# Data Mining Supercomputing with SAS JMP® Genomics

Dr. Richard S. SEGALL*
Arkansas State University, Department of Computer & Information Technology
State University, AR 72467-0130, USA, rsegall@astate.edu

Dr. Qingyu ZHANG*
Arkansas State University, Department of Computer & Information Technology
State University, AR 72467-0130, USA, qzhang@astate.edu

and

Ryan M. PIERCE
Arkansas State University, Student Affairs Technology Services,
State University, AR 72567-0348, USA, rmpierce@astate.edu

## ABSTRACT

JMP® Genomics is statistical discovery software that can uncover meaningful patterns in high-throughput genomics and proteomics data. JMP® Genomics is designed for biologists, biostatisticians, statistical geneticists, and those engaged in analyzing the vast stores of data that are common in genomic research (SAS, 2009).

Data mining was performed using JMP® Genomics on the two collections of microarray databases available from National Center for Biotechnology Information (NCBI) for lung cancer and breast cancer. The Gene Expression Omnibus (GEO) of NCBI serves as a public repository for a wide range of high-throughput experimental data, including the two collections of lung cancer and breast cancer that were used for this research. The results for applying data mining using software JMP® Genomics are shown in this paper with numerous screen shots.

**Keywords:** Microarray databases, Lung Cancer, Breast Cancer, Data Mining, Supercomputing, Gene Expression Omnibus (GEO), SAS JMP® Genomics.

## 1. BACKGROUND

The software used in this research is JMP® Genomics from SAS Institute, Inc. of Cary, NC that according to Product Brief of SAS (2009) dynamically links advanced statistics with graphics to provide a complete and comprehensive picture of results, whether the data comes from traditional microarray studies or data summarized from next-generation technologies. Preliminary work done by the authors for the visualization by supercomputing data mining using JMP® Genomics from SAS for similar data was presented in Segall et al. (2010) and (2009).

Some of the previous research that has been performed by others in the area of applications of supercomputing to data mining include those of Zaki et al. (1996) for parallel data mining, Thoennes and Weems (2003) for performance of data mining on complex microprocessors, and data mining of large datasets with geospatial information by the image spatial data analysis group (2009) and University of Illinois at Urbana-Champaign, and Wilkins-Diehr and Mirman (2009) for on-demand supercomputing for emergencies that includes discussions for applications to breast cancer diagnosis.

## 2. DATA

The Gene Expression Omnibus (GEO) is a public repository that archives and freely distributes microarray, next-generation sequencing, and other forms of high-throughput functional genomic data submitted by the scientific community. These data include single and dual channel microarray-based experiments measuring mRNA, miRNA, genomic DNA (including arrayCGH, ChIP-chip, and SNP), and protein abundance, as well as non-array techniques such as serial analysis of gene expression (SAGE), and various types of next-generation sequence data. In addition to data storage, a collection of web-based interfaces and applications are available to help users

query and download the experiments and gene expression patterns stored in GEO.

The data sets used in the research presented in this paper are those from the Gene Expression Omnibus (GEO) from the National Center of Biotechnology Information (NCBI). One set of data is that of expression data for lung cancer that was made public on August 30, 2008; and the other is that for gene expression profiling in breast cancer that was made public in February 2006.

### Lung Cancer Data Used in This Paper

According to NCBI (2007), the detection, treatment, and prediction of outcome for lung cancer patients increasingly depend on a molecular understanding of tumor development and sensitivity of lung cancer to therapeutic drugs.

NCI (2007) states that the application of genomic technologies, such as microarray, is widely used to monitor global gene expression and has built up invaluable information and knowledge, which is essential to the discovery of new insights into the mechanisms common to cancer cells, resulting in the identification of unique, identifiable signatures and specific characteristics. According to NCBI (2007) it is likely that application of microarray may revolutionize many aspects of lung cancer being diagnosed, classified, and treated in the near future. NCBI (2007) used microarrays to detail the global gene expression patterns of lung cancer.

The overall design of NCBI (2007) as used in this paper consisted of adjacent normal-tumor matched lung cancer samples that were selected at early and late stages for RNA extraction and hybridization on Affymetrix microarrays. A total of 66 samples were used for microarray analysis in NCBI (2007), including pairwise samples from 27 patients, who underwent surgery for lung cancer at the Taipei Veterans General Hospital, two tissue mixtures from the Taichung Veterans General Hospital, two commercial human normal lung tissues, one immortalized, nontumorigenic human bronchial epithelial cell line, and 7 lung cancer cell lines.

