

GeneCluster 2.0: An advanced toolset for bioarray analysis

M. Reich, K. Ohm, M. Angelo, P. Tamayo, J.P. Mesirov

Center for Genome Research, Massachusetts Institute of Technology
Cambridge, MA 02141

ABSTRACT

Summary: GeneCluster 2.0 is a software package for analyzing gene expression and other bioarray data, giving users a variety of methods to build and evaluate class predictors, visualize marker lists, cluster data, and validate results. GeneCluster 2.0 greatly expands the data analysis capabilities of GeneCluster 1.0 by adding classification, class discovery, and permutation test methods. It includes algorithms for building and testing supervised models using weighted voting (WV) and k-nearest neighbors (kNN) algorithms, a module for systematically finding and evaluating clustering via self-organizing maps (SOM), and modules for marker gene selection and heat map visualization that allow users to view and sort samples and genes by many criteria. GeneCluster 2.0 is a standalone Java application and runs on any platform that supports the Java Runtime Environment version 1.3.1 or greater.

Availability:

<http://www.broad.mit.edu/cancer/software>

Contact: gc-info@broad.mit.edu

GeneCluster 1.0, released in June of 1999, implemented the Self-Organizing Maps (SOM) algorithm popularized in Tamayo et al. 1999 as well as various preprocessing methodologies that have since become standard in microarray analysis. Since its public release, GeneCluster 1.0 has been downloaded by over 3000 users. In June 2002, a greatly enhanced version, GeneCluster 2.0, was released. It enhances the SOM clustering capability of GeneCluster 1.0 and adds modules for feature selection, predictive modeling, validation, and annotated heat-map visualization.

To enter data into GeneCluster 2.0, users provide a file containing the gene intensity values for each sample scanned along with associated annotations. Any required image processing steps such as background subtraction and dye correction should already be done.

Once loaded, an experiment file can be filtered and normalized with a variety of algorithms. These include thresholding, scaling, normalizing to a given mean and variance, fold-change, and exclusion of high and low scoring features. Users can also randomize a dataset by bootstrap sampling columns from the data set with replacement.

Additional available transformations are transposition and linear and logarithmic scaling.

Unsupervised learning, or clustering, is implemented by a SOM algorithm that has been extended to allow users to perform batch runs varying the number of clusters and cluster geometries. Results can be viewed in a visualizer that displays cluster profiles and relevant cluster member information.

GeneCluster 2.0 adds supervised learning methods (Golub et al. 1999, Slonim et al. 2000), in which a user trains an algorithm with labeled samples to obtain a model that can predict the class of an unknown sample. The user can choose between a weighted voting (WV) and a K-nearest neighbors (kNN) algorithm, as well as the number of features in the model and how those features are chosen.

GeneCluster 2.0 offers an extensive set of tools and methods to perform marker or feature selection. This module allows users to find all genes in the neighborhood of a single gene or the genes that most closely fit an expression profile. Distance metrics include Euclidean, Pearson correlation, and Manhattan distance, and statistical measures of correlation include the signal-to-noise score and t-test. Markers can be visualized as an interactive heat map (color-gram), described below, that allows users to view graphically the expression levels for a set of markers.

To determine the statistical significance of a marker, GeneCluster 2.0 provides a permutation test. This test compares the density of neighborhoods of genes correlated with a valid class distinction profile (for example, tumor vs. normal) to neighborhoods of profiles of randomly assigned distinctions (see Golub et al. (1999) and Slonim et al. (2000) for details). The user can set the significance level and number of permutations, or can use default values.

Once a predictive model is built, its accuracy can be measured through either leave-one-out or n-fold cross-validation, or via random train/test splits. The predictor's accuracy is given as the average across all the split-train-test cycles. GeneCluster 2.0 also provides a confusion matrix, which gives users a tabular view of the number of correct versus incorrect class predictions.

