

# **Knowledge Enhancement through Intra-relation, Inter-relation, and Integration of Different Levels of Genomic Data**

**Hyunjung (Helen) Shin**

Dept. of Industrial & Information Systems Engineering  
Ajou University, Korea

[shin@ajou.ac.kr](mailto:shin@ajou.ac.kr)

<http://www.alphaminers.net>

<http://www.kyb.tuebingen.mpg.de/~shin>

# Outline

---

**[Intra-Relation]**  
**Breast Cancer Survivability  
Prediction**

**Graph Representation**  
**Graph-based Semi-Supervised Learning (SSL)**

**[Integration]**  
**Protein Functional Class  
Prediction**

**Data Integration Method based on SSL**  
**MIPS Yeast Proteins/PDBselect25-GO**  
**Prediction from Multiple Protein Networks**

**Cancer Clinical Outcome  
Prediction**

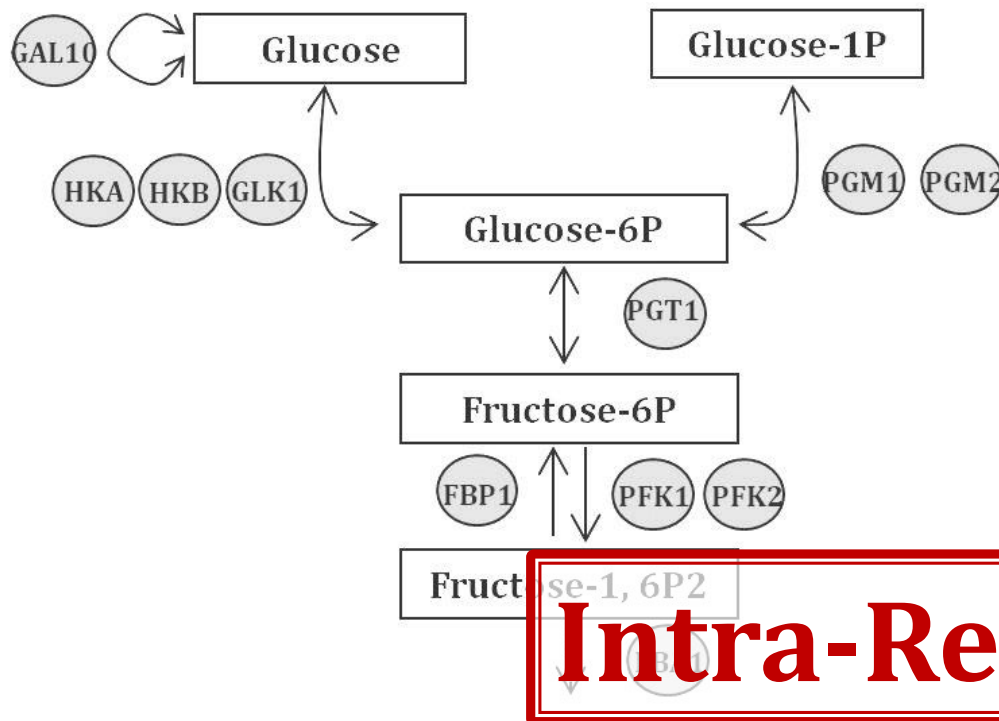
**Brain Cancer (GBM)/Ovarian Cancer (OV)**  
**Prediction from Multiple Genomic Data**

**[Inter-Relation]**  
**Cancer Clinical Outcome  
Prediction**

**Network Reconstruction using the**  
**Information From miRNA to Gene Expression**

**Closing Remarks**

**Future Work**



## Objective Function

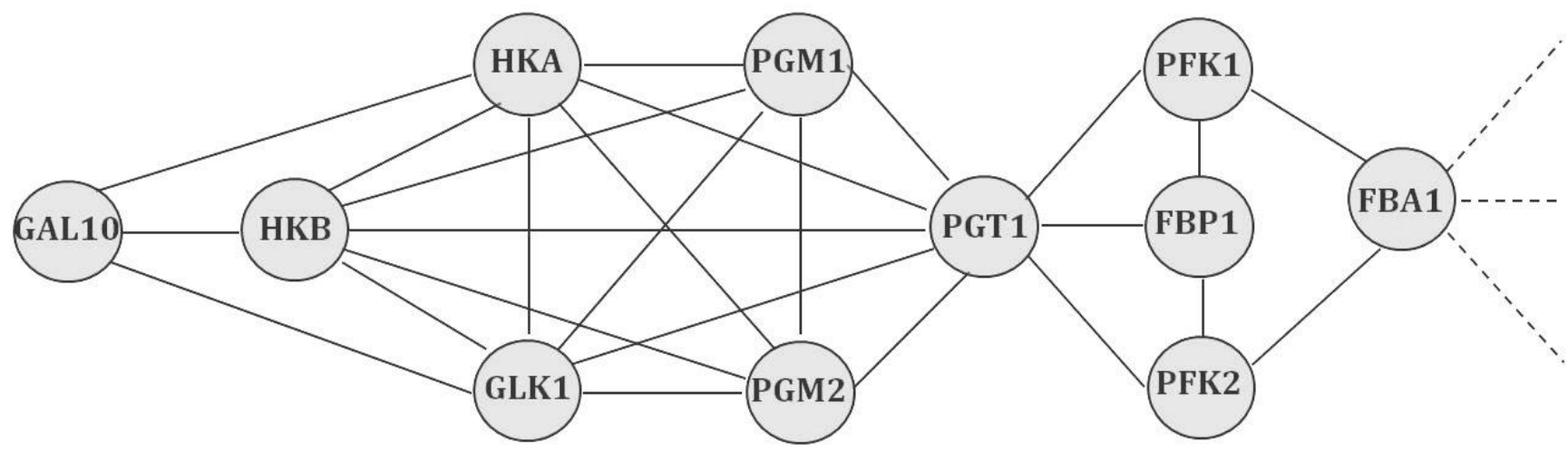
$$\min_f \mu f^T L f + (f - y)^T (f - y)$$

## Solution

$$f = \{ I + \mu L \}^{-1} y$$

# Intra-Relation

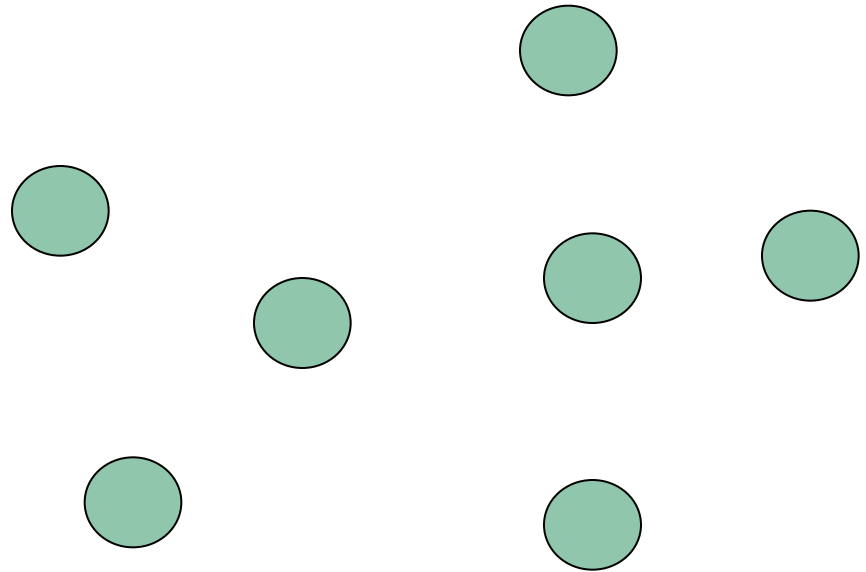
$$L = D - W, D = \text{diag}(d_i), d_i = \sum_j w_{ij}$$



# Intra-Relation: Graph Representation

**Nodes** : Protein Network  
*Proteins*

: Patient Network  
*Patient Samples*



# Intra-Relation: Graph Representation

## Labels

### Patient Network

#### *Clinical Outcome*

(ex) Breast Cancer

: Survived (+1) or Not (-1)

(ex) Brain Cancer

: Glioblastoma Multiforme

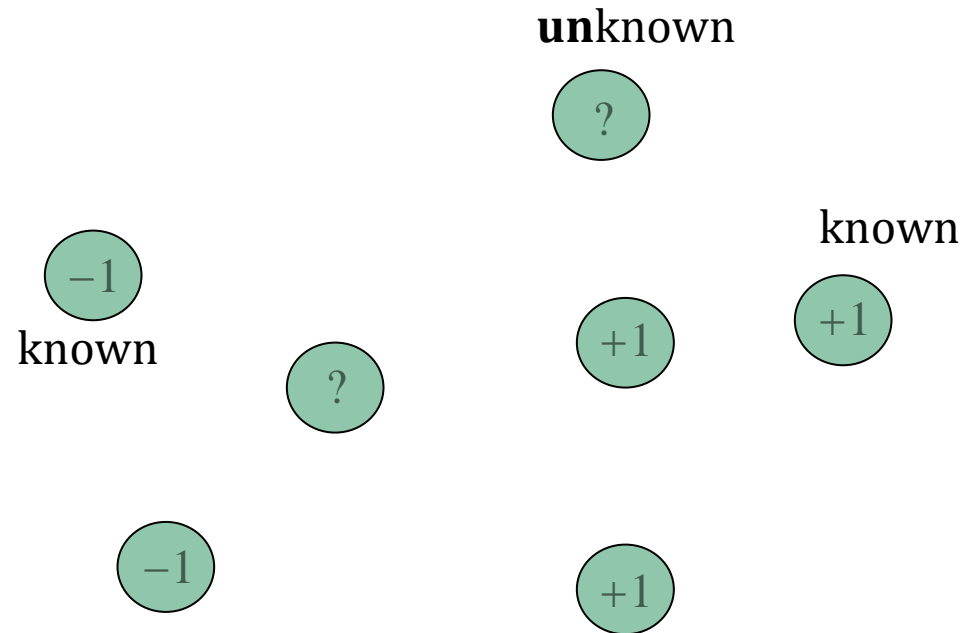
: Recurrent tumor (+1) or Initial (-1)

(ex) Ovarian Cancer

: Serous cystadenocarcinoma

: Early stage T1/T2 (+1) or

Late stage T3/T4 (-1)



