

# Gene Functional Classification from Heterogeneous Data (2001, P.Pavlidis, et.al...)

Presenter: James J. Winkle

8 May, 2014

# Outline

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

- 1 Overview
- 2 Heterogeneous Data
- 3 SVM and Kernel
- 4 Results
- 5 Conclusions

# Overview of Paper

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

- Biology seeks to understand the “molecular machinery of the cell.”
- A data-centric complementary view of this machinery is provided by the following types of (“heterogeneous”) data:
  - DNA  $\mu$ -array hybridization experiments
  - Genomic Sequences: Phylogenetic Profiles
- This paper hopes to advance computational techniques toward a long-term goal of learning about gene-function from many different types of genomic data.
- Various Kernel combinations are tested to address how best to combine heterogeneous data for genomic classification

# Genomic Sequencing Cost

Heterogeneous  
Data Gene  
Classification

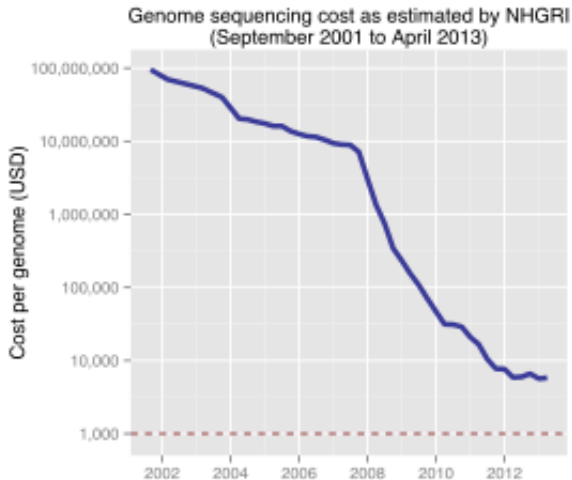
Overview

Heterogeneous  
Data

SVM

Results

Conclusions



# Illumina Stock Price (ILMN)

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions



# Heterogeneous Data Sources

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

The paper cites previous work from:

- Brown *et al.*(2000): applied SVM techniques to yeast expression data with “excellent classification performance”
- Combining heterogeneous data sets is mentioned (Marcotte, Pellegrini, . . . 1999) but with data sets considered separately rather than at once.

This paper asserts that: “the performance of SVM’s when data types are combined and a single hypothesis is formed is superior to combining two independent hypotheses.”

# Data Type Definitions

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

## DNA $\mu$ -expression data

- “The first data set derives from a collection of **DNA  $\mu$ -array hybridization** experiments. Each data point represents the logarithm of the ratio of expression levels of a particular gene under two different experimental conditions.”

$$X_i = \frac{\log(E_i/R_i)}{[\sum_{j=1}^{79} \log^2(E_j/R_j)]^{\frac{1}{2}}}$$

- A snapshot of the messenger RNA expression levels during various time points of “cell events” (diauxic shift, cell division, sporulation, “shocks”)
- If two genes have a functional link, they should be expressed together during the functional event

# Data Type Definitions

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

## Phylogenetic Profiles

- Two genes with similar phylogenetic profiles are *likely* to have similar functions, under the assumption that their similar pattern of inheritance across species is the result of a functional link.
- “In its simplest form, a **phylogenetic profile** is a bit string, in which the Boolean value of each bit indicates a close homolog in the genome.” In this paper, each genome position in the data vector is  $-\log E_{val}$  from BLAST in a search against the complete genome (negative values truncated to 0).



# Assigning protein functions by comparative genome analysis protein phylogenetic profiles

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

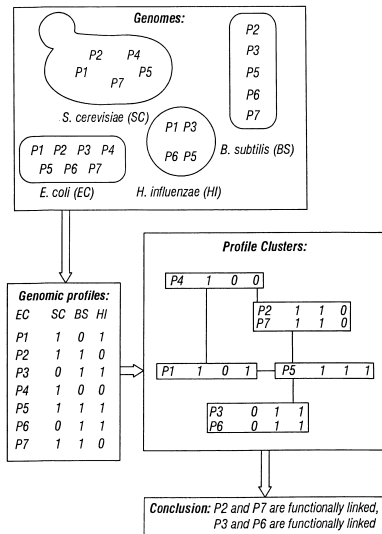


FIG. 3

US 6564151 B1

Inventors: M. Pellegrini, ...

