

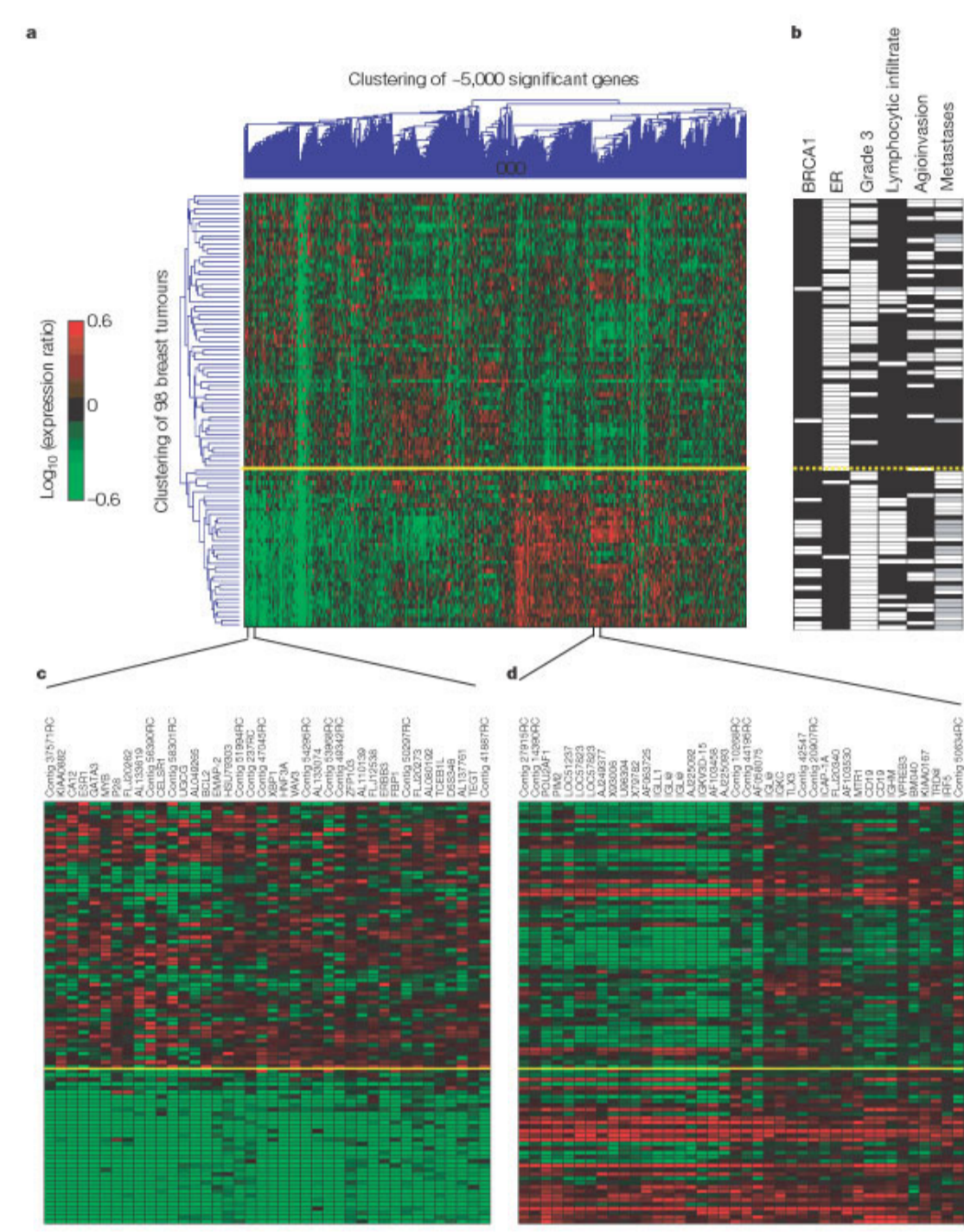
# Analyzing How Protein Interaction Networks Improve Classification Performance in Gene Expression Data Analysis



Adrin Jalali and Nico Pfeifer  
 Department of Computational Biology and Applied Algorithmics  
 Max Planck Institute for Informatics  
 Saarbrücken, Germany