### Breast Cancer Data Used in This Paper

The breast cancer data set used in this research was obtained on the web from NCBI (2006), which analyzed microarray data from 189 invasive breast carcinomas and from three published gene expression datasets from breast carcinomas. NCBI (2006) identified differentially expressed genes in a training set of 64 estrogen receptor (ER)-positive tumor samples by comparing expression profiles between histologic grade 3 tumors and histologic grade 1 tumors and used the expression of these genes to define the gene expression grade index. The data set for the figures generated in this paper consisted of over 22,000 rows representing different variables.

The breast cancer data presented by NCBI (2006) was from 597 independent tumors were used to evaluate the association between relapse-free survival and the gene expression grade index in a Kaplan-Meier analysis. All statistical tests performed by NCBI (2006) were two-sided. The overall design of NCBI (2006) was 64 microarray experiments from primary breast tumors used as training set to identify genes differentially expressed in grades 1 and 3. NCBI (2006) design included 129 microarray experiments from primary breast tumors of untreated patients used as validation set to validate the list of genes and its correlation with survival.

## 3. RESULTS

### Data Mining Performed Using Sas Jmp® Genomics For Lung Cancer Data

Figure 1 shows the window called "basic expression workflow" that is the process that runs a basic workflow for expression data used to select variables of interest.

The data used for the lung cancer and its associated tumors consisted of over 22,000 rows representing genes and 54 columns representing samples as shown in Figure 2.

Our research using SAS JMP® Genomics yielded distributions plots of conditions, patients and characteristics; correlation analysis of principle components as shown in Figure 3 which shows "normal" versus "cancer" in the scatterplots , and dendograms of hierarchical clustering as shown in Figure 4. Figure 5 shows a Volcano plot of the summary plot of individual genes and their differences in condition of cancer from normal tissues.

Our research performed some predictive modeling using SAS JMP® Genomics that yielded one-way analysis plots of fitting a selected gene number 1773 by condition and also by patient as shown in Figure 6.
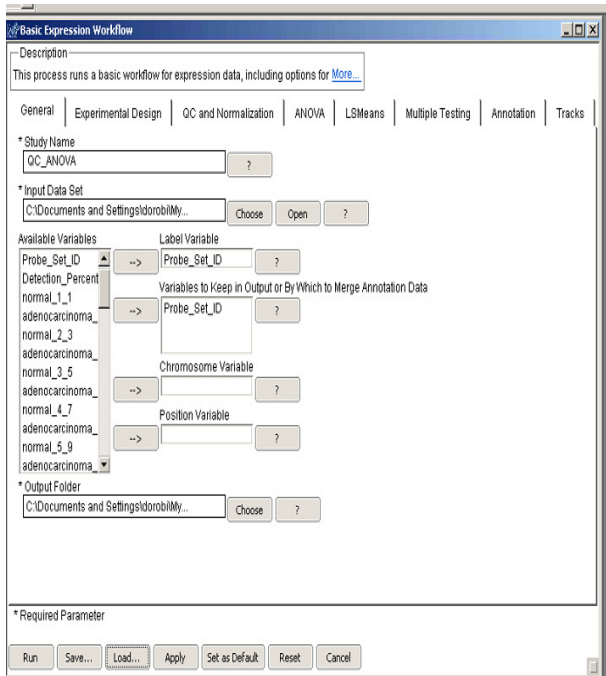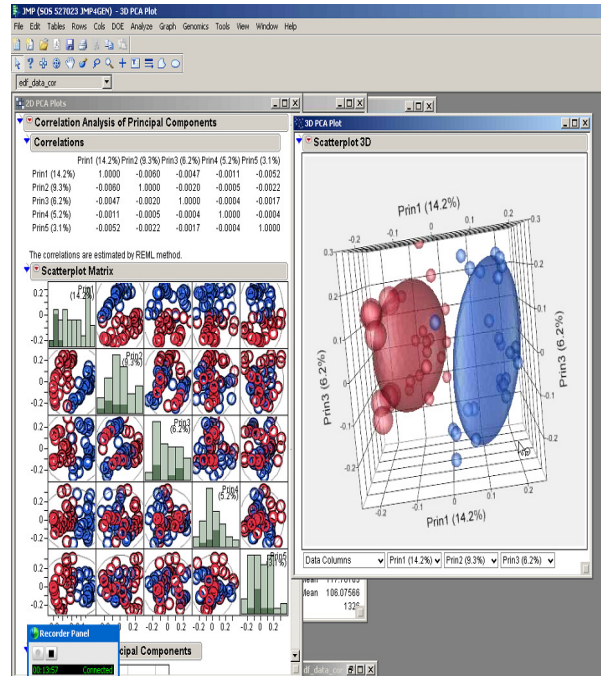
Figure 1. Basic expression workflow