+1/-1 : Labeled *patient samples* with/without a *specific clinical outcome*

? : Unlabeled patient samples

# Intra-Relation: Graph Representation

## Edges

### Protein Network

#### *Similarities between Proteins*

##### (ex) Physical Interaction

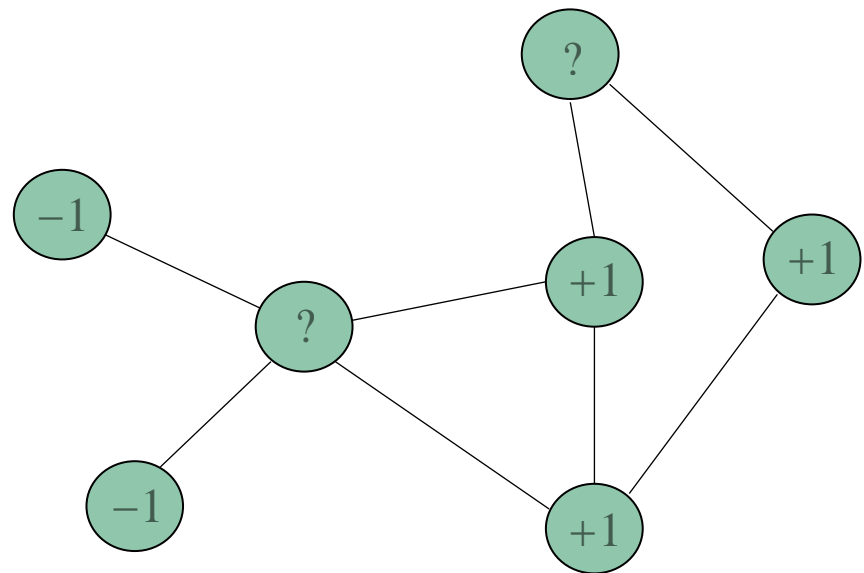
: Two proteins physically interact  
(e.g., docking)

##### (ex) Metabolic Pathway

: Two enzymes catalyzing  
successive reactions

##### (ex) Pfam domain structure

: Two proteins which show similar pattern  
in presence or absence of Pfam domains



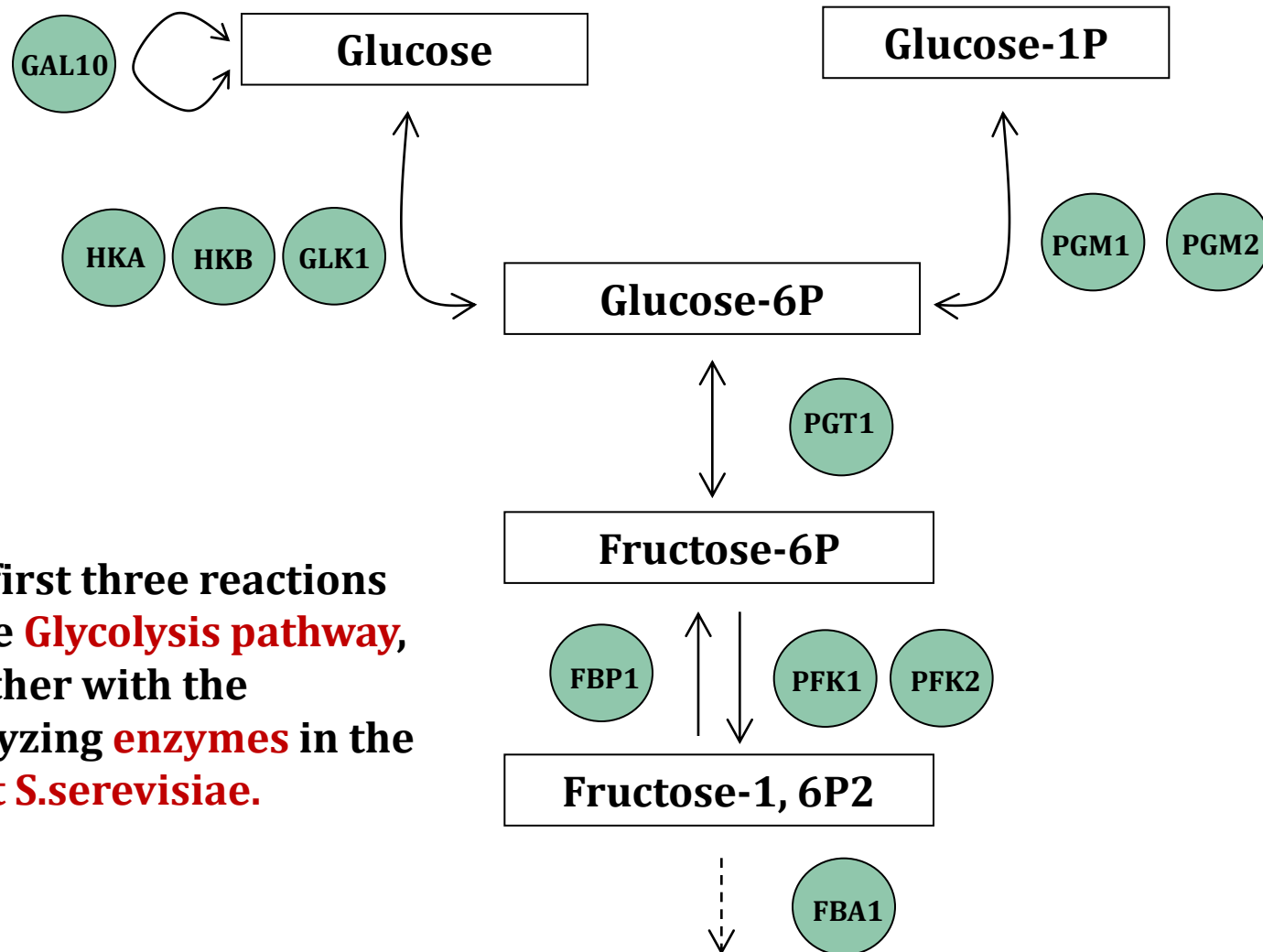
# Intra-Relation: Graph Representation – “Creating a Graph”

## Naturally Given Graphs

## Example: Metabolic Gene Network

- Graph Representation on Biological Networks
- The task is to predict (unidentified) functional classes of proteins using metabolic pathways

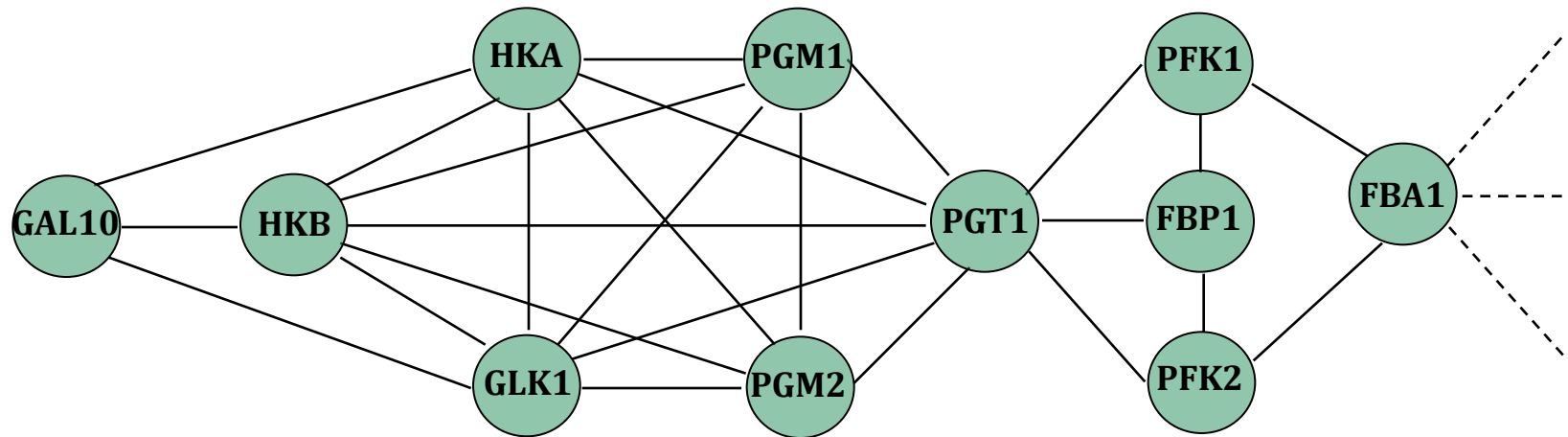
# Intra-Relation: Graph Representation – “Creating a Graph”



The first three reactions of the **Glycolysis pathway**, together with the catalyzing **enzymes** in the **Yeast *S. cerevisiae***.



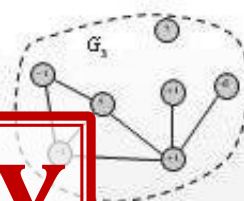
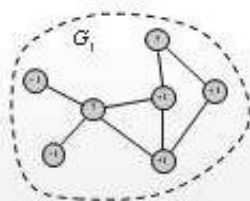
# Intra-Relation: Graph Representation – “Creating a Graph”



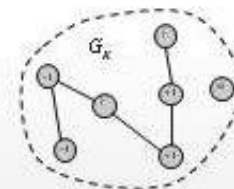
$$\min_{\beta} d(\beta) \equiv y^T (I + \sum_{k=1}^K \beta_k L_k)^{-1} y$$

$$y = \left\{ I + \sum_{k=1}^K \beta_k L_k \right\} f$$

Given Graphs



.....



