

# Chapter 3.16

## AI Methods for Analyzing Microarray Data

**Amira Djebbari**

*National Research Council Canada, Canada*

**Aedín C. Culhane**

*Harvard School of Public Health, USA*

**Alice J. Armstrong**

*The George Washington University, USA*

**John Quackenbush**

*Harvard School of Public Health, USA*

### INTRODUCTION

Biological systems can be viewed as information management systems, with a basic instruction set stored in each cell's DNA as "genes." For most genes, their information is enabled when they are transcribed into RNA which is subsequently translated into the proteins that form much of a cell's machinery. Although details of the process for individual genes are known, more complex interactions between elements are yet to be discovered. What we do know is that diseases can result

if there are changes in the genes themselves, in the proteins they encode, or if RNAs or proteins are made at the wrong time or in the wrong quantities.

Recent advances in biotechnology led to the development of DNA microarrays, which quantitatively measure the expression of thousands of genes simultaneously and provide a snapshot of a cell's response to a particular condition. Finding patterns of gene expression that provide insight into biological endpoints offers great opportunities for revolutionizing diagnostic and prognostic medicine and providing mechanistic insight in data-driven research in the life sciences, an area with a great need for advances, given the urgency

DOI: 10.4018/978-1-60960-561-2.ch316

associated with diseases. However, microarray data analysis presents a number of challenges, from noisy data to the curse of dimensionality (large number of features, small number of instances) to problems with no clear solutions (*e.g.* real world mappings of genes to traits or diseases that are not yet known).

Finding patterns of gene expression in microarray data poses problems of class discovery, comparison, prediction, and network analysis which are often approached with AI methods. Many of these methods have been successfully applied to microarray data analysis in a variety of applications ranging from clustering of yeast gene expression patterns (Eisen *et al.*, 1998) to classification of different types of leukemia (Golub *et al.*, 1999). Unsupervised learning methods (*e.g.* hierarchical clustering) explore clusters in data and have been used for class discovery of distinct forms of diffuse large B-cell lymphoma (Alizadeh *et al.*, 2000). Supervised learning methods (*e.g.* artificial neural networks) utilize a previously determined mapping between biological samples and classes (*i.e.* labels) to generate models for class prediction. A k-nearest neighbor (k-NN) approach was used to train a gene expression classifier of different forms of brain tumors and its predictions were able to distinguish biopsy samples with different prognosis suggesting that microarray profiles can predict clinical outcome and direct treatment (Nutt *et al.*, 2003). Bayesian networks constructed from microarray data hold promise for elucidating the underlying biological mechanisms of disease (Friedman *et al.*, 2000).

## BACKGROUND

Cells dynamically respond to their environment by changing the set and concentrations of active genes by altering the associated RNA expression. Thus “gene expression” is one of the main determinants of a cell’s state, or phenotype. For example, we can investigate the differences between a normal cell and a cancer cell by examining their relative gene expression profiles.

Microarrays quantify gene expression levels in various conditions (such as disease *vs.* normal) or across time points. For  $n$  genes and  $m$  instances (biological samples), microarray measurements are stored in an  $n$  by  $m$  matrix where each row is a gene, each column is a sample and each element in the matrix is the expression level of a gene in a biological sample, where samples are instances and genes are features describing those instances. Microarray data is available through many public online repositories (Table 1). In addition, the Kent-Ridge repository (<http://sdmc.i2r.a-star.edu.sg/rp/>) contains pre-formatted data ready to use with the well-known machine learning tool Weka (Witten & Frank, 2000).

Microarray data presents some unique challenges for AI such as a severe case of the curse of dimensionality due to the scarcity of biological samples (instances). Microarray studies typically measure tens of thousands of genes in only tens of samples. This low case to variable ratio increases the risk of detecting spurious relationships. This problem is exacerbated because microarray data contains multiple sources of within-class variability, both technical and biological. The high

Table 1. Some public online repositories of microarray data

Name of the repository	URL
ArrayExpress at the European Bioinformatics Institute	<a href="http://www.ebi.ac.uk/arrayexpress/">http://www.ebi.ac.uk/arrayexpress/</a>
Gene Expression Omnibus at the National Institutes of Health	<a href="http://www.ncbi.nlm.nih.gov/geo/">http://www.ncbi.nlm.nih.gov/geo/</a>
Stanford microarray database	<a href="http://smd.stanford.edu/">http://smd.stanford.edu/</a>
Oncomine	<a href="http://www.oncomine.org/main/index.jsp">http://www.oncomine.org/main/index.jsp</a>

6 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:  
[www.igi-global.com/chapter/methods-analyzing-microarray-data/53625](http://www.igi-global.com/chapter/methods-analyzing-microarray-data/53625)

## Related Content

---

### Analysis and Quantification of Motion within the Cardiovascular System: Implications for the Mechanical Strain of Cardiovascular Structures

Spyretta Golemati, John Stoitsis and Konstantina S. Nikita (2009). *Handbook of Research on Advanced Techniques in Diagnostic Imaging and Biomedical Applications* (pp. 34-47).

[www.irma-international.org/chapter/analysis-quantification-motion-within-cardiovascular/19586](http://www.irma-international.org/chapter/analysis-quantification-motion-within-cardiovascular/19586)

### Merging Different Datasets to Allow for a Complete Analysis (Inpatient, Outpatient, Physician Visits, Medications)

Patricia Cerrito and John Cerrito (2010). *Clinical Data Mining for Physician Decision Making and Investigating Health Outcomes: Methods for Prediction and Analysis* (pp. 116-153).

[www.irma-international.org/chapter/merging-different-datasets-allow-complete/44269](http://www.irma-international.org/chapter/merging-different-datasets-allow-complete/44269)

### Quantitative Analysis of Hysteroscopy Imaging in Gynecological Cancer

Marios Neofytou, Constantinos Pattichis, Vasilios Tanos, Marios Pattichis and Eftyvoulos Kyriacou (2011). *Clinical Technologies: Concepts, Methodologies, Tools and Applications* (pp. 949-964).

[www.irma-international.org/chapter/quantitative-analysis-hysteroscopy-imaging-gynecological/53630](http://www.irma-international.org/chapter/quantitative-analysis-hysteroscopy-imaging-gynecological/53630)

### General Idea of the Proposed System

Piotr Augustyniak and Ryszard Tadeusiewicz (2009). *Ubiquitous Cardiology: Emerging Wireless Telemedical Applications* (pp. 145-154).

[www.irma-international.org/chapter/general-idea-proposed-system/30489](http://www.irma-international.org/chapter/general-idea-proposed-system/30489)

### Future Perspective: Data Validity-Driven Report Optimization

Piotr Augustyniak and Ryszard Tadeusiewicz (2009). *Ubiquitous Cardiology: Emerging Wireless Telemedical Applications* (pp. 296-312).

[www.irma-international.org/chapter/future-perspective-data-validity-driven/30495](http://www.irma-international.org/chapter/future-perspective-data-validity-driven/30495)