## Machine Learning and Genetic Microarrays

### Jude Shavlik & David Page University of Wisconsin-Madison

Copyrighted © 2003 by Jude Shavlik and David Page

## Goals

 Learn about microarray technology
 See some ML problem formulations in the computational-biology literature
 A little on experimental pitfalls to avoid
 Overviews of <u>newest</u> "high-throughput" molecular-level data being gathered
 <u>Very little</u> on the details of machine learning algorithms

# Outline

Molecular Biology and Microtechnology
 Machine Learning Applications

- Technological: Designing Microarrays
- Medical: Predicting Disease (Diagnosis, Prognosis, & Treatment)
- Biological: Constructing Pathway Models
- Looking Ahead: Related Technologies

3

## Outline

Molecular Biology and Nanotechnology

- Machine Learning Applications
  - Technological: Designing Microarrays
  - Medical: Predicting Disease (Diagnosis, Prognosis, & Treatment)
  - Biological: Constructing Pathway Models
- Looking Ahead: Related Technologies



image from the DOE Human Genome Program http://www.ornl.gov/hgmis













## Cleaning Up the Data – "Controlling Variance"

Often look at *relative* expression levels

Measurement(cancerCell) / Measurement(normalCell)

Often correct for small values

MeasurementToUse(Gene1) = ActualMeasurement(Gene1) + Constant

Often use a mismatch ("near miss") probe

MeasurementToUse(Probe1) = ActualMeasurement(Probe1) – ActualMeasurement(MismatchProbe1)

19

## Outline

Molecular Biology and Microtechnology

- Machine Learning Applications
  - Technological: Designing Microarrays
  - Medical: Predicting Disease (Diagnosis, Prognosis, & Treatment)
  - Biological: Constructing Pathway Models
- Looking Ahead: Related Technologies

Need to pick we want to	k 5-15 probes for each gene monitor
<ul><li>Recall, gen</li><li>Can only cr</li></ul>	es about 1000 "bases" long reate probes about 24-bases long
Gene	
Probes —	
	21
	Goals
(	Cours
Probes shows and the second	ould bind tightly to target



## The Data

Tilings of 8 genes (from *E. coli* & *B. subtilus*)

- Every possible probe (~10,000 probes)
- Genes known to be "expressed" in sample

Gene Sequence:	GTAGCTAGCATTAGCATGGCCAGTCATG
Complement:	CATCGATCGTAATCGTACCGGTCAGTAC
Probe 1:	CATCGATCGTAATCGTACCGGTCA
Probe 2:	ATCGATCGTAATCGTACCGGTCAG
Probe 3:	TCGATCGTAATCGTACCGGTCAGT
• • •	•••

25

## Microarray that Created Examples

NimbleGen Systems, Inc.	control chip with subset	
	二日本市市市市市市市市市市市市市市市市市市市市市市市市市市市市市市市市市市市市	
	田 ししと 正論局部からから世界優勝の内の原来長手構成	
	二 2.5.日本当業業会が使一代 サラジ	
	□ □ □ □	
-		
		1
	一個自己自己以口自己自己的意思。	а.
	「「「」「「」「」「」「」「」「」「」「」「」「」「」「」「」」	
		1
		4
		а.
		_
	E CARACTER STATE	
	CARGED STATES CONTRACTOR AND ADDRESS AND ADDRESS ADDRES	

## The Features (Tobler et al., ISMB 2002)

Feature Name	Description
fracA, fracC, fracG, fracT	The fraction of A, C, G, or T in the 24-mer
fracAA, fracAC, fracAG, fracAT, fracCA, fracCC, fracCG, fracCT, fracGA, fracGC, fracGG, fracGT, fracTA, fracTC, fracTG, fracTT	The fraction of each of these dimers in the 24-mer
n1, n2,, n24	The particular nucleotide (A, C, G, or T) at the specified position in the 24-mer
d1, d2,, d23	The particular dimer (AA, AC,TT) at the specified position in the 24-mer











## Outline

Molecular Biology and Microtechnology

- Machine Learning Applications
  - Technological: Designing Microarrays
  - Medical: Predicting Disease (Diagnosis, Prognosis, & Treatment)
  - Biological: Constructing Pathway Models
- Looking Ahead: Related Technologies



## Two Ways to View The Data

Person Gene►	A28202_ac	AB00014_at	AB00015_at	
Person 1	1142.0	321.0	2567.2	
Person 2	586.3	586.1	759.0	
Person 3	105.2	559.3	3210.7	
Person 4	42.8	692.1	812.0	