**THEORY**

SDP/SVM

Kernel matrix

**Prediction from Intra-Relation  
Mathematical Formulation**

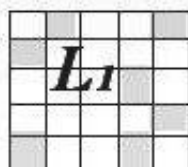
*Dense*  
 $K(\mu)$

$K_K$

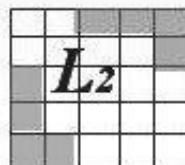
SSL

Laplacian matrix  $L$  (or Similarity matrix  $W$ )

*Sparse*  
 $L(\beta) = \beta_1$



$+ \beta_2$



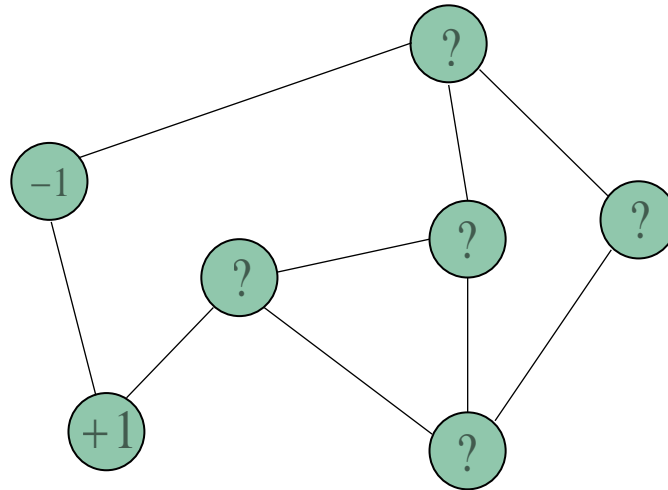
$+ \beta_3$



$+ \dots + \beta_k$



# Prediction from Intra-Relation: Graph-based SSL



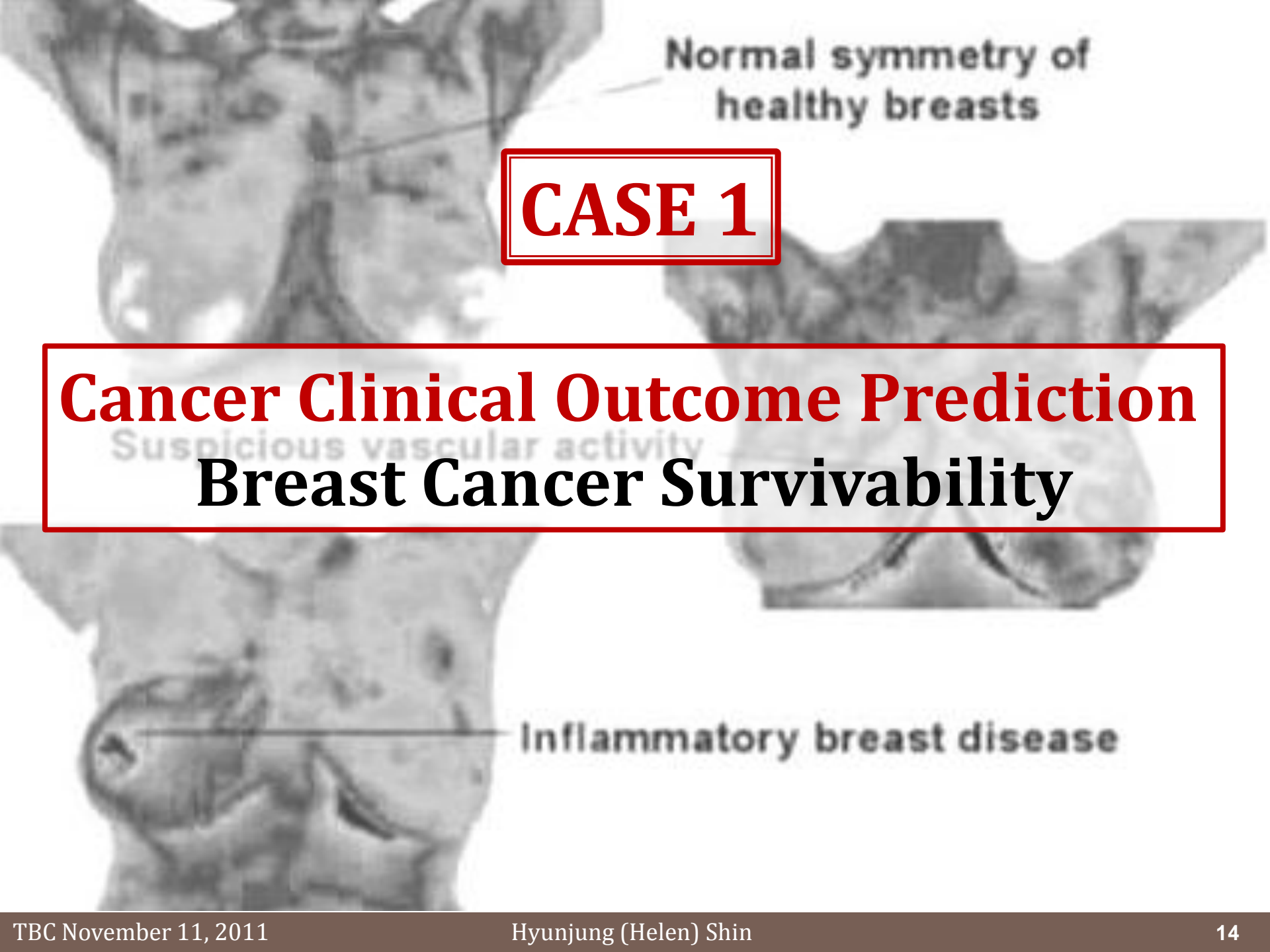
**Objective Function**

$$\min_f \mu f^T L f + (f - y)^T (f - y)$$

**Solution**

$$f = \{ I + \mu L \}^{-1} y$$

where  $L = D - W$ ,  $D = \text{diag}(d_i)$ ,  $d_i = \sum_j w_{ij}$



Normal symmetry of  
healthy breasts

## CASE 1

# Cancer Clinical Outcome Prediction Breast Cancer Survivability

Suspicious vascular activity

Inflammatory breast disease

# Background

According to *American Cancer Society*...

- Estimated **new cases** of breast cancer in 2010 in US..
  - **Females: 207,090**
  - **Males: 1,970**
- Estimated **deaths**
  - **Females: 39,840**
  - **Males: 390**

## Prediction of Breast Cancer Survivability

- **“Survival”** is defined as patient remaining alive for a specified period of time after the diagnosis of cancer
- **Cancer Prognosis** helps in establishing a treatment plan by predicting the outcome of a disease

# Data

---

## **Surveillance, Epidemiology, End Results (SEER) cancer incident data**

**162,500 Breast cancer patient records**

**16 attributes**

**1 class label (Survivability)**

: +1 (not survive)

: - 1 (survived)

# Data: Attributes

<b>Stage</b>	Defined by the size of cancer tumor and its spread
<b>Grade</b>	How does the tumor looks like and its resemblance to more or less aggressive tumors
<b>Lymph Node Involvement</b>	None, (1-3) Minimal, (4-9) Significant etc
<b>Race</b>	Ethnicity like White, Black, Chinese etc.
<b>Age at Diagnosis</b>	Actual age of patient in years
<b>Marital Status</b>	Married, Single, Divorced, Widowed, Separated
<b>Primary Site</b>	Presence of tumor at a particular location in body. Topographical classification of cancer
<b>Tumor Size</b>	2-5 cm, at 5cm prognosis worsens
<b>Site Specific Surgery</b>	Information on surgery during first course of therapy whether it was cancer directed or not.
<b>Radiation</b>	None, Beam Radiation, Radioisotopes, Refused, Recommended etc.
<b>Histological Type</b>	The form and structure of tumor
<b>Behavior Code</b>	Normal or aggressive behaviors of tumor have been defined in codes.
<b># of Positive Nodes Examined</b>	When the lymph nodes are involved in the cancer, they are called "positive."
<b># of Nodes Examined</b>	Total nodes (positive/negative) examined
<b># of Primaries</b>	Number of primary tumors (1-6)
<b>Clinical Extension of Tumor</b>	Defines the spread of tumor relative to breast
<b>Survivability</b>	Target binary variable defines the class of survival of patient.



# Model Comparison: Predictive models

---

- ❖ **Artificial Neural Network (ANN)**
- ❖ **Support Vector Machine (SVM)**
- ❖ **Semi-Supervised Learning (SSL)**  
*with a patient network*

# Aspects of Comparison

Let the **oncologists** (medical specialists) run a predictive model **by himself** and **interpret** the results with his medical domain knowledge !

Then, a **predictive model** has the properties of

**Reasonable  
Accuracy**

&

**Robustness**  
over  
Parameter  
Variation

# Aspects of Comparison

## Model Parameters (to be tuned)

### ANN

Random Seed = {1, 3, 5, 7, 10}

Hidden Node = {3, 6, 9, 12, 15}

### SVM

C = {0.2, 0.4, 0.6, 0.8, 1}

Gamma = {0.0001, 0.001, 0.01, 0.1, 1}

### SSL

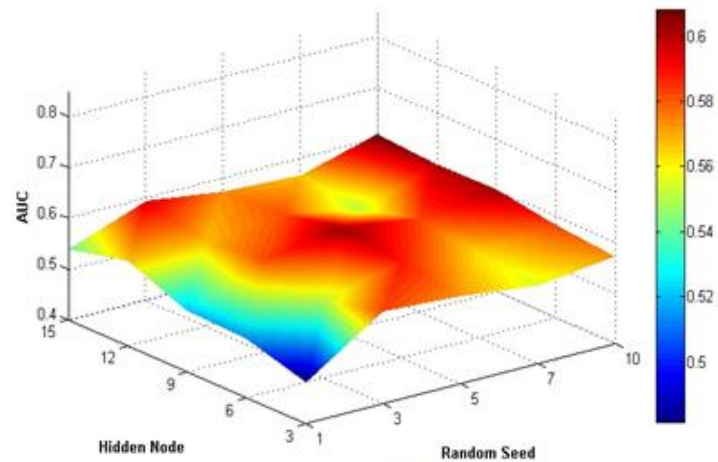
Mu = {0.0001, 0.01, 1, 100, 1000}

K = {3, 7, 15, 20, 30}

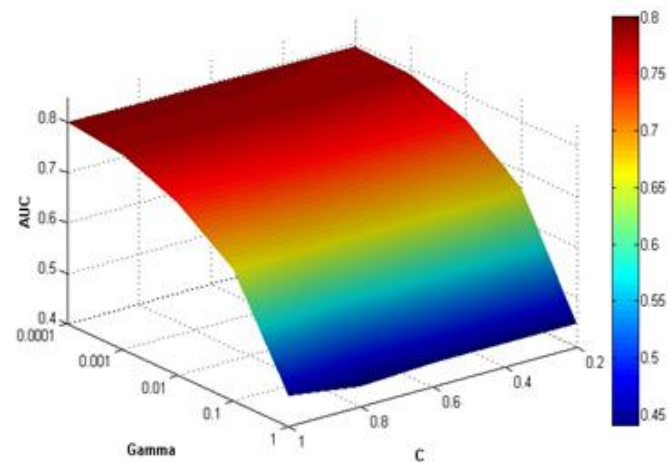
# Experimental Results: Accuracy

Dataset	ANN	SVM	SSL
	Avg_AUC	Avg_AUC	Avg_AUC
1	0.59	0.68	0.76
2	0.56	0.69	0.77
3	0.55	0.68	0.75
4	0.56	0.68	0.75
5	0.56	0.70	0.77
6	0.54	0.71	0.75
7	0.57	0.67	0.75
8	0.58	0.69	0.78
9	0.56	0.70	0.76
10	0.59	0.71	0.76
Mean (St Dev)	<b>0.57</b> <b>(<math>\pm 0.07</math>)</b>	<b>0.69</b> <b>(<math>\pm 0.13</math>)</b>	<b>0.76</b> <b>(<math>\pm 0.03</math>)</b>