Figure 3. Correlation analysis of principle components
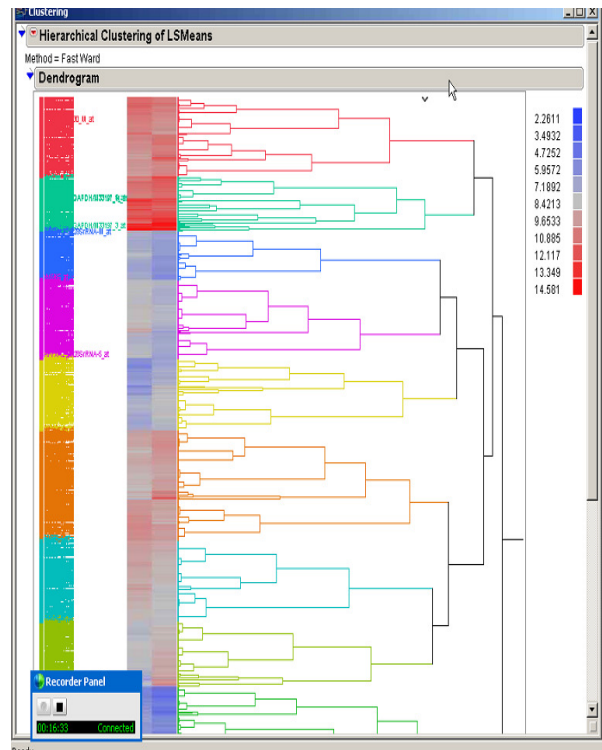


Figure 2. Adenocarcinoma Cancer Data



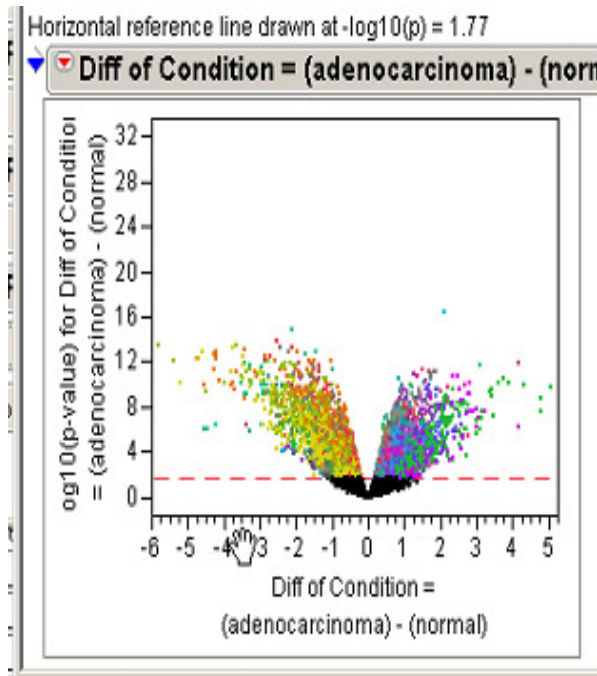Figure 4. Dendograms of hierarchical clustering

Figure 5. Volcano plot



Figure 6. One-way analysis plots

## Data Mining Performed Using Sas Jmp® Genomics For Breast Cancer Data

Box plots of a 50-iteration simple random cross-validation root mean square error (RMSE) are shown I Figure 7 for five different models. In this Figure 7, the dependent variables is "grade" for level of severity of cancer tumors in breast cancer, and the predictor continuous variables is "age". Cross validation was performed that on predictive model settings selected and compares the results.

Figure 8 shows the 235 predictors ranked for each of the models used as training set data. Figure 9 shows the Heat Map and Dendograms for breast cancer data which uses colors to indicate the intensity of correlation. The lower right corner of Figure 9 Heat Map is in red indicating highly correlated microarrays.