## Data Points are Genes

Person Gene►	A28202_ac	AB00014_at	AB00015_at	
Person 1	1142.0	321.0	2567.2	
Person 2	586.3	586.1	759.0	
Person 3	105.2	559.3	3210.7	
Person 4	42.8	692.1	812.0	
	-		-	

39

## Data Points are Samples

Person Gene►	A28202_ac	AB00014_at	AB00015_at	
Person 1	1142.0	321.0	2567.2	
Person 2	586.3	586.1	759.0	
Person 3	105.2	559.3	3210.7	
Person 4	42.8	692.1	812.0	
		-		
			_	

## Supervision: Add Class Values

-				
Person Gene►	A28202_ac	AB00014_at	AB00015_at	 Class
Person 1	1142.0	321.0	2567.2	 normal
Person 2	586.3	586.1	759.0	 cancer
Person 3	105.2	559.3	3210.7	 normal
Person 4	42.8	692.1	812.0	 cancer
	•		•	
I				

41

### Supervised Learning Task 2

 Given: a set of microarray experiments, each done with mRNA from a <u>different patient</u> (same cell type from every patient)

Patient's expression values for each gene constitute the <u>features</u>, and patient's disease constitutes the <u>class</u>

Do: Learn a model that accurately predicts <u>class</u> based on features





## X-Val Accuracies for Multiple Myeloma (74 мм vs. 31 Normal)

Trees	98.1
Boosted Trees	99.0
SVMs	100.0
Vote	100.0
Bayes Nets	97.0





## Some Future Directions

Using gene-chip data to <u>select therapy</u> Predict which <u>therapy</u> gives best prognosis for patient

Comparing cancer with related benign conditions, rather than with normal Tougher, but may give more insight

51

## **Unsupervised Learning Task 1**

Given: a set of microarray experiments under different conditions

Do: <u>cluster</u> the genes, where a gene described by its expression levels in different experiments







# Outline

Molecular Biology and Microtechnology
 Machine Learning Applications

- Technological: Designing Microarrays
- Medical: Predicting Disease (Diagnosis, Prognosis, & Treatment)
- Biological: Constructing Pathway Models
- Looking Ahead: Related Technologies

59

## Some Biological Pathways

Regulatory pathways

- Nodes are labeled by genes
- Arcs denote influence on transcription
- G1 codes for P1, P1 inhibits G2's transcription

#### Metabolic pathways

- Nodes are metabolites, large biomolecules (eg, sugars, lipids, proteins and modified proteins)
- Arcs from biochemical reaction inputs to outputs
- Arcs labeled by enzymes (themselves proteins)



## Using Microarray Data Only

#### Regulatory pathways

- Nodes are labeled by genes
- Arcs denote influence on transcription
- G1 codes for P1, P1 inhibits G2's transcription

#### Metabolic pathways

- Nodes are metabolites, large biomolecules (eg, sugars, lipids, proteins, and modified proteins)
- Arcs from biochemical reaction inputs to outputs
- Arcs labeled by enzymes (themselves proteins)

63

## Supervised Learning Task 3

### Given: a set of microarray experiments for <u>same organism under different</u> <u>conditions</u>

Do: Learn graphical model that accurately predicts expression of some genes in terms of others











# Outline

Molecular Biology and Microtechnology
 Machine Learning Applications

- Technological: Designing Microarrays
- Medical: Predicting Disease (Diagnosis, Prognosis, & Treatment)
- Biological: Constructing Pathway Models
- Looking Ahead: Related Technologies

75

### Single-Nucleotide Polymorphisms

SNPs: Individual positions in DNA where variation is common

- Now 1.8 million known SNPs in humans
- Easier/faster/cheaper to measure SNPs than to completely sequence everyone

Motivation ...