# Heterogeneous Data Sources

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

- Classification comes from the CYGD (MIPS Comprehensive Yeast Genome Database), which contains “several hundred functional classes”
- Classes containing 10 or more genes are selected
- 108 Classes are initially selected (but later narrowed to 27 “learnable classes”)
- The two genomic data vectors are of length:

$$\mathbf{x}_g = [\mu\text{-array expression data}] \quad n = 79$$

$$\mathbf{x}_p = [\text{phlyo}] \quad n = 24$$

- There are  $N = 2465$  yeast genes used as the data set (selected for “accurate functional annotations”)

# Polynomial Kernel

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

- The Kernel function selected is polynomial degree 3 (data vector is projected on the unit sphere)

$$K(\mathbf{x}, \mathbf{y}) = \left( \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|} + 1 \right)^3$$

- The polynomial “takes into account pairwise and tertiary correlations...”
- For a fixed vector  $\mathbf{x}$ , the *level sets* of  $K$  are radial in  $\mathbf{y}$
- This would clearly present difficulty for radially symmetric (about  $\mathbf{0}$ ) classes

- Each class is trained with one-against-others (binary) SVMs
- The two types of data are integrated in 3 different ways:
  - 1 Early: concatenate the vectors
  - 2 **Intermediate**: add kernel values for each separately
  - 3 Late: one SVM for each type of data
- The Intermediate integration can be expressed as a new (“Heterogeneous”) Kernel:

$$K(\cdot, \cdot) = K(\mathbf{x}_g, \mathbf{y}_g) + K(\mathbf{x}_p, \mathbf{y}_p)$$

# Integration

Heterogeneous  
Data Gene  
Classification

Overview

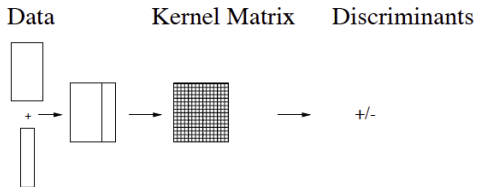
Heterogeneous  
Data

SVM

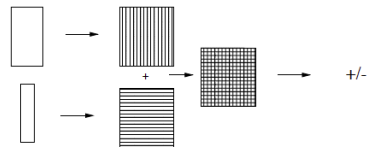
Results

Conclusions

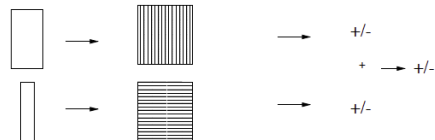
Early  
integration



Intermediate  
integration



Late  
integration



# Restricting Correlations

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

Argument for the use of the Intermediate Integration:

- The heterogeneous kernel creates *local* features of polynomial relationships of one type of data only
- The local features are combined linearly
- Thus, polynomial relationships between different *types* of data are ignored
- Removal of these correlations reduces overfitting

# Validation

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

- For most of the experiments, 3-fold cross-validation is used
- A method-cost is used to evaluate the performance of a method  $M$  (early, intermediate, late)

$$C(M) = (f_p(M) + 2 \cdot f_n(M))/n$$

- False negatives  $f_n$  are given more weight than false positives  $f_p$  ( $n$  = number in class)
- Failing to recognize a limited class member is worse than recognizing a non-member
- The method-cost is normalized to  $[0,1]$ , with 1 being a perfect classifier, as follows:  $S(M) = (C(N) - C(M))/2$

# Example Normalization Calculation

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

- For a method that classifies perfectly, we have  $C(M) = 0$  ( $f_p = f_n = 0$ )
- The cost of classifying all data as negative is:

$$C(N) = (0 + 2 \cdot n)/n = 2$$

- Thus, a perfect classifier is normalized to  $(2 - 0)/2 = 1$
- ...and a null classifier is normalized to  $(2 - 2)/2 = 0$

*\*the formula shown in the paper does not work as written and is fixed here (it varies from the previously published Brown paper also).*



