

Analyzing How Protein Interaction Networks Improve Classification Performance in Gene Expression Data Analysis

Adrin Jalali

Supervised by: Nico Pfeifer

MPI for Informatics

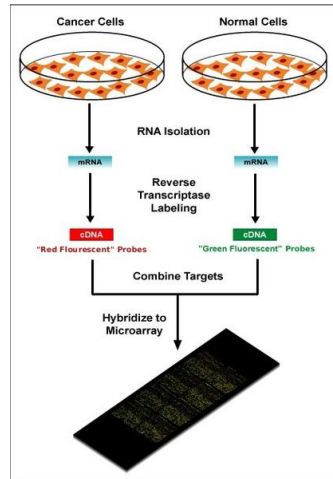
July 9, 2013



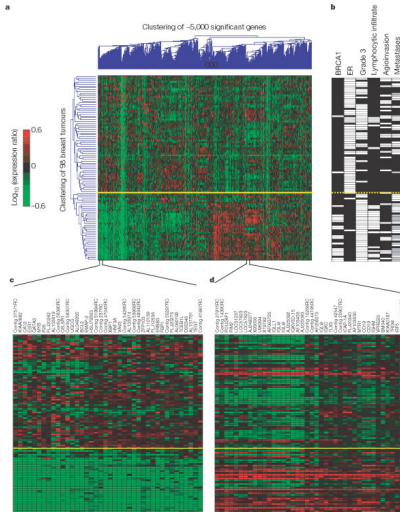
Overview

- ➊ Introduction and data
- ➋ Formulate the approach - incorporate PPI-Network in SVM
- ➌ Results
- ➍ Work in progress - graph kernels

Microarray Gene Expression

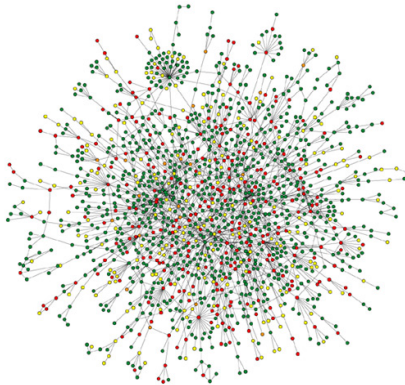


van 't Veer Breast-Cancer Data



- 1 Input: Gene expression data
- 2 Output: Prognosis (Poor vs. Good), Metastases
- 3 Goal: Classify and find important genes
- 4 Issue: Hard to classify due to huge number of features (genes) compared to number of samples ($\sim 22000 \gg 98$)

Yeast Protein Interaction Network



Nature Reviews | Genetics

The colour of a node indicates the phenotypic effect of removing the corresponding protein (red = lethal, green = non-lethal, orange = slow growth, yellow = unknown)

A. Barabasi, Z. Oltvai, Nature Reviews
Genetics, (2004)

JOURNAL OF COMPUTATIONAL BIOLOGY
Volume 19, Number 6, 2012
© Mary Ann Liebert, Inc.
Pp. 694–709
DOI: 10.1089/cmb.2012.0065

Network-Induced Classification Kernels for Gene Expression Profile Analysis

OFER LAVI,^{1,3} GIDEON DROR,² and RON SHAMIR¹

ABSTRACT

Computational classification of gene expression profiles into distinct disease phenotypes has been highly successful to date. Still, robustness, accuracy, and biological interpretation of the results have been limited, and it was suggested that use of protein interaction information jointly with the expression profiles can improve the results. Here, we study three aspects of this problem. First, we show that interactions are indeed relevant by showing that co-expressed genes tend to be closer in the network of interactions. Second, we show that the improved performance of one extant method utilizing expression and interactions is not really due to the biological information in the network, while in another method this is not the case. Finally, we develop a new kernel method—called NICK—that integrates network and expression data for SVM classification, and demonstrate that overall it achieves better results than extant methods while running two orders of magnitude faster.

Key word: algorithms.

- ① It's shown:
 - Co-expressed genes tend to be close in the PPI-Network.
 - Exploit this fact to modify the SVM objective function - called NICK
- ② What can be done:
 - Reverse engineer the learned machine to extract important genes after using the network information.