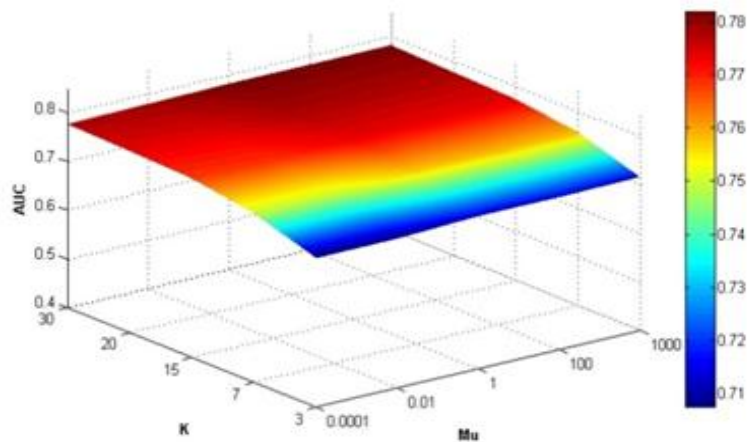
# Experimental Results: Robustness over parameter variation



(a) ANN



(b) SVM



(c) SSL

# Result Wrap-up

---

Building a **patient network** (Graph Representation) from patient samples is **straightforward**.

**Prediction algorithms** based on **Intra-relation** is well established.

The algorithm shows reasonably **high accuracies**, **stability (or robustness)** over model parameter variation, and is easy to use !

$$L(\beta) = \beta_1 \text{ (graph } G_1) + \beta_2 \text{ (graph } G_2) + \beta_3 \text{ (graph } G_3) + \dots + \beta_k \text{ (graph } G_k)$$

$$L(\beta) = \sum_{k=1}^K \beta_k L_k$$

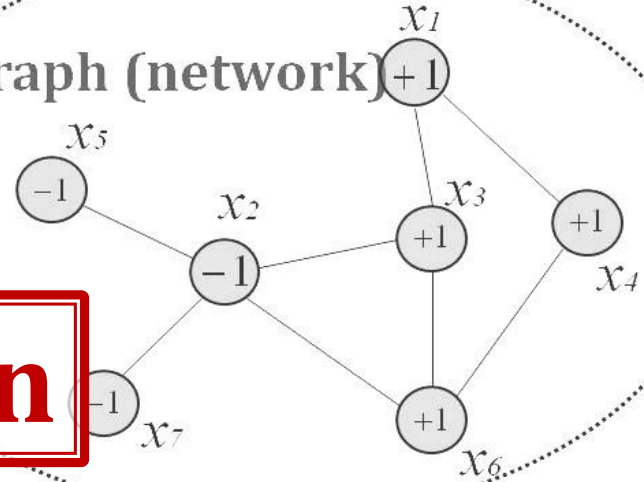
## Heterogeneous Data Sources

### Vectorial Data

	$A_1$	$A_2$	...	$A_{10}$	
$x_1$	10	5	...	1000	1
$x_2$	6	6	...	3500	-1
$x_3$	7	7	...	400	1
...	...	...	...	...	...
$x_7$	3	88	...	700	-1

**Integration**

### Graph (network)



### Sequence (string)

$x_1$	agctgttagctatatgcgtataggget	1
$x_2$	cagtgtcgaatagecgetegaaaaa a	-1
...	...	...
$x_7$	catgetgtatgcccgatagegtgatcg	-1

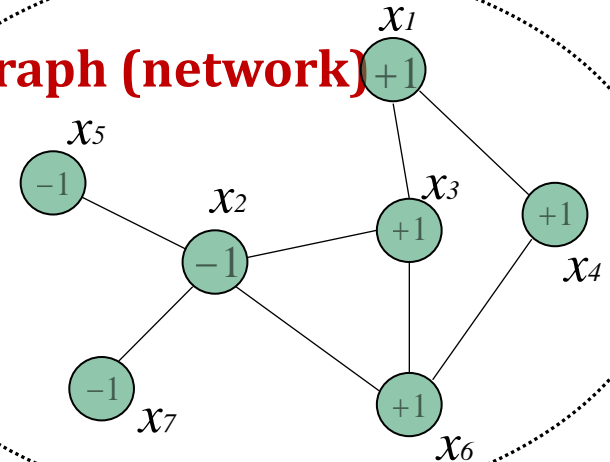
# Abstract: Data Integration

## Heterogeneous Data Sources

### Vectorial Data

	$A_1$	$A_2$	...	$A_{10}$	$y$
$x_1$	10	5	...	1000	1
$x_2$	6	6	...	3500	-1
$x_3$	7	7	...	400	1
...	...	...	...	...	...
$x_7$	3	88	...	700	-1

### Graph (network)



+

### Sequence (string)

$x_1$	agctgttagctatatgcgtatagggct	1
$x_2$	cagtgtcgaatagccgctcgaaaaa a	-1
...	...	...
$x_7$	catgctgtatgccgatagcgtgatcg	-1

+



# **Abstract: Data Integration**

---

**Data Integration** is concerned with the **integration of different or heterogeneous data sources** in order to **enhance the total information** about the problem at hand.

**Each** of data sources contains **partly independent** and **partly complementary** pieces of information about the problem...

# Outline

---

**[CASE 2-1]  
Protein Function  
Prediction**

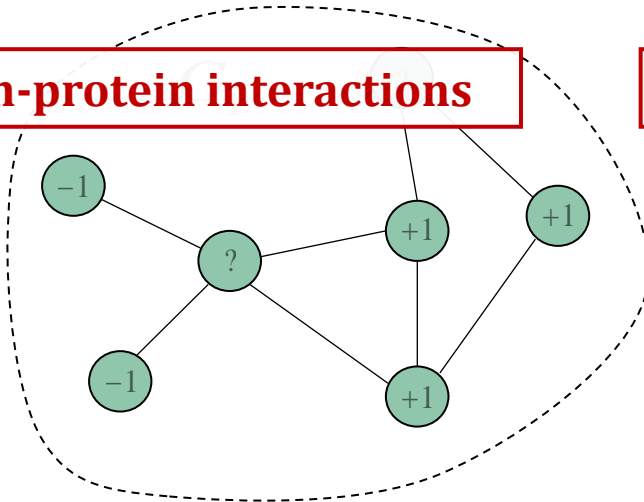
**MIPS Yeast Proteins/PDBselect25-GO  
Prediction from Multiple Protein Networks**

**[CASE 2- 2]  
Cancer Clinical Outcome  
Prediction**

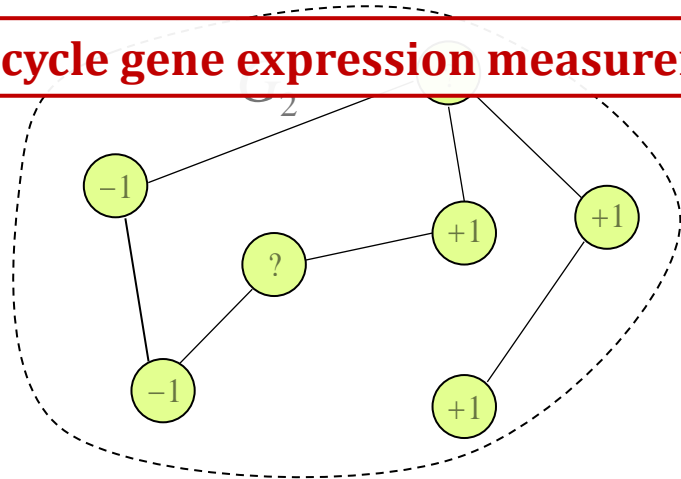
**Brain Cancer (GBM)/Ovarian Cancer (OV)  
Prediction from Multiple Genomic Data**

# If Multiple Graphs are Given?

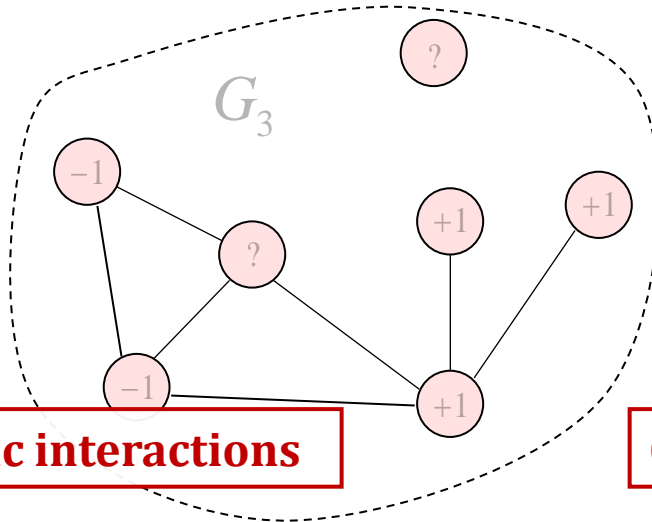
**Protein-protein interactions**



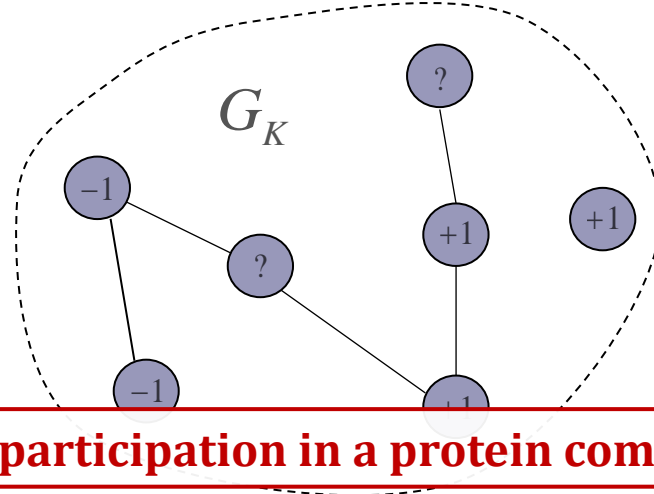
**Cell cycle gene expression measurements**



**Genetic interactions**



**Co-participation in a protein complex**



# If Multiple Graphs are Given?

## Example: Multiple Graph Sources on Proteins



### Physical interactions of the proteins

[Schwikowski,et al., 2000, Uetz et al., 2000, von Mering et al., 2002]

### Gene regulatory relationships

[Lee et al., 2002, Ihmels et al., 2002, Segal et al., 2003]

### Edges in a metabolic pathway

[Kanehisa et al., 2004]

### Similarities between protein sequences

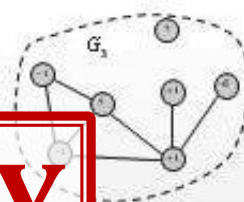
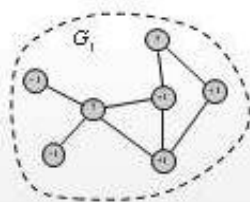
[Yona et al., 1999]

etc.