# Results

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

Class	Exp	Phylo	Early	Intermediate	Late
amino acid transporters	0.05 ± 0.04	<b>0.77 ± 0.10</b>	0.50 ± 0.04	<b>0.71 ± 0.08</b>	0.49 ± 0.07
ribosomal proteins	0.71 ± 0.02	0.09 ± 0.03	<b>0.76 ± 0.01</b>	0.71 ± 0.01	0.69 ± 0.01
sugar and carbohydrate transporters	0.33 ± 0.07	<b>0.67 ± 0.02</b>	<b>0.68 ± 0.06</b>	<b>0.70 ± 0.01</b>	0.63 ± 0.03
glycolysis and gluconeogenesis	0.21 ± 0.03	<b>0.43 ± 0.05</b>	0.28 ± 0.02	<b>0.39 ± 0.05</b>	<b>0.39 ± 0.04</b>
mitochondrial organization	<b>0.40 ± 0.03</b>	0.15 ± 0.01	<b>0.43 ± 0.03</b>	<b>0.42 ± 0.02</b>	0.35 ± 0.02
tricarboxylic acid pathway	0.21 ± 0.11	0.15 ± 0.07	<b>0.32 ± 0.08</b>	<b>0.42 ± 0.07</b>	<b>0.25 ± 0.13</b>
deoxyribonucleotide metabolism	0.07 ± 0.05	<b>0.31 ± 0.11</b>	<b>0.24 ± 0.15</b>	<b>0.39 ± 0.11</b>	<b>0.31 ± 0.12</b>
organization of cytoplasm	0.35 ± 0.01	0.18 ± 0.01	<b>0.38 ± 0.01</b>	0.34 ± 0.02	<b>0.35 ± 0.02</b>
transport ATPases	0.13 ± 0.04	<b>0.37 ± 0.05</b>	0.23 ± 0.05	<b>0.32 ± 0.04</b>	0.22 ± 0.03
amino acid biosynthesis	0.18 ± 0.02	0.28 ± 0.02	0.29 ± 0.03	<b>0.36 ± 0.04</b>	0.27 ± 0.02
purine ribonucleotide metabolism	0.17 ± 0.03	<b>0.26 ± 0.05</b>	0.20 ± 0.04	<b>0.33 ± 0.04</b>	0.19 ± 0.03
pyrimidine ribonucleotide metabolism	0.03 ± 0.02	<b>0.33 ± 0.06</b>	0.11 ± 0.04	<b>0.28 ± 0.03</b>	0.17 ± 0.03
cytoplasmic degradation	<b>0.32 ± 0.01</b>		<b>0.32 ± 0.06</b>	<b>0.30 ± 0.03</b>	0.17 ± 0.02
respiration	<b>0.32 ± 0.02</b>		<b>0.30 ± 0.04</b>	0.23 ± 0.04	0.17 ± 0.03
organization of chromosome structure	<b>0.31 ± 0.01</b>		<b>0.30 ± 0.01</b>	<b>0.29 ± 0.02</b>	0.13 ± 0.03
phosphate utilization	<b>0.22 ± 0.04</b>	0.08 ± 0.05	<b>0.26 ± 0.05</b>	<b>0.21 ± 0.04</b>	<b>0.22 ± 0.04</b>
organization of plasma membrane	0.07 ± 0.02	<b>0.25 ± 0.01</b>	<b>0.24 ± 0.03</b>	<b>0.26 ± 0.03</b>	<b>0.26 ± 0.02</b>
pentose phosphate pathway		<b>0.20 ± 0.15</b>		<b>0.26 ± 0.07</b>	<b>0.15 ± 0.10</b>
cellular import	0.04 ± 0.02	<b>0.25 ± 0.04</b>	<b>0.18 ± 0.05</b>	0.17 ± 0.03	<b>0.21 ± 0.04</b>
protein folding and stabilization		<b>0.24 ± 0.04</b>	<b>0.20 ± 0.04</b>	<b>0.23 ± 0.05</b>	0.14 ± 0.04
proteolysis	<b>0.23 ± 0.02</b>		<b>0.24 ± 0.02</b>	<b>0.18 ± 0.06</b>	0.17 ± 0.01
pheromone response generation	<b>0.24 ± 0.05</b>		0.15 ± 0.03	<b>0.14 ± 0.08</b>	
nuclear organization	<b>0.21 ± 0.01</b>	0.07 ± 0.01	<b>0.24 ± 0.03</b>	<b>0.24 ± 0.02</b>	0.17 ± 0.02
drug transporters		<b>0.23 ± 0.09</b>			
organization of endoplasmic reticulum	<b>0.20 ± 0.02</b>		<b>0.22 ± 0.03</b>	<b>0.19 ± 0.05</b>	0.13 ± 0.03
organization of cell wall	<b>0.12 ± 0.04</b>	<b>0.19 ± 0.06</b>	<b>0.14 ± 0.08</b>	<b>0.16 ± 0.07</b>	<b>0.21 ± 0.08</b>
anion transporters		<b>0.21 ± 0.02</b>			
Mean cost savings	0.19 ± 0.02	0.21 ± 0.04	0.27 ± 0.03	0.31 ± 0.03	0.24 ± 0.03
Number of best-performing	10	12	17	21	8
Number of non-learnable	4	6	3	2	3