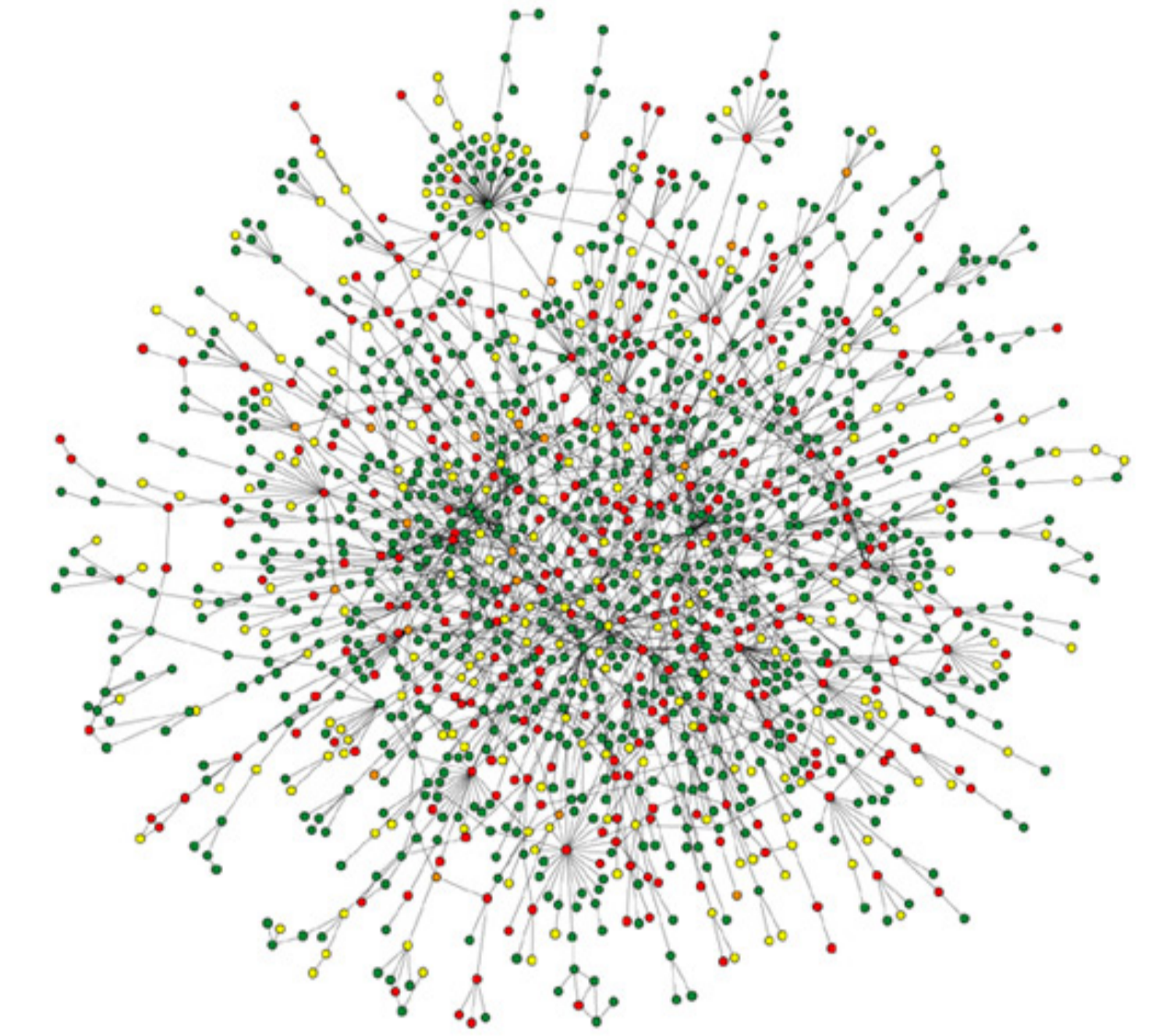
## Motivation



Unsupervised two-dimensional cluster analysis of 98 breast tumours [1]

Problem statement:

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



Nature Reviews | Genetics

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) [2]

## Method

### 1 SVM modified objective function [3]

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

$$\text{s.t.: } \forall i \in \{1, \dots, n\} : (wx_i + w_0)y_i \geq 1$$

### 2 Dual Problem

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

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

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

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

Laplacian matrix:

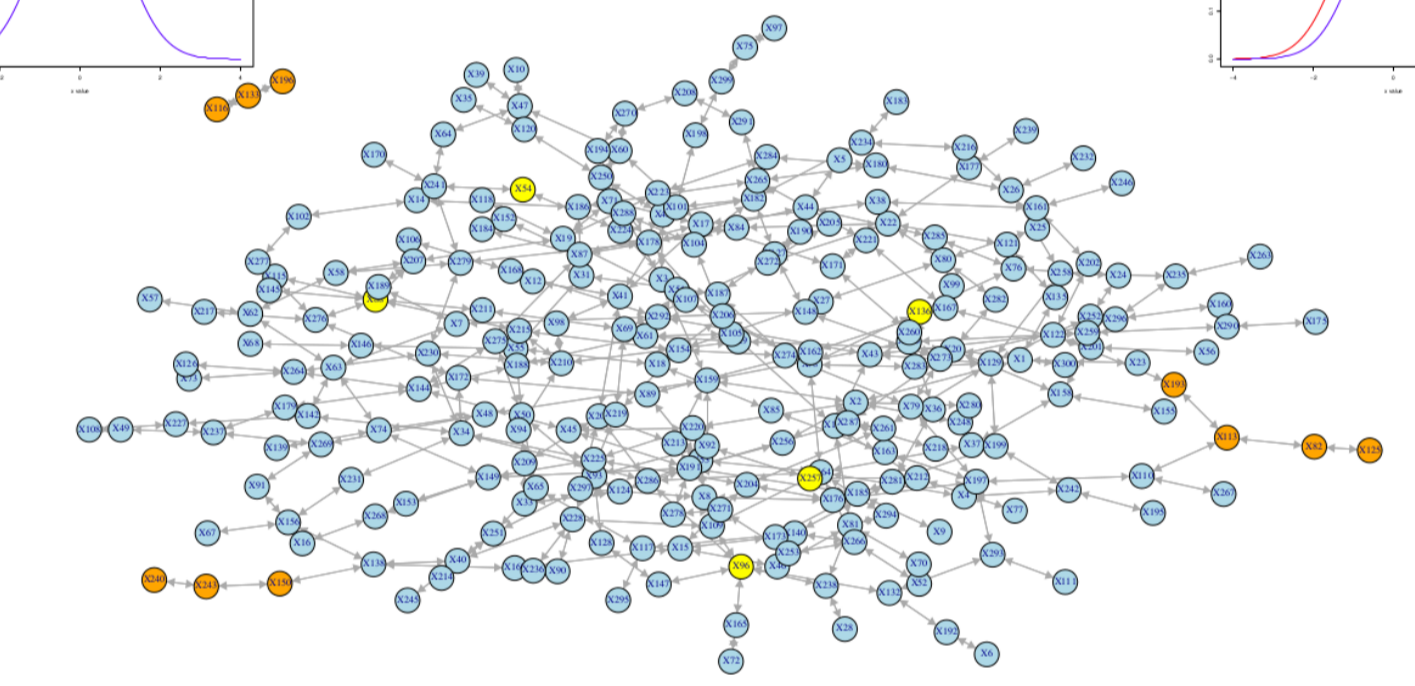
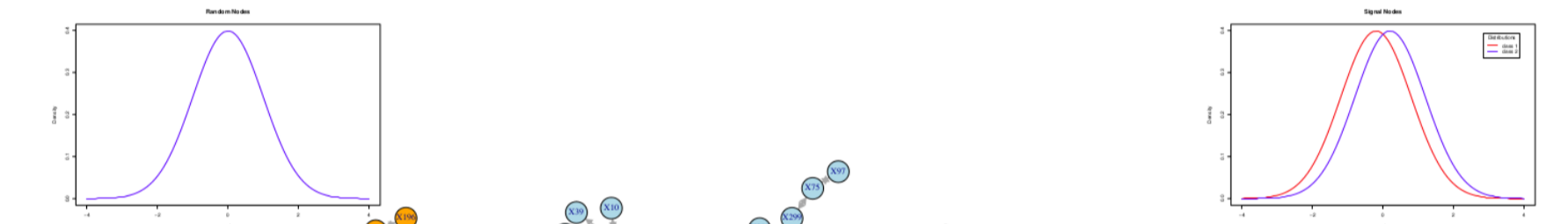
$$B = D - A$$

### 3 Dual to Primal

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

What we do:

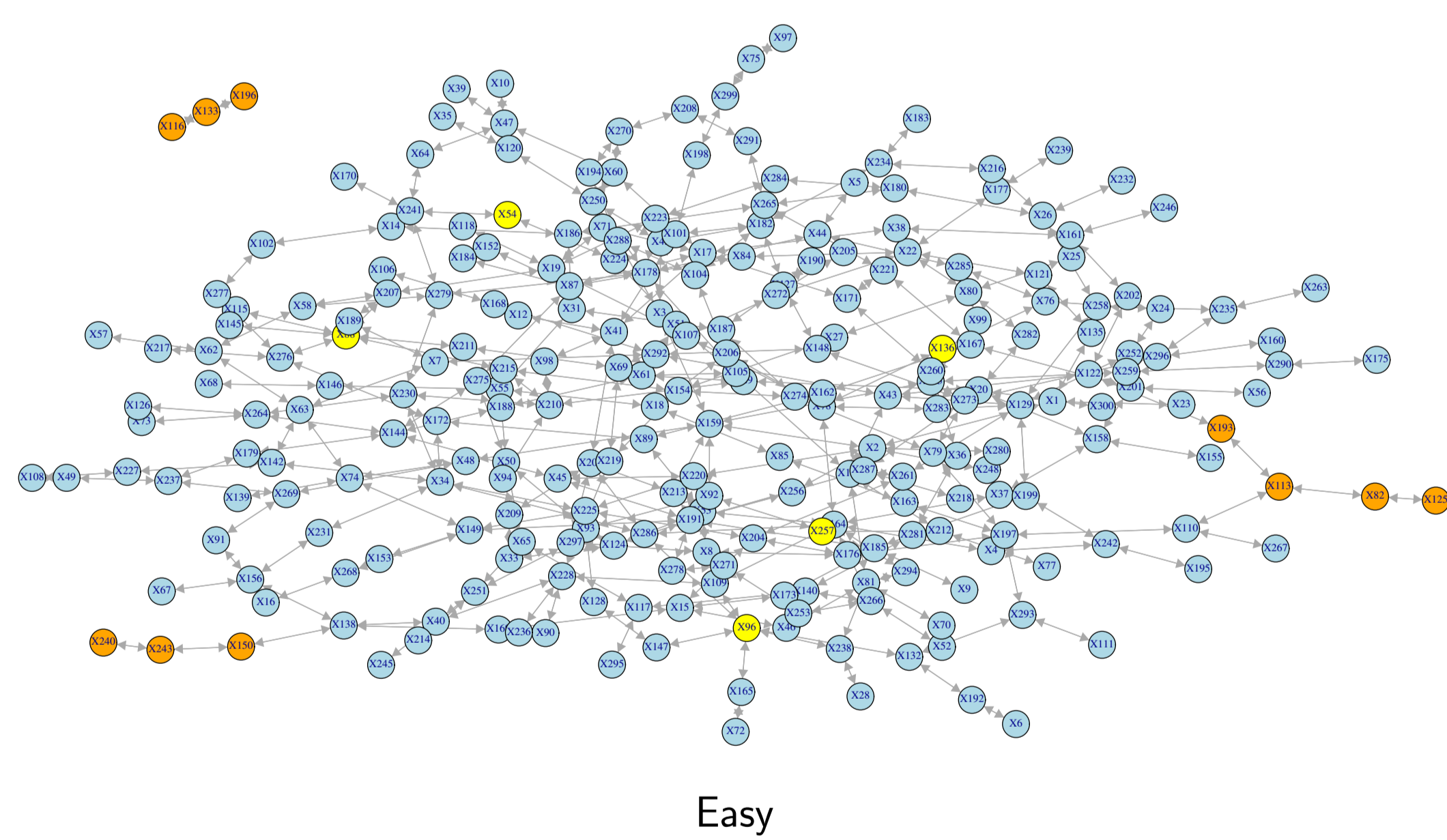
- Reverse engineer the learned machine to extract important genes after using the network information.
- Solve SVM problem for original and transformed data.
- Calculate  $w$  for both models.
- Compute for each pair of nodes, for each model:  
 $Score(i, j) = \frac{|w_i| + |w_j|}{2} \times e^{-max(d_G(i,j), 1)}$
- Report pairs with highest scores for both trained models.