NICK

1. SVM modified objective function

$$\min_{\mathbf{w}, w_0} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{2} \beta \sum_{(j,k) \in E} (w_j - w_k)^2 \right\}$$

s.t.:

$$\forall i \in \{1, \dots, n\} : (\mathbf{w}\mathbf{x}_i + w_0)y_i \geq 1$$

3. Dual to Primal

$$\mathbf{w} = (\mathbf{I} + \beta \mathbf{B})^{-1} \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

2. Dual Problem

$$\max_{\alpha} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{L})(\mathbf{L}^T \mathbf{x}_j) \right\}$$

$$\mathbf{L}\mathbf{L}^T = (\mathbf{I} + \beta \mathbf{B})^{-1}$$

s.t.:

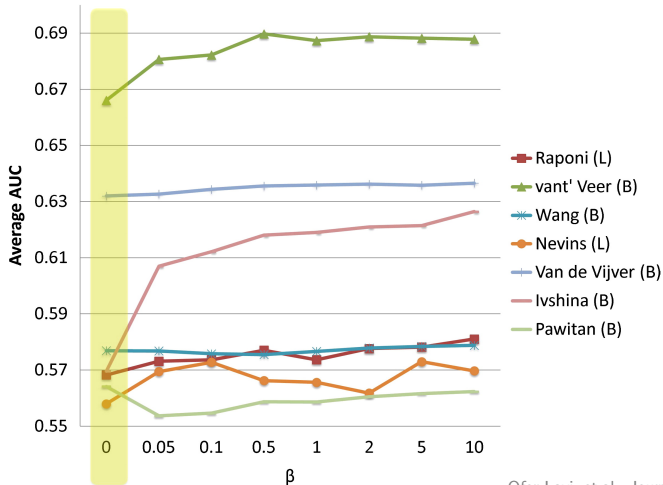
$$\forall i \in \{1, \dots, n\} : \sum_{i=1}^n \alpha_i y_i = 0$$

$$\forall i \in \{1, \dots, n\} : \alpha_i \geq 0$$

Laplacian matrix:

$$\mathbf{B} = \mathbf{D} - \mathbf{A}$$

NICK Performance Summary



Verify What Has To Be Done

- ① To be done: extract important genes.
- ② There is no gold standard for it.
- ③ Synthesized data for the purpose of method verification.

Synthesize Data

- 1 A random graph (PPI-Network)
- 2 Signal nodes (genes):

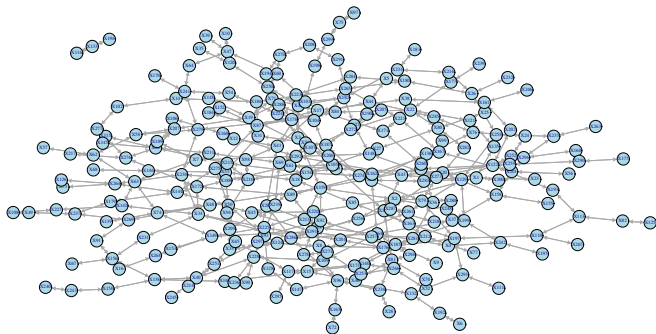
$$f(n) = \begin{cases} N(-\mu, 1) & \text{if } n \text{ is in class 1} \\ N(\mu, 1) & \text{if } n \text{ is in class 2} \end{cases}$$

- 3 Random nodes (non-informative genes):

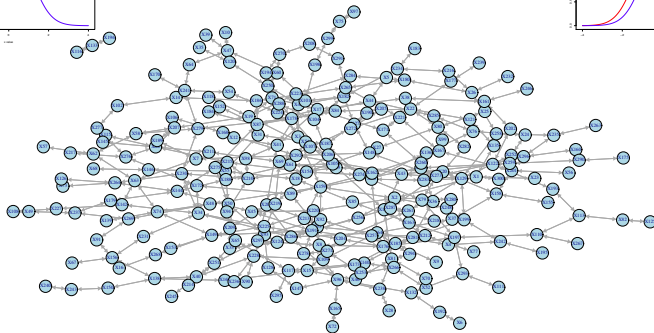
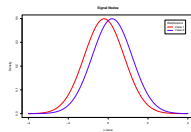
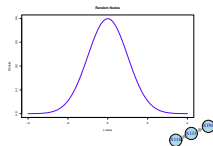
$$f(n) = N(0, 1)$$

- 4 Pathway: 2, 3, or 4 connected signal nodes.