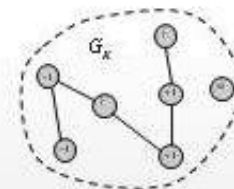
$$\min_{\beta} d(\beta) \equiv y^T (I + \sum_{k=1}^K \beta_k L_k)^{-1} y$$

$$y = \left\{ I + \sum_{k=1}^K \beta_k L_k \right\} f$$

Given Graphs



.....



**THEORY**

SDP/SVM

Kernel matrix

**Prediction from Integration  
Mathematical Formulation**

*Dense*  
 $K(\mu_1, \mu_2)$

$\mu_k$



SSL

Laplacian matrix  $L$  (or Similarity matrix  $W$ )

*Sparse*

$$L(\beta) = \beta_1 \begin{bmatrix} \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \end{bmatrix} \mathbf{L}_1 + \beta_2 \begin{bmatrix} \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \end{bmatrix} \mathbf{L}_2 + \beta_3 \begin{bmatrix} \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \end{bmatrix} \mathbf{L}_3 + \dots + \beta_k \begin{bmatrix} \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \end{bmatrix} \mathbf{L}_k$$

# Graph Integration using SSL

## Multiple Graph (Data) Integration

$$L(\beta) = \beta_1 \begin{array}{c} \text{Graph } G_1 \\ \text{(green nodes)} \end{array} + \beta_2 \begin{array}{c} \text{Graph } G_2 \\ \text{(yellow nodes)} \end{array} + \beta_3 \begin{array}{c} \text{Graph } G_3 \\ \text{(pink nodes)} \end{array} + \dots + \beta_k \begin{array}{c} \text{Graph } G_k \\ \text{(blue nodes)} \end{array}$$

**Shin, H., Lisewski, A.M. and Lichtarge, O. (2007)**

Graph sharpening plus graph integration: a synergy that improves protein functional classification, *Bioinformatics*, 23, 3217-3224.

**Shin, H. and Tsuda, K. (2006)**

Prediction of Protein Function from Networks, in *Book: Semi-Supervised Learning*, MIT press, Chapter 20, 339-352.

**Tsuda, K., Shin, H. and Scholkopf, B. (2005)**

Fast protein classification with multiple networks, *Bioinformatics*, 21 Suppl 2, ii59-65.



The background of the slide features a detailed diagram of the central dogma of molecular biology. At the top left, a DNA double helix is shown undergoing transcription, with a polymerase enzyme moving along it to synthesize an mRNA strand. The mRNA is labeled 'mRNA' and has a sequence of bases visible. Below this, the process of translation is depicted. An mRNA strand with the sequence 'AUGCGUAGACCUACUGACA' is shown. A ribosome is attached to the mRNA, and a tRNA molecule with an anticodon 'GCA' is shown pairing with a codon 'GCA' on the mRNA. The ribosome is labeled 'ribosome'. A polypeptide chain is shown emerging from the ribosome, with individual amino acids represented as beads. The chain is labeled 'polypeptide chain' and 'amino acid'. The entire process is labeled 'Translation'.

## CASE 2-1

# Protein Function Prediction

## MIPS Yeast Proteins/PDBselect25-GO

# Experiment I

(H.Shin, K.Tsuda, and B.Schoelkopf, *Bioinformatics*, 2005)

Task : Protein Functional Class Classification  
Model : Graph Integration based on SSL  
Data : MIPS Comprehensive Yeast Genome Database



# Protein Function Prediction: Experiment I

## MIPS Comprehensive Yeast Genome Database ([CYGD-mips.gsf.de/proj/yeast](http://CYGD-mips.gsf.de/proj/yeast))

**Data**

**3588** yeast proteins

**Output**

**13** functional categories

Binary classification for each category

**Input**

**5** networks

**Setting**

5 fold cross validation

5 times repetition

# Protein Function Prediction: Experiment I

MIPS Comprehensive Yeast Genome Database ([CYGD-mips.gsf.de/proj/yeast](http://CYGD-mips.gsf.de/proj/yeast))

## 13 CYGD functional Classes

1. metabolism
2. energy
3. cell cycle and DNA processing
4. transcription
5. protein synthesis
6. protein fate
7. cellular transportation and transportation mechanism
8. cell rescue, defense and virulence
9. interaction with cell environment
10. cell fate
11. control of cell organization
12. transport facilitation
13. others

# Protein Function Prediction: Experiment I

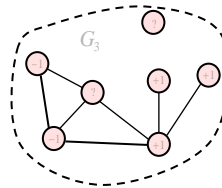
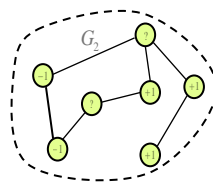
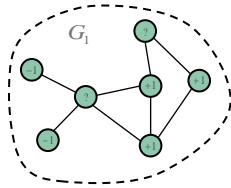
## Input Data Sources (5 networks)

- $W_1$  Network created from Pfam domain structure. A protein is represented by a **4950-dimensional binary vector**, in which each bit represents the **presence or absence of one Pfam domain**. An edge is created if the inner product between two vectors exceeds 0.06. The edge weight corresponds to the inner product.
- $W_2$  Co-participation in a protein complex (determined by tandem affinity purification, TAP). An edge is created if there is a **bait-prey relationship** between two proteins.
- $W_3$  Protein-protein interactions (MIPS physical interactions)
- $W_4$  Genetic interactions (MIPS genetic interactions)
- $W_5$  Network created from the cell cycle gene expression measurements [Spellman et al., 1998]. An edge is created if the **Pearson coefficient** of two profiles exceeds 0.8. The edge weight is set to 1. This is identical with the network used in [Deng et al., 2003]

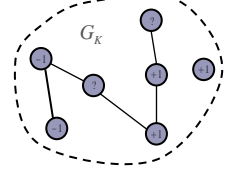
# Protein Function Prediction: Experiment I

## Density of Working Matrices

Given Graphs



.....



**SDP/SVM**

Kernel matrix

Dense

$$K(\mu) = \mu_1 \begin{matrix} & & K1 \\ & & \begin{matrix} \text{5x5 grid} \end{matrix} \\ & & \end{matrix} + \mu_2 \begin{matrix} & & K2 \\ & & \begin{matrix} \text{5x5 grid} \end{matrix} \\ & & \end{matrix} + \mu_3 \begin{matrix} & & K3 \\ & & \begin{matrix} \text{5x5 grid} \end{matrix} \\ & & \end{matrix} + \dots + \mu_k \begin{matrix} & & Kk \\ & & \begin{matrix} \text{5x5 grid} \end{matrix} \\ & & \end{matrix}$$

**SSL**

Laplacian matrix L (or Similarity matrix W)

Sparse

$$L(\beta) = \beta_1 \begin{matrix} & & L1 \\ & & \begin{matrix} \text{5x5 grid with sparse green blocks} \end{matrix} \\ & & \end{matrix} + \beta_2 \begin{matrix} & & L2 \\ & & \begin{matrix} \text{5x5 grid with sparse green blocks} \end{matrix} \\ & & \end{matrix} + \beta_3 \begin{matrix} & & L3 \\ & & \begin{matrix} \text{5x5 grid with sparse purple blocks} \end{matrix} \\ & & \end{matrix} + \dots + \beta_k \begin{matrix} & & Lk \\ & & \begin{matrix} \text{5x5 grid with sparse orange blocks} \end{matrix} \\ & & \end{matrix}$$

# Protein Function Prediction: Experiment I

## Methods in Comparison

$L_k$	Label propagation with an <b>Individual</b> Graphs (k=1...5)
$L_{opt}$	Laplacian of Combined Graph with <b>Optimized Weights</b>
$L_{fix}$	Label propagation with <b>Equal Weights</b>
MRF	<b>Markov Random Field</b> , proposed by Deng et al [2003]
SDP/SVM	<b>Semi-definite Programming</b> based <b>Support Vector Machines</b> , proposed by Lanckriet et al [2004]

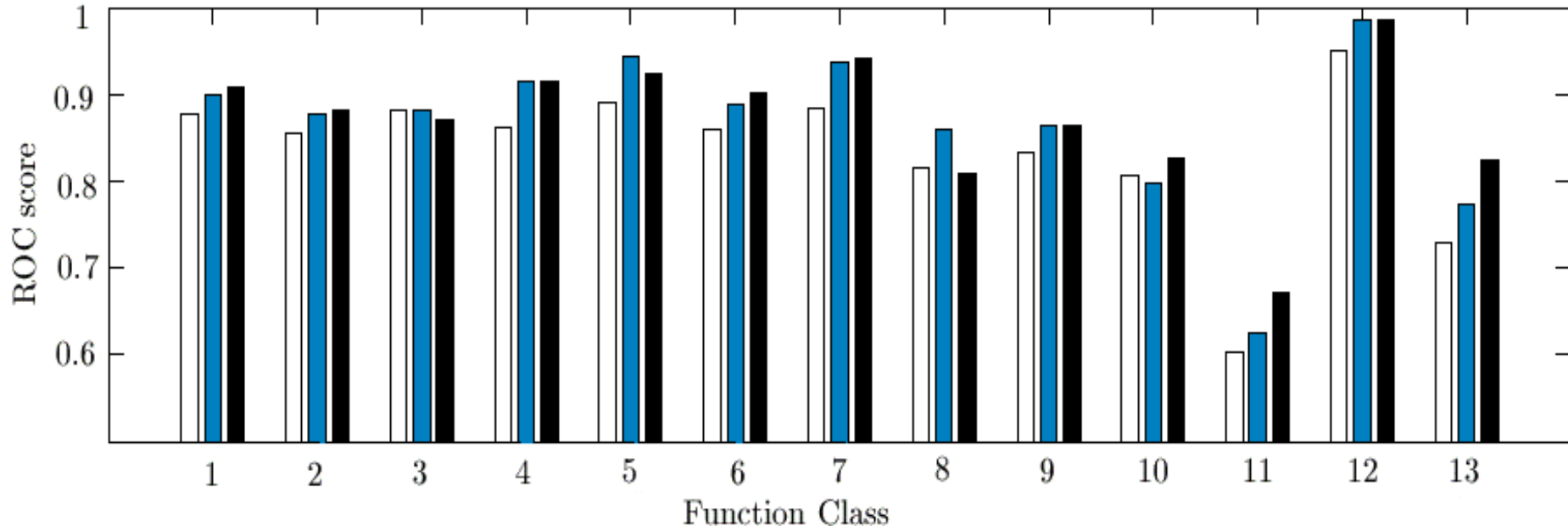
# Protein Function Prediction: Experiment I