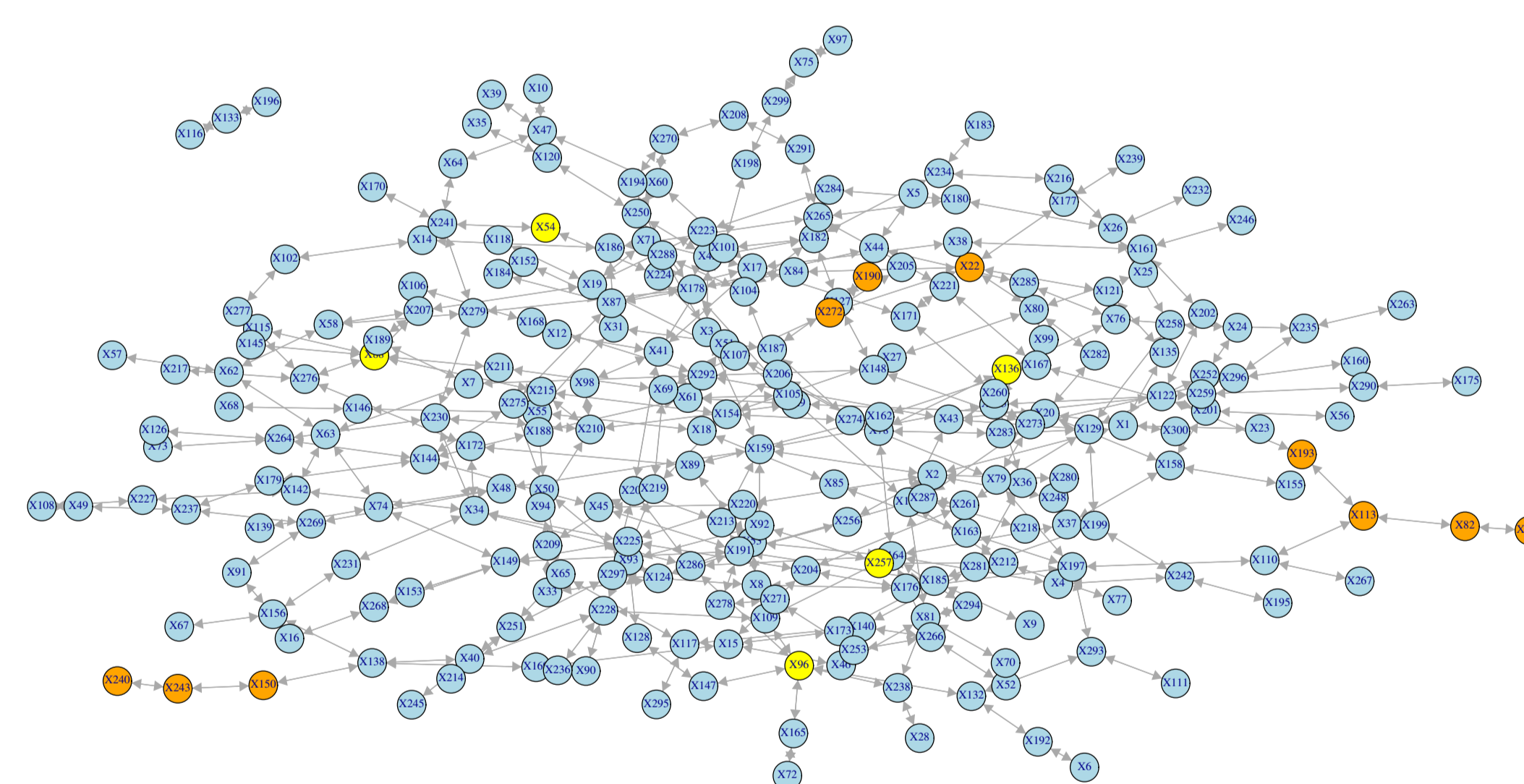
Blue: random gene, Orange: Signal node being a member of a pathway of signal nodes, Yellow: A lonely signal node

- Signal nodes (genes):  $f(n) = \begin{cases} N(-\mu, 1) & \text{if } n \text{ is in class 1} \\ N(\mu, 1) & \text{if } n \text{ is in class 2} \end{cases}$
- Random nodes (non-informative genes):  $f(n) = N(0, 1)$