# Results

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

Exp	Phylo	Early	Intermediate	Late
$0.05 \pm 0.04$	<b><math>0.77 \pm 0.10</math></b>	$0.50 \pm 0.04$	<b><math>0.71 \pm 0.08</math></b>	$0.49 \pm 0.07$
$0.71 \pm 0.02$	$0.09 \pm 0.03$	<b><math>0.76 \pm 0.01</math></b>	$0.71 \pm 0.01$	$0.69 \pm 0.01$
$0.33 \pm 0.07$	<b><math>0.67 \pm 0.02</math></b>	<b><math>0.68 \pm 0.06</math></b>	<b><math>0.70 \pm 0.01</math></b>	$0.63 \pm 0.03$
$0.21 \pm 0.03$	<b><math>0.43 \pm 0.05</math></b>	$0.28 \pm 0.02$	<b><math>0.39 \pm 0.05</math></b>	<b><math>0.39 \pm 0.04</math></b>
<b><math>0.40 \pm 0.03</math></b>	$0.15 \pm 0.01$	<b><math>0.43 \pm 0.03</math></b>	<b><math>0.42 \pm 0.02</math></b>	$0.35 \pm 0.02$
$0.21 \pm 0.11$	$0.15 \pm 0.07$	<b><math>0.32 \pm 0.08</math></b>	<b><math>0.42 \pm 0.07</math></b>	<b><math>0.25 \pm 0.13</math></b>
$0.07 \pm 0.05$	<b><math>0.31 \pm 0.11</math></b>	<b><math>0.24 \pm 0.15</math></b>	<b><math>0.39 \pm 0.11</math></b>	<b><math>0.31 \pm 0.12</math></b>
$0.35 \pm 0.01$	$0.18 \pm 0.01$	<b><math>0.38 \pm 0.01</math></b>	$0.34 \pm 0.02$	<b><math>0.35 \pm 0.02</math></b>
$0.13 \pm 0.04$	<b><math>0.37 \pm 0.05</math></b>	$0.23 \pm 0.05$	<b><math>0.32 \pm 0.04</math></b>	$0.22 \pm 0.03$
$0.18 \pm 0.02$	$0.28 \pm 0.02$	$0.29 \pm 0.03$	<b><math>0.36 \pm 0.04</math></b>	$0.27 \pm 0.02$
$0.17 \pm 0.03$	<b><math>0.26 \pm 0.05</math></b>	$0.20 \pm 0.04$	<b><math>0.33 \pm 0.04</math></b>	$0.19 \pm 0.03$
$0.03 \pm 0.02$	<b><math>0.33 \pm 0.06</math></b>	$0.11 \pm 0.04$	<b><math>0.28 \pm 0.03</math></b>	$0.17 \pm 0.03$
<b><math>0.32 \pm 0.01</math></b>		<b><math>0.32 \pm 0.06</math></b>	<b><math>0.30 \pm 0.03</math></b>	$0.17 \pm 0.02$
<b><math>0.32 \pm 0.02</math></b>		<b><math>0.30 \pm 0.04</math></b>	$0.23 \pm 0.04$	$0.17 \pm 0.03$
<b><math>0.31 \pm 0.01</math></b>		<b><math>0.30 \pm 0.01</math></b>	<b><math>0.29 \pm 0.02</math></b>	$0.13 \pm 0.03$