## Results : Integrated Network vs. the Best Performing Individual (ROC scores)

White: the best performing individual network

Blue:  $L_{\text{fix}}$

Black:  $L_{\text{opt}}$



**Across the 13 classes, the proposed integrated network outperforms the best performing individual.**

# Protein Function Prediction: Experiment I

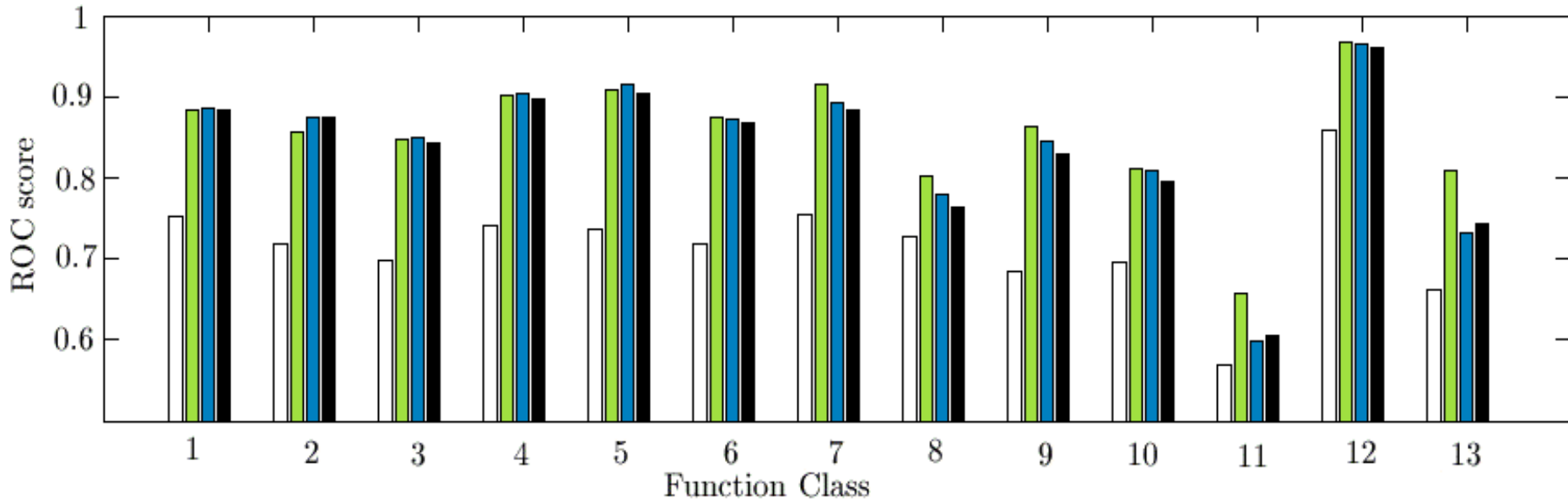
## Results : the proposed vs. others integration methods (ROC scores)

White: MRF

Green: SDP/SVM

Blue:  $L_{\text{fix}}$

Black:  $L_{\text{opt}}$

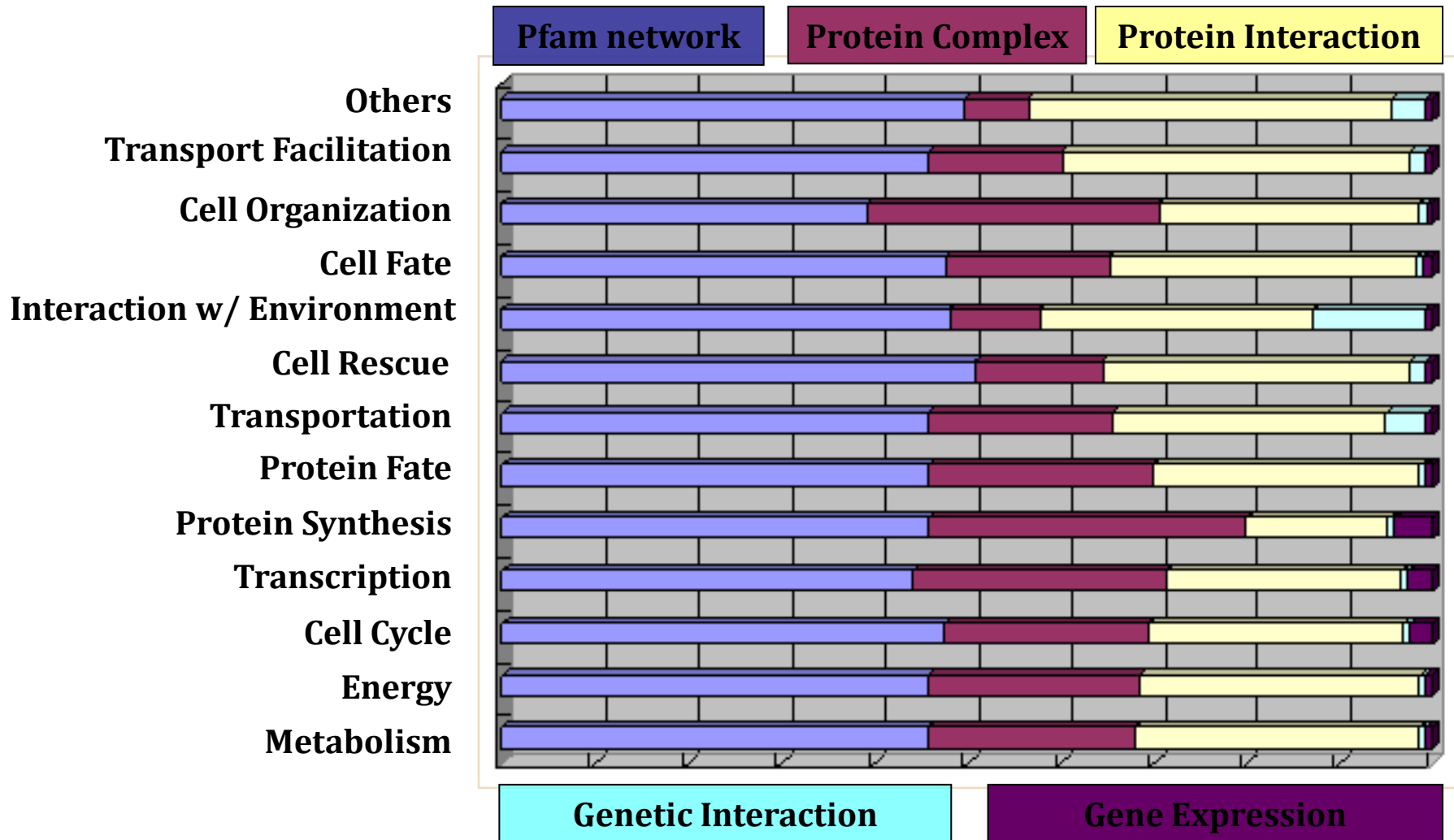


For most classes, the proposed method achieves high scores, which are **similar to the SDP/SVM methods**.

In classes 11 and 13, the proposed method performs poor (but still better than the MRF method), However, taking into account the Simplicity and Efficiency the method shows the promising results

# Protein Function Prediction: Experiment I

**Results :** Which data source is more informative?





# Protein Function Prediction: Experiment I

## Results : Computational Time

### The proposed:

**49.3 seconds (std. 14.8)**

### SDP/SVM :

**Approx. Several CPU days**

(G. Lanckriet, personal communication)

\* Measured in a standard 2.2Ghz PC with 1GByte memory

# Protein Function Prediction: Experiment I

## Results : Computational Time

### The proposed:

Nearly linearly proportional to the number of non-zero entries of sparse matrices

### SDP/SVM :

$$O((m+n)^2 n^{2.5})$$

# Protein Function Prediction: Experiment I

## Wrap-Up

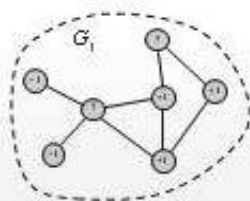
The proposed integrated network of multiple data sources **outperforms** the best performing **individuals**.

The proposed integration method is **simple, computationally efficient, scalable** when compared with the existing integration method such as SDP/SVM.

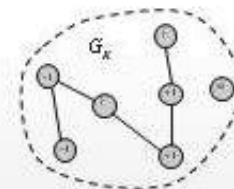
$$\min_{\beta} d(\beta) \equiv y^T (I + \sum_{k=1}^K \beta_k L_k)^{-1} y$$

$$y = \left\{ I + \sum_{k=1}^K \beta_k L_k \right\} f$$

Given Graphs



.....



