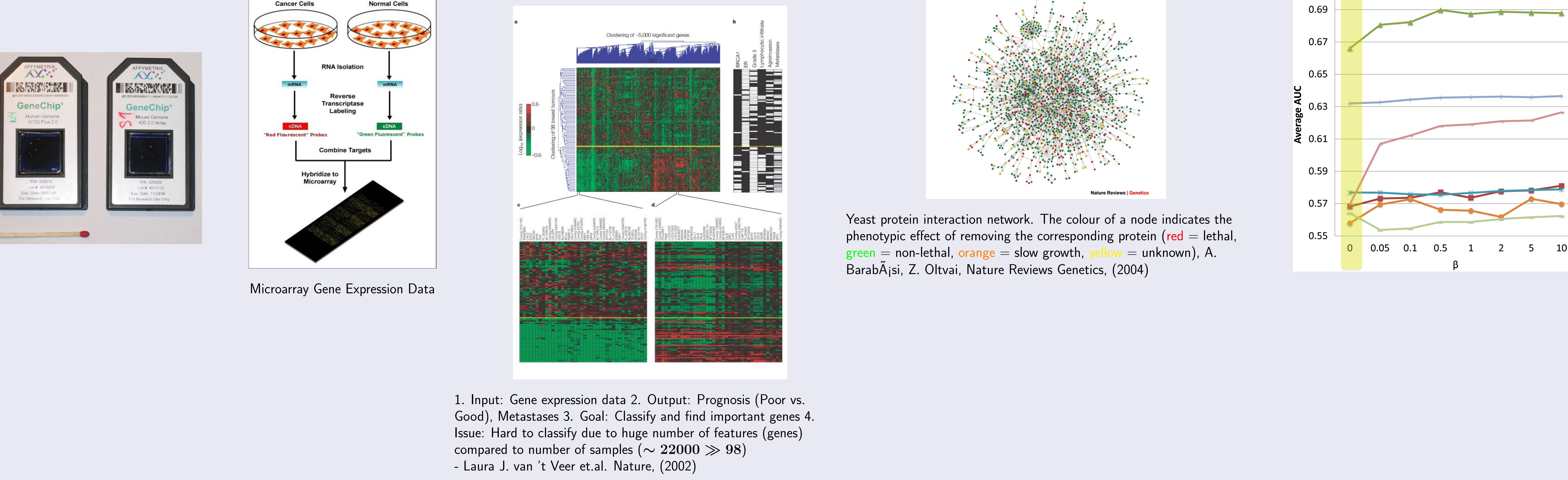# Analyzing How Protein Interaction Networks Improve Classification Performance in Gene Expression Data Analysis

Adrin Jalali

Supervised by: Nico Pfeifer

MPI for Informatics
July 12, 2013

## Data



Microarray Gene Expression Data



Yeast protein interaction network. 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. BarabÃ¡si, Z. Oltvai, Nature Reviews Genetics, (2004)



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 ($\sim 22000 \gg 98$)
- Laura J. van 't Veer et.al. Nature, (2002)

## Method

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

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

s.t.:

$$\forall i \in \{1,\cdots,n\} : (wx_i + w_0)y_i \geq 1$$

3. Dual to Primal

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

2. Dual Problem

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

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

s.t.:

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

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

Laplacian matrix:

$$B = D - A$$

## 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.



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

- Solve SVM problem for original and transformed data.
- Calculate $w$ 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.

## Results



Easy

| Original | | | | Transformed | | | |
|---|---|---|---|---|---|---|---|
| X196 | X196 | X53 | X53 | X196 | X196 | X233 | X233 |
| X233 | X233 | X39 | X39 | X196 | X133 | X133 | X133 |
| X88 | X88 | X196 | X133 | X133 | X116 | X116 | X116 |
| X116 | X116 | X127 | X127 | X95 | X95 | X240 | X240 |
| X197 | X197 | X127 | X148 | X39 | X39 | X240 | X243 |
| X148 | X148 | X150 | X150 | X59 | X59 | X106 | X106 |
| X243 | X273 | X116 | X133 | X243 | X243 | X106 | X168 |

Medium

| Original | | | | Transformed | | | |
|---|---|---|---|---|---|---|---|
| X190 | X190 | X104 | X104 | X233 | X233 | X190 | X190 |
| X233 | X233 | X190 | X272 | X112 | X112 | X240 | X240 |
| X277 | X277 | X88 | X88 | X190 | X272 | X240 | X243 |
| X190 | X127 | X165 | X165 | X86 | X86 | X243 | X243 |
| X272 | X272 | X272 | X22 | X243 | X150 | X190 | X127 |
| X106 | X106 | X165 | X96 | X150 | X150 | X272 | X272 |
| X150 | X150 | X250 | X250 | X246 | X246 | X298 | X298 |

Hard

| Original | | | | Transformed | | | |
|---|---|---|---|---|---|---|---|
| X190 | X190 | X101 | X101 | X233 | X233 | X190 | X190 |
| X233 | X233 | X190 | X272 | X112 | X112 | X190 | X272 |
| X88 | X88 | X297 | X297 | X86 | X86 | X190 | X127 |
| X190 | X127 | X93 | X93 | X272 | X272 | X272 | X205 |
| X26 | X26 | X138 | X138 | X205 | X205 | X146 | X146 |
| X272 | X272 | X272 | X272 | X146 | X68 | X68 | X68 |
| X101 | X41 | X123 | X123 | X298 | X298 | X272 | X22 |

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

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

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