The frequency distributions are shown in Figure 10 that were obtained by highlighting the selected portion of Figure 9 Heat Map and indicate no grade 3 tumors. Partitioning the decision trees as shown in Figure 11 shows contingency analysis of predicted class by grade of tumor, and also the distribution data by true grade of tumor, actual probabilities, and correct predictions.
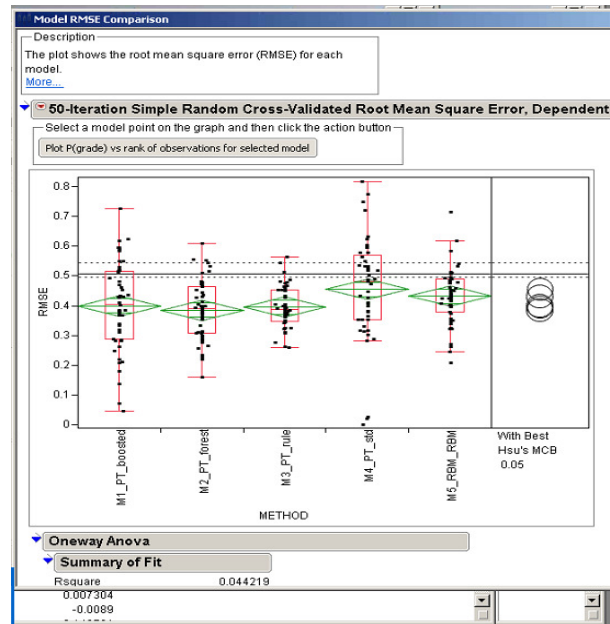


Figure 7 Five different models
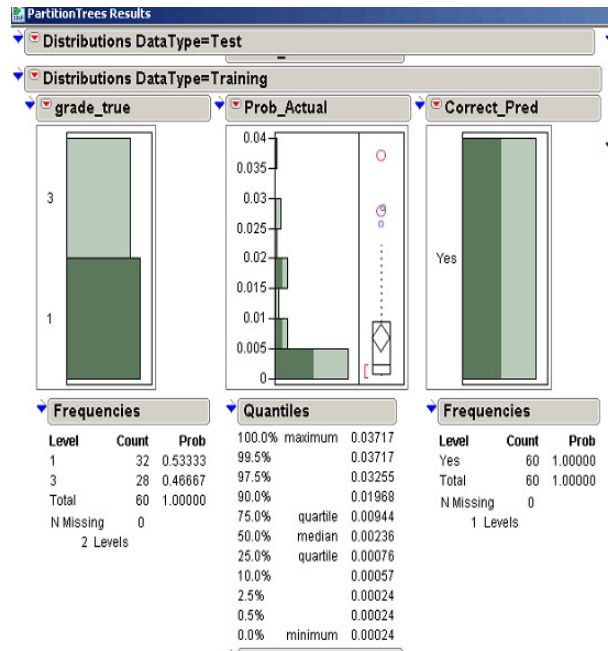
Figure 8 Training set data
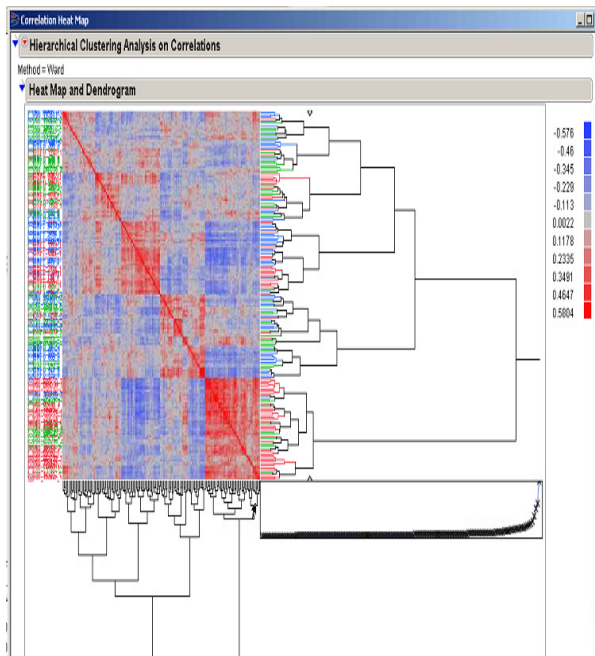