**THEORY**

SDP/SVM

Kernel matrix

**Prediction from “Inter”-Relation  
Method & Mathematical Formulation**

SSL

Laplacian matrix  $L$  (or Similarity matrix  $W$ )

$$L(\beta) = \beta_1 \begin{bmatrix} \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \end{bmatrix} \mathbf{L}_1 + \beta_2 \begin{bmatrix} \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \end{bmatrix} \mathbf{L}_2 + \beta_3 \begin{bmatrix} \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \end{bmatrix} \mathbf{L}_3 + \dots + \beta_k \begin{bmatrix} \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } & \text{ } \end{bmatrix} \mathbf{L}_k$$

# Inter-Relation: Method/Mathematical Formulation

$G_O$  Patient Graph from Gene Expression Data (Original )

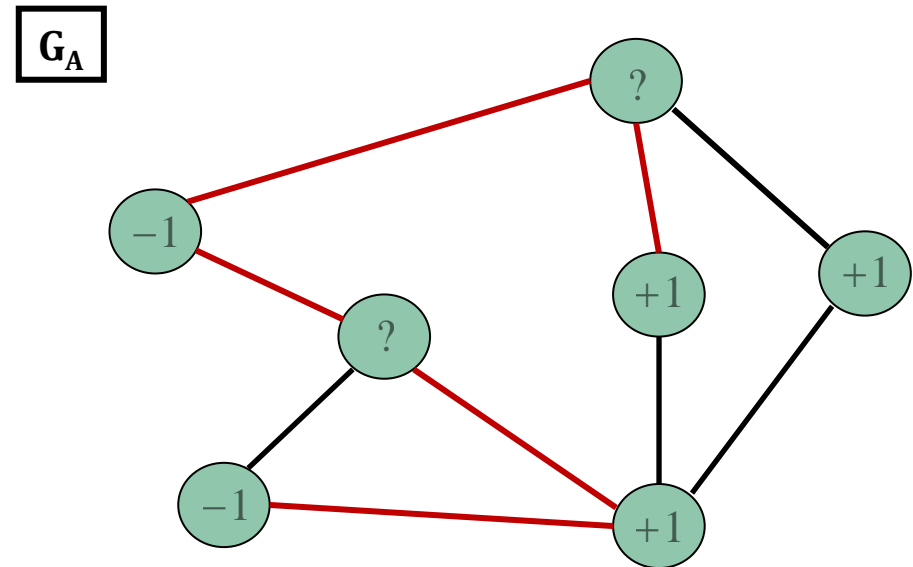
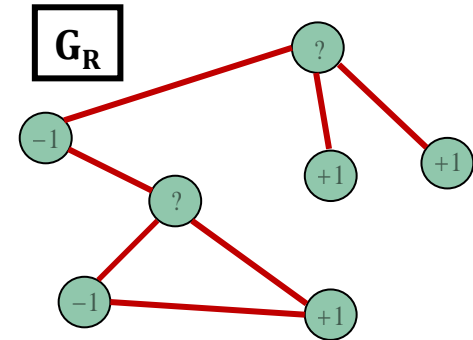
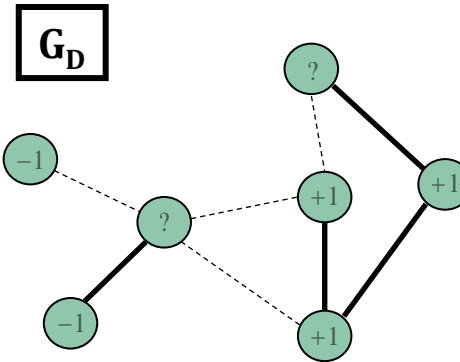
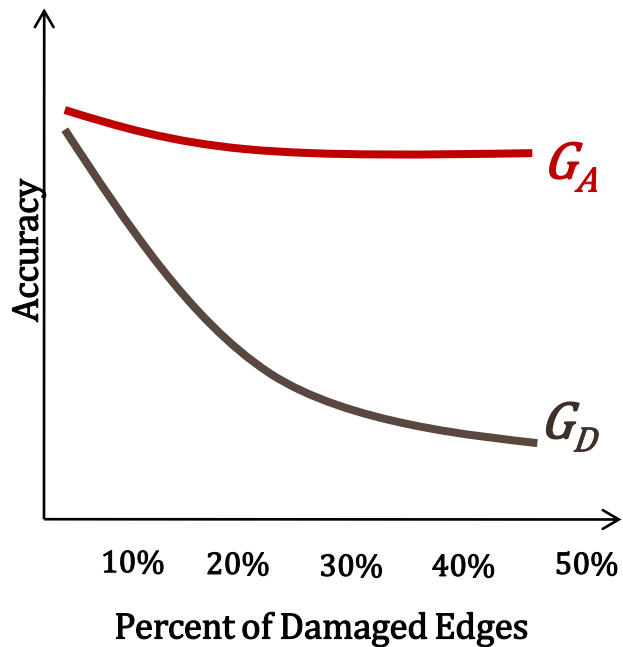
$G_D$  Damaged Graph

$G_R$  Reconstructed Graph (via inter-relationship)

$G_A$  Augmented graph ( $G_D + G_R$ )

# Inter-Relation: Method/Mathematical Formulation

**$G_A$**  Augmented graph ( $G_D + G_R$ )



## CASE 3

# Cancer Clinical Outcome Prediction Brain Cancer

A

# Experiment IV

(D.Kim, H. Shin, S. Lee and J. Kim, *TBC*, 2011)

Task : Brain Cancer Clinical Outcome Classification  
Model : Inter-Relation + SSL  
Data : The Cancer Genomic Atlas (TCGA database)



# Inter-Relation: Experiment - Data

## TCGA: Gene Expression & miRNA

**82 patient samples of Brain Cancer (GBM)**

**1 class label (Survivability)**

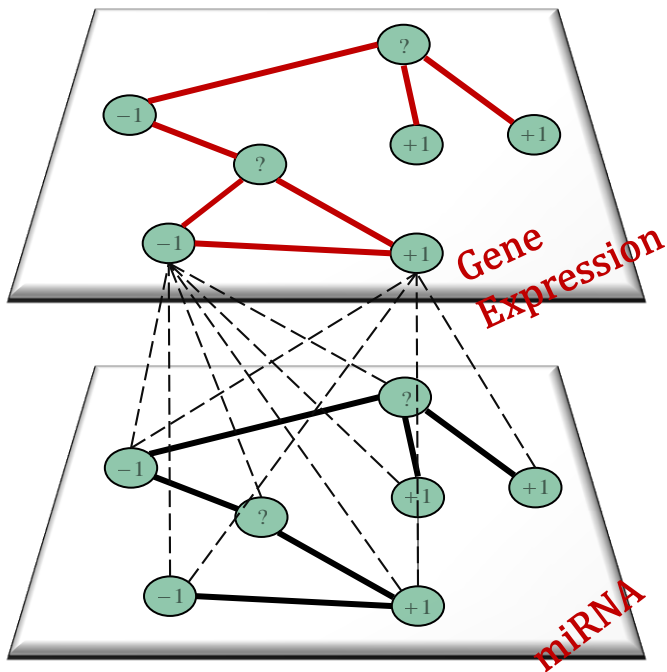
: - 1 (Short-term survival: #54)

: + 1 (Long-term survival: #28)

Data type	Platform	Num of Attributes
Gene Expression	Affymetrix HT Human Genome U133 Array Plate Set	12,043
miRNA	Agilent Human miRNA Microarray Rel12.0	799

# Inter-Relation: Experiment - Data

## miRNA & Target Gene (mRNA) Relation

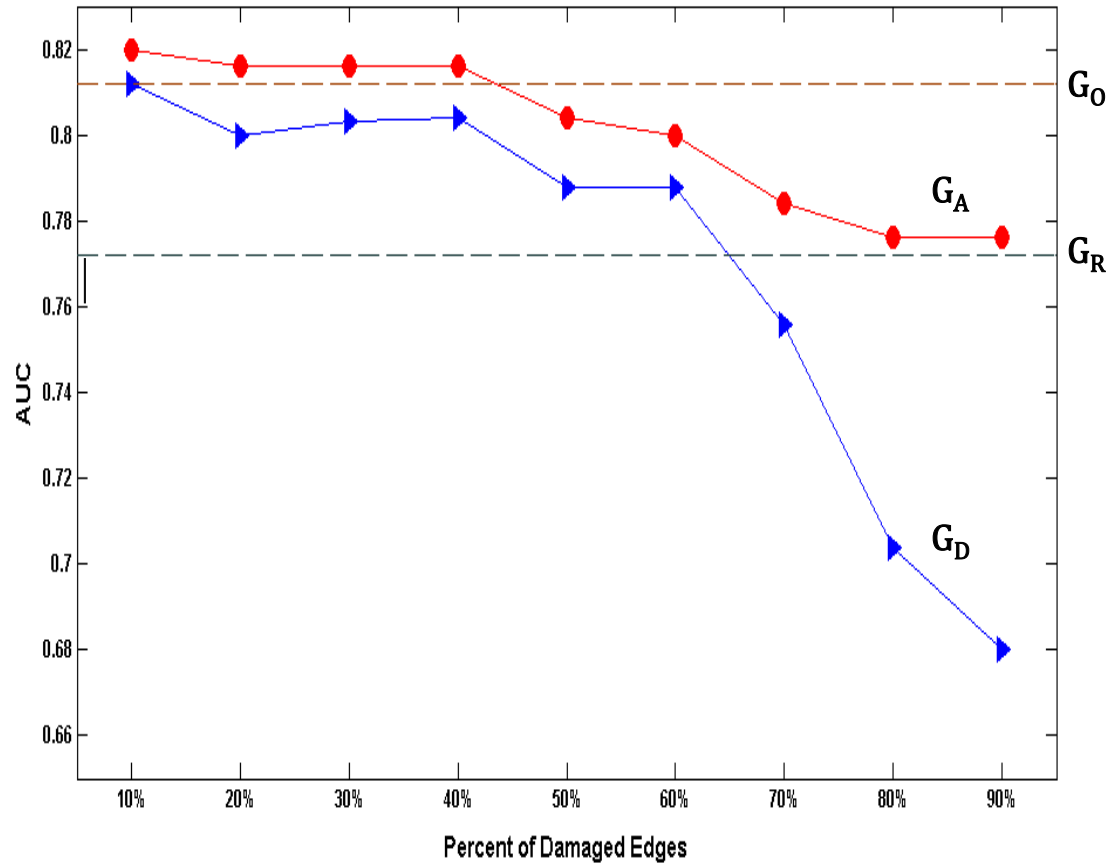


**miRecords** which is integrated resources of miRNA that store **target interactions** produced by **11 established miRNA target prediction program** (Xiao et al., 2009)

Among 11 algorithms, **a binary relation** between **miRNA** and **mRNA** was set when **more than 3 algorithms** provide the **target relation**