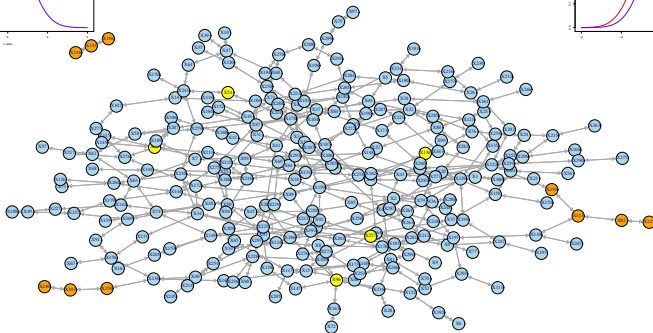
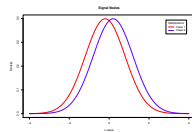
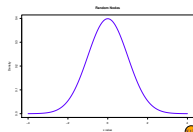
Synthesized Data



Synthesized Data



Synthesized Data



Blue: random gene, Orange: Signal node being a member of a pathway of signal nodes, Yellow: A lonely signal node

Extract Important Genes

- Solve SVM problem for original and transformed data.
- Calculate \mathbf{w} for both models.
- Compute for each pair of nodes, for each model:

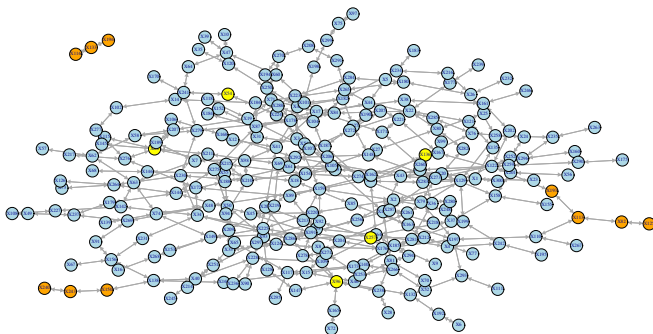
$$Score(i, j) = \frac{|w_i| + |w_j|}{2} \times e^{\max(d_G(i, j), 1)}$$

- Report pairs with highest scores for both trained models.

Synthesized Data Easy Scenario

Original

X196	X196
X53	X53
X233	X233
X39	X39
X88	X88
X196	X133
X116	X116
X127	X127
X197	X197
X127	X148
X148	X148
X150	X150
X148	X273
X116	X133
X160	X160
X96	X96
X95	X95
X273	X273
X88	X115
X40	X40
X53	X8
X53	X164
X195	X195
X56	X56



AUC (Original):	60.6
AUC (Transformed):	62.4
wc p-value (paired):	5.669e-09

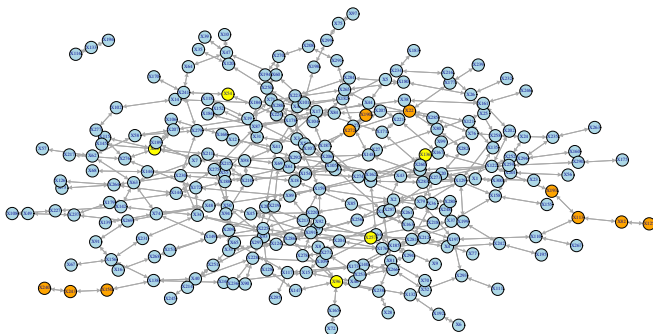
Transformed

X196	X196
X233	X233
X196	X133
X133	X133
X133	X116
X116	X116
X95	X95
X240	X240
X39	X39
X240	X243
X59	X59
X106	X106
X243	X243
X106	X168
X114	X114
X168	X168
X243	X150
X56	X56
X39	X47
X298	X298
X150	X150
X247	X247
X125	X125
X83	X83

Synthesized Data Medium Scenario

Original

X190	X190
X104	X104
X233	X233
X190	X272
X277	X277
X88	X88
X190	X127
X165	X165
X272	X272
X272	X22
X106	X106
X165	X96
X150	X150
X250	X250
X88	X215
X22	X22
X51	X51
X28	X28
X73	X73
X35	X35
X162	X162
X113	X113
X112	X112
X277	X102



AUC (Original):	60.1
AUC (Transformed):	61.5
wc p-value (paired):	1.383e-06

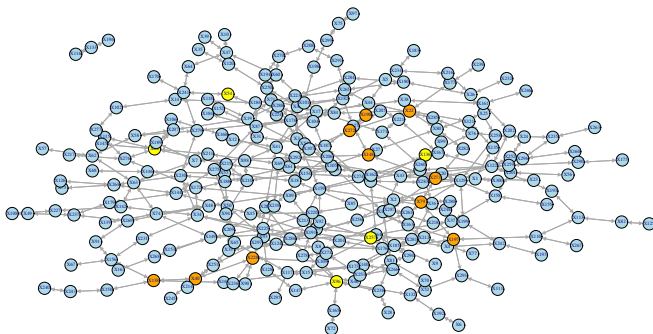
Transformed

X233	X233
X190	X190
X112	X112
X240	X240
X190	X272
X240	X243
X86	X86
X243	X243
X243	X150
X190	X127
X150	X150
X272	X272
X246	X246
X298	X298
X106	X106
X125	X125
X35	X35
X125	X82
X247	X247
X272	X69
X272	X22
X82	X82
X100	X100
X257	X257