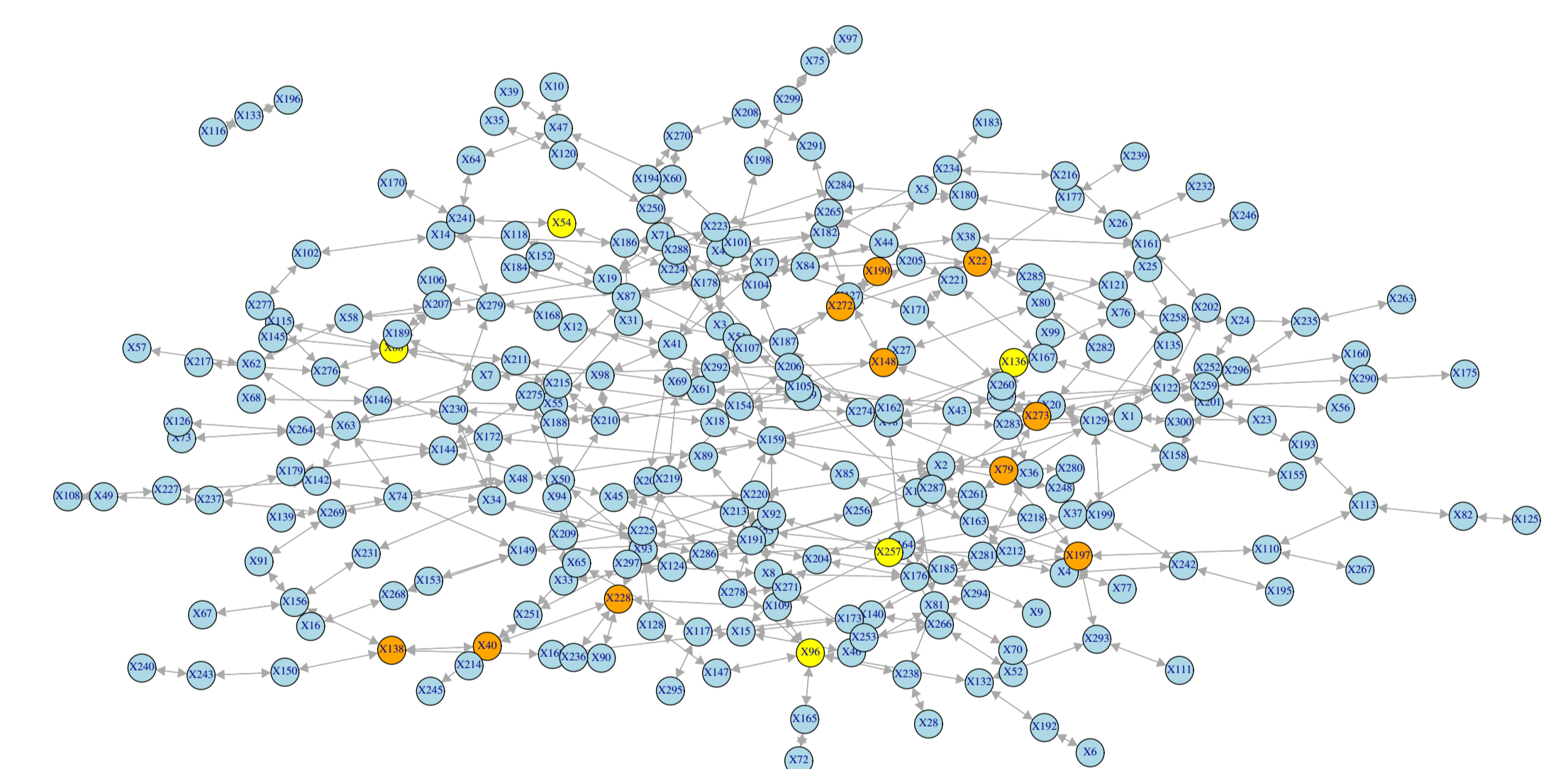
## Results



Easy



Medium



Hard

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

Transformed			
X196	X196	X233	X233
X196	X133	X133	X133
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

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

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

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

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

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

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

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

## References and Acknowledgment

Acknowledgment:

Prof. Dr. Dr. Thomas Lengauer, Sarvesh Nikumbh, Nora Speicher, Anna Feldmann.

References:

1. van't Veer, Laura J., et al. "Gene expression profiling predicts clinical outcome of breast cancer." *nature* 415.6871 (2002): 530-536.
2. Barabási, Albert-László, and Zoltan N. Oltvai. "Network biology: understanding the cell's functional organization." *Nature Reviews Genetics* 5.2 (2004): 101-113.
3. Lavi, Ofer, Gideon Dror, and Ron Shamir. "Network-induced classification kernels for gene expression profile analysis." *Journal of Computational Biology* 19.6 (2012): 694-709.