# Results

Heterogeneous  
Data Gene  
Classification

<b><math>0.22 \pm 0.04</math></b>	$0.08 \pm 0.05$	<b><math>0.26 \pm 0.05</math></b>	<b><math>0.21 \pm 0.04</math></b>	<b><math>0.22 \pm 0.04</math></b>	
$0.07 \pm 0.02$	<b><math>0.25 \pm 0.01</math></b>	<b><math>0.24 \pm 0.03</math></b>	<b><math>0.26 \pm 0.03</math></b>	<b><math>0.26 \pm 0.02</math></b>	
Overview	<b><math>0.20 \pm 0.15</math></b>		<b><math>0.26 \pm 0.07</math></b>	<b><math>0.15 \pm 0.10</math></b>	
Heterogeneous Data	$0.04 \pm 0.02$	<b><math>0.25 \pm 0.04</math></b>	<b><math>0.18 \pm 0.05</math></b>	$0.17 \pm 0.03$	<b><math>0.21 \pm 0.04</math></b>
SVM		<b><math>0.24 \pm 0.04</math></b>	<b><math>0.20 \pm 0.04</math></b>	<b><math>0.23 \pm 0.05</math></b>	$0.14 \pm 0.04$
Results	<b><math>0.23 \pm 0.02</math></b>		<b><math>0.24 \pm 0.02</math></b>	<b><math>0.18 \pm 0.06</math></b>	$0.17 \pm 0.01$
Conclusions	<b><math>0.24 \pm 0.05</math></b>		$0.15 \pm 0.03$	<b><math>0.14 \pm 0.08</math></b>	
	<b><math>0.21 \pm 0.01</math></b>	$0.07 \pm 0.01$	<b><math>0.24 \pm 0.03</math></b>	<b><math>0.24 \pm 0.02</math></b>	$0.17 \pm 0.02$
		<b><math>0.23 \pm 0.09</math></b>			
	<b><math>0.20 \pm 0.02</math></b>		<b><math>0.22 \pm 0.03</math></b>	<b><math>0.19 \pm 0.05</math></b>	$0.13 \pm 0.03$
	<b><math>0.12 \pm 0.04</math></b>	<b><math>0.19 \pm 0.06</math></b>	<b><math>0.14 \pm 0.08</math></b>	<b><math>0.16 \pm 0.07</math></b>	<b><math>0.21 \pm 0.08</math></b>
		<b><math>0.21 \pm 0.02</math></b>			
<hr/>					
	$0.19 \pm 0.02$	$0.21 \pm 0.04$	$0.27 \pm 0.03$	$0.31 \pm 0.03$	$0.24 \pm 0.03$
	10	12	17	21	8
	4	6	3	2	3

# 5 Most Learnable Classes

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

Class	Size	FP	FN
amino acid transporters	22	$2.0 \pm 0.4$	$5.6 \pm 0.2$
ribosomal proteins	173	$26.6 \pm 1.2$	$34.2 \pm 1.1$
sugar and carbohydrate transporters	32	$2.4 \pm 0.7$	$9.0 \pm 0.0$
deoxyribonucleotide metabolism	9	$0.2 \pm 0.2$	$4.6 \pm 0.7$
mitochondrial organization	296	$84.8 \pm 1.8$	$128.4 \pm 1.7$

# Conclusions

Heterogeneous  
Data Gene  
Classification

Overview

Heterogeneous  
Data

SVM

Results

Conclusions

- SVM's have extended to other data types in this domain (phylogenetic profiles)
- The results of intermediate integration do not show *radical* improvement
- But *some* improvement can be worth a lot
- No analysis of other kernels was made (but claimed no expectation of helping one method over another)
- There is no claim that gene functional ID wants to be perfect (cf. digit recognition); the domain here is to be better (via SVMs).