Synthesized Data Hard Scenario

Original

X190	X190
X101	X101
X233	X233
X190	X272
X88	X88
X297	X297
X190	X127
X93	X93
X26	X26
X138	X138
X272	X272
X272	X22
X101	X41
X123	X123
X22	X22
X101	X198
X146	X146
X228	X228
X278	X278
X72	X72
X88	X115
X96	X96
X148	X148
X112	X112



wc p-value (paired):	8.151e-13
AUC (Original):	60.2
AUC (Transformed):	62.5

Transformed

X233	X233
X190	X190
X112	X112
X190	X272
X86	X86
X190	X127
X272	X272
X272	X205
X205	X205
X146	X146
X146	X68
X68	X68
X298	X298
X272	X22
X90	X90
X127	X127
X100	X100
X272	X69
X297	X297
X72	X72
X127	X148
X155	X155
X247	X247
X196	X196

van 't Veer

Original

X85453	X85453
X85453	X92140
X6605	X6605
X56886	X56886
X10640	X10640
X8817	X8817
X56894	X56894
X6605	X332
X5733	X5733
X57758	X57758
X7532	X7532
X51	X51
X7566	X7566
X3267	X3267
X89953	X89953
X5713	X5713
X5193	X5193
X5365	X5365
X10874	X10874
X5982	X5982

Transformed

X9917	X9917
X84279	X84279
X197370	X197370
X51143	X51143
X58475	X58475
X55585	X55585
X25949	X25949
X54892	X54892
X126695	X126695
X57168	X57168
X10456	X10456
X148223	X148223
X9742	X9742
X253558	X253558
X342527	X342527
X10175	X10175
X83930	X83930
X57035	X57035
X145482	X145482
X57465	X57465

wc p-value (paired):	0.006
AUC (Original):	72.9
AUC (Transformed):	73.6

van 't Veer

Original

X85453	X85453
X85453	X92140
X6605	X6605
X56886	X56886
X10640	X10640
X8817	X8817
X56894	X56894
X6605	X332
X5733	X5733
X57758	X57758
X7532	X7532
X51	X51
X7566	X7566
X3267	X3267
X89953	X89953
X5713	X5713
X5193	X5193
X5365	X5365
X10874	X10874
X5982	X5982

Node	Degree
X85453	12
X6605	98
X56886	26
X10640	16
X8817	152
X56894	28
X5733	150
X57758	8
X7532	86
X51	172
X7566	16
X3267	56
X89953	4
X5713	126
X5193	32
X5365	70
X10874	132
X5982	172
X92140	20
X332	328

Node	Degree
X9917	0
X84279	0
X197370	0
X51143	0
X58475	0
X55585	0
X25949	0
X54892	0
X126695	0
X57168	0
X10456	0
X148223	0
X9742	0
X253558	0
X342527	0
X10175	0
X83930	0
X57035	0
X145482	0
X57465	0

Transformed

X9917	X9917
X84279	X84279
X197370	X197370
X51143	X51143
X58475	X58475
X55585	X55585
X25949	X25949
X54892	X54892
X126695	X126695
X57168	X57168
X10456	X10456
X148223	X148223
X9742	X9742
X253558	X253558
X342527	X342527
X10175	X10175
X83930	X83930
X57035	X57035
X145482	X145482
X57465	X57465

wc p-value (paired):	0.006
AUC (Original):	72.9
AUC (Transformed):	73.6

Idea

- 1 Estimate gene expression value probability distributions for samples of class A.

Idea

- ① Estimate gene expression value probability distributions for samples of class A.
- ② For each sample in class A and B:
 - ① Extract abnormal genes according to above estimated distributions.
 - ② Extract the part of PPI network induced by extracted genes (almost)

Idea

- ➊ Estimate gene expression value probability distributions for samples of class A.
- ➋ For each sample in class A and B:
 - ➊ Extract abnormal genes according to above estimated distributions.
 - ➋ Extract the part of PPI network induced by extracted genes (almost)
- ➌ Use a graph kernel for labeled graphs to classify extracted graphs.