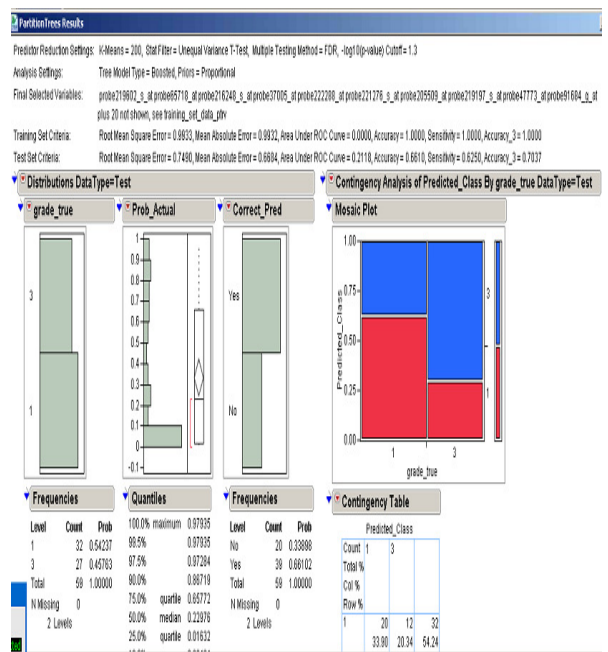Figure 10 Frequency distributions



Figure 9 Heat Map



Figure 11 Partitioning the decision trees

## 4. CONCLUSIONS AND SUMMARY

This paper emphasizes the usefulness of SAS JMP® Genomics with supercomputing and data mining. This research illustrates genetic visualization for the analysis and modeling of microarray databases for both lung and breast cancer as a tool for better understanding of the consequences of these diseases and for potential strategies for their treatments

## 5. ACKNOWLEDGEMENTS

## 6. REFERENCES

1. Image Spatial Data Analysis Group (2009), National Center for Supercomputing, University of Illinois at Urbana-Champaign, http://isda.ncsa.illinois.edu

2. NCBI (2007), "Expression data from Lung Cancer", Gene Expression Omnibus (GEO), Series GSE7670, **National Center for Biotechnology Information**, http://www.ncbi.nlm.nih.gov/geo/

3. NCBI (2006), "Gene Expression Profiling in Breast Cancer: Understanding the Molecular Basis of Histologic Grade to Improve Prognosis, Gene Expression Omnibus (GEO), Series GSE2990, **National Center for Biotechnology Information**, http://www.ncbi.nlm.nih.gov/geo/

4. SAS (2009), **JMP® Genomics 4.0 Product Brief**, http://www.jmp.com/software/genomics/pdf/103112_j mpg4_prodbrief.pdf

5. Segall, Richard S., Zhang, Qingyu and Pierce, Ryan M.(2009), "Visualization by Supercomputing Data Mining", **Proceedings of the 4th INFORMS Workshop on Data Mining and System Informatics**, San Diego, CA, October 10, 2009.

6. Segall, Richard S., Zhang, Qingyu and Pierce, Ryan M.(2010), "Data Mining Supercomputing with SAS JMP® Genomics: Research-in-Progress" submitted to **Proceedings of 2010 Conference on Applied Research in Information Technology**, sponsored by Acxiom Laboratory of Applied Research (ALAR),University of Central Arkansas (UCA), Conway, AR, April 9, 2010.

7. Sotiriou C, Wirapati P, Loi S, Harris A et al. (2006), Gene expression profiling in breast cancer: understanding the molecular basis of histologic grade to improve prognosis. **J National Cancer Institute**, February 15; volume 98, number 4, pp. 262-72. PMID: 16478745

8. Su LJ, Chang CW, Wu YC, Chen KC et al.(2007), "Selection of DDX5 as a novel internal control for Q-RT-PCR from microarray data using a block bootstrap re-sampling scheme", **BMC Genomics** 2007 June 1; volume 8, number 140. PMID: 17540040

9. Thoennes, M.S., and Weems, C.C. (2003), "Exploration of the performance of a data mining application via hardware based monitoring", **The Journal of Supercomputing**, volume 26, pp. 25-42.

10. Wilkins-Diehr, N. and Mirman, I. (2008), "On-Demand supercomputing for emergencies", **Design Engineering Technology News Magazine**, February 1, http://www.deskeng.com/articles/aaagtk.htm

11. Zaki, M.J., Ogihara, M., Parthasararthy, S., and Li, W. (1996), "Parallel data mining for association rules on shared-memory multi-processors", **Proceedings of the 1996 ACM/IEEE Conference on Supercomputing**, Pittsburgh, PA, Article Number 43, ISBN 0-89791-854-1.