### Example of SNP Data

Person SNP►	1			2		3	 CLASS
Person 1	С	Т	A	G	Т	Т	 old
Person 2	С	С	Α	G	С	Т	 young
Person 3	Т	Т	A	Α	С	С	 old
Person 4	С	Т	G	G	Т	Т	 young



## Advantages of SNP Data

Person's SNP pattern does not change with time or disease, so it can give more insight into susceptibility

Easier to collect samples (can simply use blood rather than affected tissue)



















## **Peak Picking**

# Want to pick peaks that are statistically significant from the noise signal

- Fortunately, data from Duke had peaks picked from spectra already
- Page Group working on a peak-picking algorithm
- Want sensitivity to peaks, while filtering out peaks tdue to noise



### **Metabolomics**

Measures concentration of each lowmolecular weight molecule in sample

- These typically are "metabolites," or small molecules produced or consumed by reactions in biochemical pathways
- These reactions typically catalyzed by proteins (specifically, enzymes)

## Lipomics

Analogous to metabolomics, but measuring concentrations of <u>lipids</u> rather than metabolites

Potentially help induce biochemical pathway information or to help disease diagnosis or treatment choice

101

## **Final Wrapup**

- Molecular biology collecting lots and lots of data in post-genome era
- Opportunity to "connect" molecular-level information to diseases and treatment
- Need analysis tools to interpret
- Machine learning opportunities abound
- Hopefully this tutorial provided solid start toward applying ML to biological data

## Some Additional Readings

Molla, Waddell, Page & Shavlik, Using Machine Learning to Design and Interpret Gene-Expression Microarrays (to appear in the *AI Magazine* special issue on Bioinformatics)

Special issue of *Machine Learning* journal (Volume 52:1/2, 2003) on Machine Learning in the Genomics Era

103

## Thanks To

Mark Craven
Michael Molla
Michael Waddell
Sean McIlwain
Irene Ong
Roland Green
John Tobler

#### Some Useful Datasets

#### **Brief Description**

www.ebi.ac.uk/arrayexpress/

www.ncbi.nlm.nih.gov/geo/

genome-www5.stanford.edu/MicroArray/SMD/

rana.lbl.gov/EisenData.htm

www.genome.wisc.edu/functional/microarray.htm

llmpp.nih.gov/lymphoma/data.shtml

llmpp.nih.gov/DLBCL/

www.rii.com/publications/2002/vantveer.htm

www-genome.wi.mit.edu/cgi-bin/cancer/datasets.cgi

lambertlab.uams.edu/publicdata.htm

www.cs.wisc.edu/~dpage/kddcup2001/

clinicalproteomics.steem.com/

snp.cshl.org/

EBI microarray data repository NCBI microarray data repository Stanford microarray database Eisen-lab's yeast data, (Spellman et al. 1998) University of Wisconsin E. coli Genome Project Diffuse large B-cell lymphoma (Alizadeh et al. 2000) Molecular profiling (Rosenwald et al. 2002) Breast cancer prognosis (Van't Veer et al. 2002) **MIT Whitehead Center for Genome** Research, including data in Golub et al. (1999) Lambert Laboratory data for multiple myeloma

KDD Cup 2001 data; Task 2 includes correlations in genes' expression levels

Proteomics data (mass spectrometry of proteins)

