A Set Correlation Model for Partitional Clustering

Michael E. HOULE

NII

Nguyen Xuan VINH

National ICT Australia, University of New South Wales

Why

Traditional clustering-based data discovery tools often fail when the data sets are very large or have many attributes, or have performance parameters that are very difficult to tune in practice. We have developed a novel formulation for partitional data clustering that overcomes many of these difficulties, based on the generic Relevant Set Correlation (RSC) model.

What

The new model, GlobalRSC, resembles the famous Kmeans clustering heuristic, but with a sharedneighbor similarity measure instead of the Euclidean distance. Unlike K-means and most other clustering heuristics that can only work with real-valued data and distance measures taken from specific families, GlobalRSC has the advantage that it can work with any distance measure, and any data representation.

