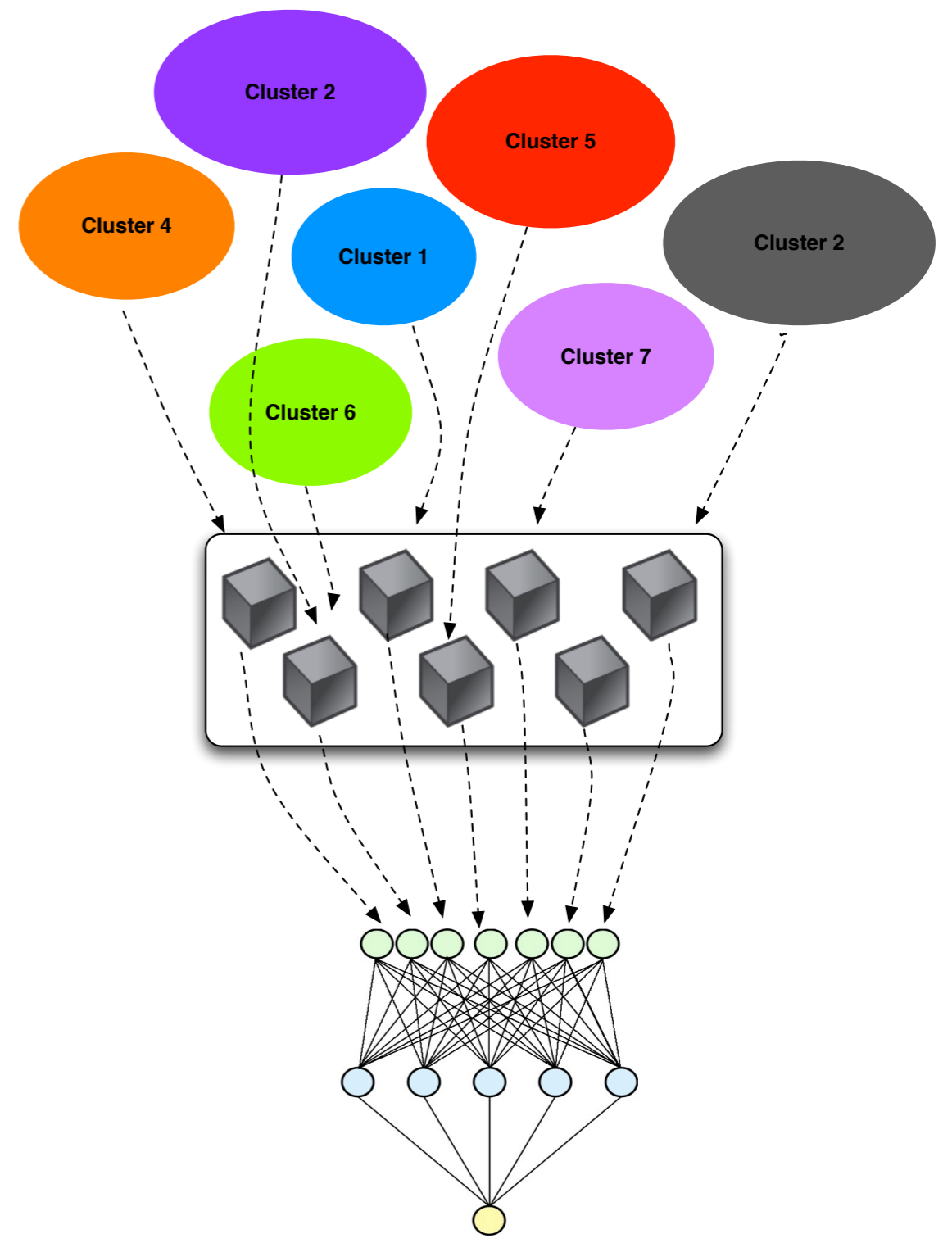
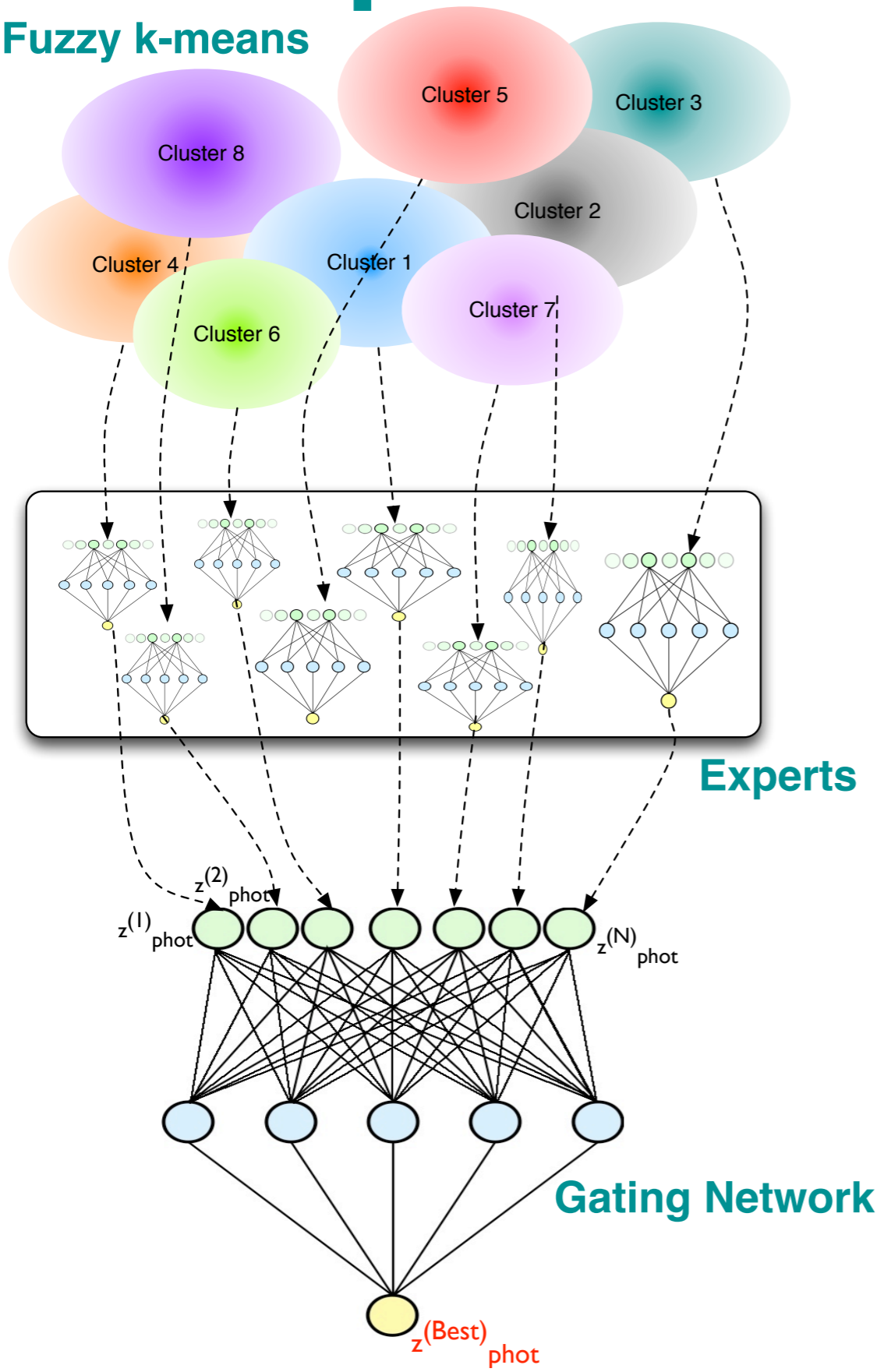


# Experts for Z<sub>phot</sub>



# Experts for other problems

Fuzzy k-means



# Conclusions

## Results for the first dataset and Chandra COSMOS survey with HC

### Open issues:

clustering with upper limits/no detection/missing data (model based? simulations?)

labels “binning”: is there a data-driven way to accomplish this? Co-clustering is being explored.

### Integration of different methods in a Gated Expert model:

example 1: template fitting and machine learning experts for  $z_{\text{phot}}$  calculation.

example 2: extraction of candidate quasars from optical photometric quasars.