Single nucleotide polymorphism (SNP) data

## Bibliography from AI Mag Article

- Alizadeh, A.; Eisen, M.; Davis, R.; Ma, C.; Lossos, I., Rosenwald, A.; Boldrick, J.; Hajeer, S.;Tran, T.; Yu, X.; Powell, J.; Yang, L.; Marti, G.; Moore, T.; Hudson, J. Jr; Lu, L.; Lewis, D.; Tibshirani, R.; Sherlock, G; Chan, W.; Greiner, T.; Weisenburger, D.; Armitage, J.; Warnke, R.; Levy, R.; Wyndham Wilson, W.; Grever, M.; Byrd, J.; Botstein, D.; Brown, P.; and Staudt, L. 2000. Distinct Types of Diffuse Large B-cell Lymphoma Identified by Gene Expression Profiling. *Nature* 403:503-511. ۲
- Bairoch, A. and Apweiler, R. 2000. The SWISS-PROT Protein Sequence Database and its Supplement TrEMBL in 2000. Nucleic Acids Research 28:45-48
- Breslauer, K.; Frank, R.; Blocker, H.; and Marky, L. 1986. Predicting DNA Duplex Stability from the Base Sequence. *Proceedings* of the National Academy of Science USA 83:3746-3750.
- Brown, M.; Grundy, W.; Lin, D.; Cristianini, N.; Sugnet, C.; Furey, T.; Ares M. Jr.; and Haussler, D. 2000. Knowledge-based Analysis of Microarray Gene Expression Data by using Support Vector Machines. *Proceedings of the National Academy of Science USA* 97(1):262-267.
- Cheng, J.; Hatzis, C.; Hayashi, H.; Krogel, M.; Morishita, S.; Page, D. and Sese, J. 2002. Report on KDD Cup 2001. SIGKDD Explorations 3(2):47-64. ۲
- Craven, M.; Page, D.; Shavlik, J.; Bockhorst J.; and Glasner J. 2000. Using Multiple Levels of Learning and Diverse Evidence Sources to Uncover Coordinately Controlled Genes. *Proceedings of the 17th International Conference on Machine Learning*, Morgan Kaufmann, Palo Alto, CA. ۲
- Notgair Kaumanin, Paio Alio, CA.
  Davidson, E.; Rast, J.; Oliveri, P.; Ransik, A.; Calestani, C.; Yuh, C.; Amore, G.; Minokawa, T.; Hynman, V.; Arenas-Mena, C.; Otim, O.; Brown, C.; Livi, C.; Lee, P.; Revilla, R.; Alistair R.; Pan Z.; Schilstra M.; Clarke, P.; Arnone, M.; Rowen, L.; Cameron, R.; McClay, D.; Hood, L. and Bolouri, H. 2002. A Genomic Regulatory Network for Development. *Science* 295:1669-1678.
  Eisen M.; Spellman P.; Brown P.; and Botstein D. 1998. Cluster Analysis and Display of Genome-Wide Expression Patterns. *Proceedings of the National Academy of Science* USA 95:14863-14863.
- Friedman, N. and Halpern J. 1999. Modeling Beliefs in Dynamic Systems. Part II: Revision and Update. *Journal of AI Research* 10:117-167.
- Golub T.; Slonim D.; Tamayo, P.; Huard, C.; Gaasenbeek, M.; Mesirov, J.; Coller, H.; Loh, M.; Downing, J.; Caligiuri, M.; Bloomfield, C; and Lander, E. 1999. Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring. *Science* 286:531-537.
- Hanisch, D.; Zien, A.; Zimmer, R.; and Lengauer, T. 2002. Co-Clustering of Biological Networks and Gene Expression Data. *Bioinformatics* 18:S145-S1554. Hood, L. and Galas, D. 2003. The Digital Code of DNA. *Nature* 421:444-448. ۲

- Hunter, L. 2003. An Introduction to Molecular Biology for Computer Scientists. *AI Magazine*, this issue. Khodursky, A.: Peter, B.: Cozzarelli, N.: Botstein, D.: Brown, P. and Yanofsky, C. 2000. DNA Microarray Analysis of Gene Expression in Response to Physiological and Genetic Changes that Affect Tryptophan in *Escheria Coll. Proceedings of the National Academy of Science USA* 97:12170-12175.
- Lazarou, J.; Pomeranz, B. and Corey, P. 1998. Incidence of Adverse Drug Reactions in Hospitalized Patients. *Journal of the American Medical Association* 279(15):1200-1205.