Graph Kernel

SVM Dual Problem

$$\max_{\alpha} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right\}$$

s.t.:

$$\forall i \in \{1, \dots, n\} : \sum_{i=1}^n \alpha_i y_i = 0$$

$$\forall i \in \{1, \dots, n\} : \alpha_i \geq 0$$

Graph Kernel

SVM Dual Problem

$$\max_{\alpha} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \right\}$$

s.t.:

$$\forall i \in \{1, \dots, n\} : \sum_{i=1}^n \alpha_i y_i = 0$$

$$\forall i \in \{1, \dots, n\} : \alpha_i \geq 0$$

Graph Kernel

SVM Dual Problem

$$\max_{\alpha} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \right\}$$

s.t.:

$$\forall i \in \{1, \dots, n\} : \sum_{i=1}^n \alpha_i y_i = 0$$

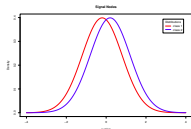
$$\forall i \in \{1, \dots, n\} : \alpha_i \geq 0$$

Graph Kernel

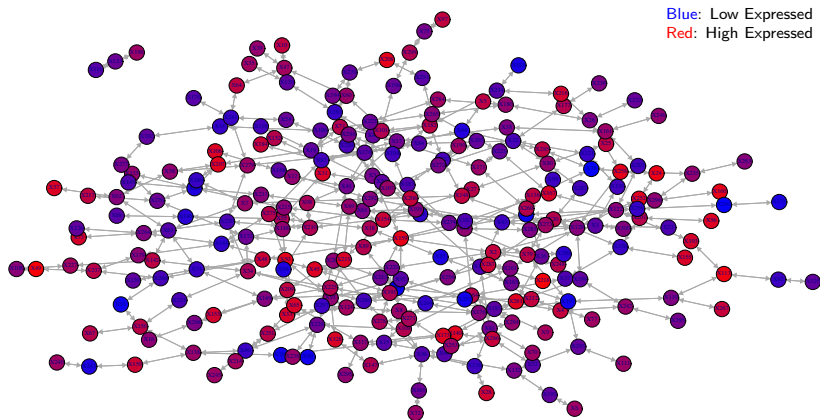
- Random Walks (A^k : number of walks of length k)
- Sub-Tree Kernels
- Shortest Paths Kernels (1-shortest path, k -shortest paths)
- Graphlet Kernels (Isomorphism proved for $n = k + 1, n \leq 11$)
- Laplacian Matrix - Eigenvalues
($\text{Isomorph}(G_1, G_2) \iff \exists P : B(G_1) = P^t B(G_2) P$),
number of connected components = number of $\lambda_i = 0$)

Idea

- 1 Estimate gene expression value probability distributions for samples of class A.
- 2 For each sample in class A and B:
 - 1 Extract abnormal genes according to above estimated distributions.
 - 2 Extract the part of PPI network induced by extracted genes (almost)
- 3 Use a graph kernel for labeled graphs to classify extracted graphs.

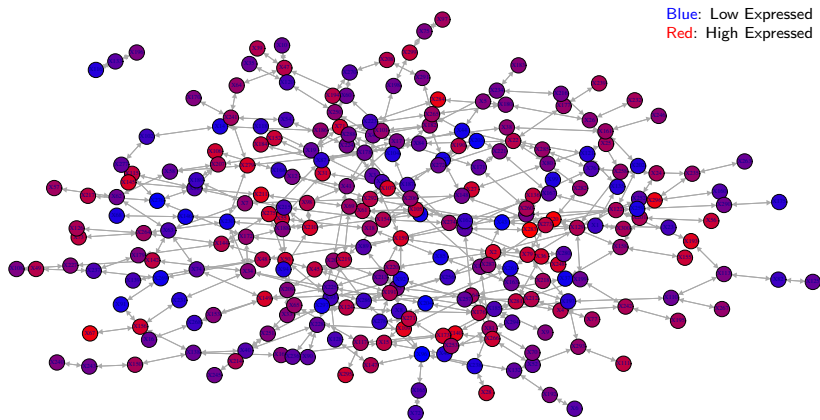


Weighted Idea



Class A member

Weighted Idea



Class B member

Next ...

- Same structured graphs for each sample - even same edge weights.
- Most kernels:
 - Detect structural differences on graphs.
 - Assume nodes are not labeled.
- Design/Find/Change a kernel for our graphs.
- Reverse engineer the kernel on support vectors to detect common substructures on them.

Acknowledgment

- Thomas Lengauer
- Nico Pfeifer
- Sarvesh Nimkubh
- Nora Speicher
- Anna Feldmann

Thank You!
Questions?