

Poster K-10

Recursive Top-Down Quantum Clustering of Biological Data



Authors:

Varshavsky Roy (*School of Computer Science and Engineering, The Hebrew University of Jerusalem, Israel*)

Horn David (*School of Physics and Astronomy, Tel Aviv University, Israel*)

Linial Michal (*Department of Biological Chemistry, The Hebrew University of Jerusalem, Israel*)

Short Abstract: We present a top-down, recursive clustering algorithm, for various datasets (e.g., gene-expression, protein sequences). The algorithm is successfully tested on several benchmarks. This novel algorithm identifies some internal structures, while other algorithms fail. Investigating these new structures suggests a new biological interpretation.

Long Abstract:

Motivation: Hierarchical tree is a natural way to present different granularity in gene expression, proteins, function annotations and more. In most of the cases where such representation is desired, the number of clusters may be large and a priori unknown, hence global clustering is insufficient (e.g., K-Means, QC). Many of the available known hierarchical clustering algorithms are bottom-up (agglomerative). The main drawback is that arbitrary considerations (e.g., dependency on the similarity and metric representations, determining the number of clusters) are applied. We present a new procedure, Top-Down Quantum Clustering (TDQC), taking advantage of the potential value assigned to each data point by QC. It overcomes the tendency of QC to generate a small number of clusters and miss some internal structures.

Method: Top-down, recursive clustering

The Algorithm:

1. [Optionally] To the original dataset apply preprocessing:
2. Run QC
3. Divide the data to:
 - a. The cluster with the global minimum
 - b. The remaining clusters
4. Recursively go to 1

Stop dividing when a set includes ≤ 2 elements

Results:

TDQC was applied to gene expression data and tested on several benchmarks. Here we illustrate the results for the Spellman experiment that analyzes gene expression stages in the yeast genome. According to their expert view there are 5 major partitions of the 800 genes that correspond to the phases of the cell cycle. Fig. 2 shows a graphic result of the TDQC algorithm. We were able to automatically partition the data to the main 5 groups. Interestingly, our results suggest some refinement of the expert view. A small subset of the assigned genes in S/G2 (cyan) is more likely to be associated with G2/M (red). The biological relevance of our results and the power of the method and its application to other protein sets and to gene expression data will be discussed.