## More Bibliography

- Li, C. and Wong, W. 2001. Model-based Analysis of Oligonucleotide Arrays: Expression Index Computation and Outlier Detection. *Proceedings of the National Academy of Science USA* 98(1):31-36. Mancinelli, L.; Cronin, M. and Sadee W. 2000. Pharmacogenomics: The Promise of Personalized Medicine. *AAPS PharmSci* 2(1):
- ۲ article 4.
- Molla, M; Andrae, P; Glasner, J; Blattner, F. and Shavlik, J. 2002. Interpreting Microarray Expression Data Using Text Annotating the Genes. *Information Sciences* 146:75-88. Mitchell, T. 1997. *Machine Learning*. McGraw-Hill, Boston, MA. Oliver, S.; Winson, M.; Kell, D. and Baganz, F. 1998. Systematic Functional Analysis of the Yeast Genome. *Trends in Biotechnology* 16(9):373-378. ۲
- ۲
- ۲
- ۲ Ong, I.; Glassner, J. and Page, D. 2002. Modelling Regulatory Pathways in *E.coli* from Time Series Expression Profiles. *Bioinformatics* 18:241S-248S.
- ۲
- Doublinding 16:2413-2463.
   Newton, M.; Kendziorski C.; Richmond, C.; Blattner, F. and Tsui, K. 2001. On Differential Variability of Expression Ratios: Improving Statistical Inference about Gene Expression Changes from Microarray Data. *Journal of Computational Biology* 8:37-52.
   Nuwaysir, E. F.; Huang, W.; Albert, T.; Singh, J.; Nuwaysir, K.; Pitas, A.; Richmond, T.; Gorski, T.; Berg, J.; Ballin, J.; McCormick, M.; Norton, J.; Pollock, T.; Sumwalt, T.; Butcher, L.; Porter, D.; Molla, M.; Hall, C.; Blattner, F.; Sussman, M.; Wallace, R.; Cerrina, F. and Green, R. 2002. Gene Expression Analysis Using Oligonucleotide Arrays Produced by Maskless Lithography. *Genome Research* 12(11):1749-1755. ۲
- ۲ Pe'er, D.; Regev, A.; Elidan, G. and Friedman, N. 2001. Inferring Subnetworks from Perturbed Expression Profiles. *Bioinformatics* 17:S215-S224
- Rosenwald, A.; Wright, G.; Chan, W.; Connors, J.; Campo, E.; Fisher, R.; Gascoyne, R.; Muller-Hermelink, H.; Smeland, E. and Staudt, L. 2002. The Use of Molecular Profiling to Predict Survival after Chemotherapy for Diffuse Large-B-Cell Lymphoma. *New England Journal of Medicine* 346(25):1937-1947. ۲
- Segal, E.; Taskar, B.; Gasch, A.; Friedman, N. and Koller, D. 2001. Rich Probabilistic Models for Gene Expression. *Bioinformatics* 1(1):1-10. ۲ ۲
- ۲
- Shrager, J.; Langley, P.; and Pohorille, A. 2002. Guiding Revision of Regulatory Models with Expression Data. *Proceedings of the Pacific Symposium on Biocomputing*, 486-497, World Scientific, Lihue, Hawaii. Spellman, P.; Sherlock, G.; Zhang, M.; Iyer, V.; Anders, K.; Elsen, M.; Brown, P.; Botstein, D. and Futcher, B. 1998. Comprehensive Identification of Cell Cycle-Regulated Genes of the Yeast *Saccharomyces Cerevisiae* by Microarray Hybridization. *Molecular Biology of the Cell* 9:3273-3297.
- Thomas, R.; Rank, D.; Penn, S.; Zastrow, G.; Hayes, K.; Pande, K.; Glover, E.; Silander, T.; Craven, M.; Reddy, J.; Jovanovich, S. and Bradfield, C. 2001. Identification of Toxicologically Predictive Gene Sets using cDNA Microarrays. *Molecular Pharmacology* 60:1189-1194.
- Doler J; Wolla M.; Nuwaysir, E.; Green R. and Shavlik J. 2002. Evaluating Machine Learning Approaches for Aiding Probe Selection for Gene-Expression Arrays. *Bioinformatics*, 18:S164-S171.
   Van 't Veer, L.; Dai, H.; van de Vijver, M.; He, Y.; Hart, A.; Mao, M.; Peterse, H.; van der Kooy, K.; Marton, M.; Witteveen, A.; Schreiber, G.; Kerkhoven, R.; Roberts, C.; Linsley, P.; Bernords, R. and Friend, S. 2002. Gene Expression Profiling Predicts Clinical Outcome of Breast Cancer. *Nature* 415:530-536.

## Additional Citations

- Akutsu, T.; Kuhara, S.; Maruyama, O. and Miyano, S. 1998. Identification of Gene Regulatory Networks by Strategic Gene Disruptions and Gene Overexpressions. ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 695-702
- Chrisman, L.; Langley, P.; Bay, S. and Pohorille, A. 2003. Incorporating Biological Knowledge into Evaluation of Causal Regulatory Hypotheses. *Pacific Symposium on* Biocomputing, pp. 128-139.
- Ideker, T.; Thorsson, V. and Karp, R. 2000. Discovery of Regulatory Interactions Through Perturbation: Inference and Experimental Design. Pacific Symposium on Biocomputing, pp. 302-313.
- Segal, E.; Taskar, B.; Gasch, A.; Friedman, N. and Koller, D. 2002. Rich Probabilistic Models for Gene Expression. *Proc. Ninth International Conference on Intelligent Systems for Molecular Biology (ISMB), Bioinformatics*, 17 (Suppl 1), pp. 243--252.
- Shamir, R. and Sharan, R. 2000. CLICK: A Clustering Algorithm with Applications to Gene Expression Analysis. *Currents in Computational Molecular Biology*, pages 6--7, S. Miyano, R. Shamir and T. Takagi (editors) Universal Academy Press, 2000). Proc. ISMB '00, pp., 307--316, AAAI Press, Menlo Park, CA.
- Tanay, A. and Shamir, R. 2001. Computational Expansion of Genetic Networks. Proc. Ninth International Conference on Intelligent Systems for Molecular Biology (ISMB), pp. 270-278.