# Inter-Relation: Experiment – Comparison Results

**Improved Performance** from  
Augmented Knowledge via  
**INTER-RELATION**  
between miRNA and Gene  
Expression

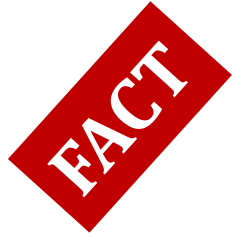


# Inter-Relation: Experiment – Comparison Results

**Significance for Differences in Performance**

Percent of damaged edges	AUC of $G_D$	AUC of $G_A$	P-value
10%	0.812	0.820	1.87e-02
30%	0.803	0.816	2.09e-03
50%	0.788	0.804	3.43e-05
70%	0.756	0.784	9.59e-08
90%	0.680	0.776	1.24e-13

**Improved Performance** from Augmented Knowledge via **INTER-RELATION** between miRNA and Gene Expression



## Inter-Relation between Different Levels of Biological Data

There exist **Interactions** between **two or more layers** in the hierarchy of different biological levels

**Ex) miRNAs regulate target genes**

CONTRIBUTION

## Knowledge Reconstruction/Augmentation via Inter-Relation

This work shows **how to extract the Knowledge** between two **layers** of biological process and **how to use it to complete the incomplete knowledge** in other levels

**A Method (or Mathematical frame work)** incorporating **Inter-Relation** and **Intra-Relation** is **proposed** and **validated** through a case example of **Cancer Phenotype Prediction** based on **miRNAs** and **Genes**

**Inter-Relation** from **miRNAs** to **Genes** **augments** the **Intra-Relation** among Genes, which leads to **better accuracies** and **perception** in cancer phenotype prediction

# Inter-Relation: Wrap-up

---

**CONCLUSION**

**Knowledge from Inter-Relation helps to Complete the Incomplete Knowing about Intra-Relation**

# Future works



## **Biological Data**

**Heterogeneity**

**&**

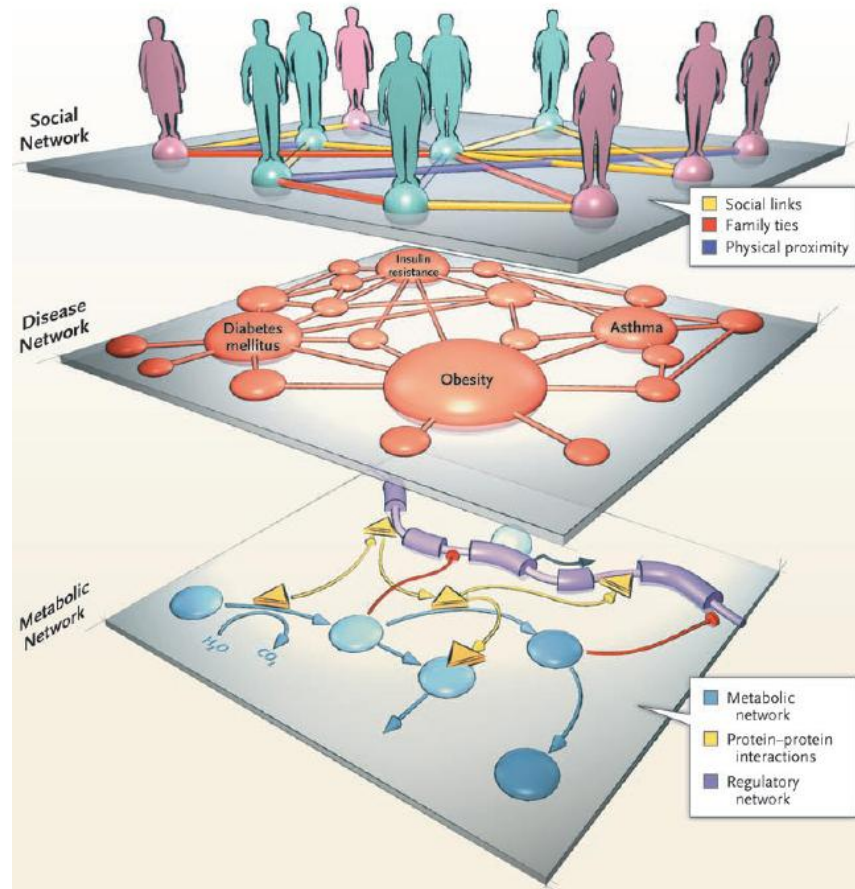
**Hierarchy**

# Future Works: Heterogeneity & Hierarchy

## Heterogeneous types and Hierarchical Structure of **Biological Data**

[Network Medicine]

Complex Networks  
of Direct Relevance



Barabasi, NEJM, 2007

# Acknowledgements

---

The authors would like to gratefully acknowledge support from Post Brain Korea 21 and the research grant from National Research Foundation of the Korean Government (2009-0065043/2011-0018257)

# References

- Bach, F., Lanckriet, G. and Jordan, M. (2004) Multiple kernel learning, conic duality, and the SMO algorithm, *In Proceedings of the Twenty-first International Conference on Machine Learning (ICML), Banff, Canada, ACM Press*, pp. 6-13.
- Belkin, M. (2004) Regularization and Semi-supervised Learning on Large Graphs, *In Proceedings of the 17th Annual Conference on Learning Theory (COLT) 3120. Lecture Notes in Computer Science*, 624-638.
- Ben-Hur, A. and Noble, W.S. (2005) Kernel methods for predicting protein-protein interactions, *Bioinformatics*, **21 Suppl 1**, i38-46.
- Berchuck, A., *et al.* (2005) Patterns of gene expression that characterize long-term survival in advanced stage serous ovarian cancers, *Clin Cancer Res*, **11**, 3686-3696.
- Beroukhi, R., *et al.* (2010) The landscape of somatic copy-number alteration across human cancers, *Nature*, **463**, 899-905.
- Bild, A.H., *et al.* (2006) Oncogenic pathway signatures in human cancers as a guide to targeted therapies, *Nature*, **439**, 353-357.
- Chapelle, O., Weston, J. and Scholkopf, B. (2003) Cluster kernels for semi-supervised learning, *Advances in Neural Information Processing Systems (NIPS)*, **15**, 585-592.
- Chin, L. and Gray, J.W. (2008) Translating insights from the cancer genome into clinical practice, *Nature*, **452**, 553-563.
- Chung, F.R.K. (1997) Spectral Graph Theory, *Number 92 in Regional Conference Series in Mathematics*.
- Demsar, J. (2006) Statistical comparisons of classifiers over multiple data sets, *Journal of Machine Learning Research*, **7**, 1-30.
- Fan, X., *et al.* (2010) DNA microarrays are predictive of cancer prognosis: a re-evaluation, *Clin Cancer Res*, **16**, 629-636.
- Furnari, F.B., *et al.* (2007) Malignant astrocytic glioma: genetics, biology, and paths to treatment, *Genes Dev*, **21**, 2683-2710.
- Golub, T.R., *et al.* (1999) Molecular classification of cancer: class discovery and class prediction by gene expression monitoring, *Science*, **286**, 531-537.
- Gribskov, M. and Robinson, N.L. (1996) Use of receiver operating characteristic (ROC) analysis to evaluate sequence matching, *Comput Chem*, **20**, 25-33.
- Hanash, S. (2004) Integrated global profiling of cancer, *Nat Rev Cancer*, **4**, 638-644.
- Huang, E., *et al.* (2003) Gene expression predictors of breast cancer outcomes, *Lancet*, **361**, 1590-1596.
- Jafari, P. and Azuaje, F. (2006) An assessment of recently published gene expression data analyses: reporting experimental design and statistical factors, *BMC Med Inform Decis Mak*, **6**, 27.
- Jansen, R., *et al.* (2003) A Bayesian networks approach for predicting protein-protein interactions from genomic data, *Science*, **302**, 449-453.
- Jemal, A., *et al.* (2009) Cancer statistics, 2009, *CA Cancer J Clin*, **59**, 225-249.
- Kondor, I. and Lafferty, J. (2002) Diffusion kernels on graphs and other discrete structures, *In Sammut, C. and Hoffmann, A.G. (eds), Proceedings of the Nineteenth International Conference on Machine Learning (ICML 2002), Sydney, Australia, Morgan Kaufmann*, pp. 315-322.
- Lanckriet, G.R., *et al.* (2004) A statistical framework for genomic data fusion, *Bioinformatics*, **20**, 2626-2635.
- Lu, J., *et al.* (2005) MicroRNA expression profiles classify human cancers, *Nature*, **435**, 834-838.
- Marko, N.F., *et al.* (2008) Genomic expression patterns distinguish long-term from short-term glioblastoma survivors: a preliminary feasibility study, *Genomics*, **91**, 395-406.

# References

- Mischel, P.S., Cloughesy, T.F. and Nelson, S.F. (2004) DNA-microarray analysis of brain cancer: molecular classification for therapy, *Nat Rev Neurosci*, **5**, 782-792.
- Myllykangas, S., *et al.* (2008) Classification of human cancers based on DNA copy number amplification modeling, *BMC Med Genomics*, **1**, 15.
- Ohn, J.H., Kim, J. and Kim, J.H. (2007) Genomic characterization of perturbation sensitivity, *Bioinformatics*, **23**, i354-358.
- Qiu, J. and Noble, W.S. (2008) Predicting co-complexed protein pairs from heterogeneous data, *PLoS Comput Biol*, **4**, e1000054.
- Roepman, P., *et al.* (2005) An expression profile for diagnosis of lymph node metastases from primary head and neck squamous cell carcinoma s, *Nat Genet*, **37**, 182-186.
- Salzman, M. and Kaplan, R. (1991) Intracranial tumors in adults, *In : Salzman M (ed) Neurology of brain tumors. Williams & Wilkins, Baltimore*, 1339-1352.
- Saxena, A., Robertson, J.T. and Ali, I.U. (1996) Abnormalities of p16, p15 and CDK4 genes in recurrent malignant astrocytomas, *Oncogene*, **13**, 661-664.
- Segal, E., *et al.* (2003) Module networks: identifying regulatory modules and their condition-specific regulators from gene expression data, *Nat Genet*, **34**, 166-176.
- Shin, H., Lisewski, A.M. and Lichtarge, O. (2007) Graph sharpening plus graph integration: a synergy that improves protein functional classification, *Bioinformatics*, **23**, 3217-3224.
- Shin, H. and Tsuda, K. (2006) Prediction of Protein Function from Networks, *in Book: Semi-Supervised Learning, Edited by Olivier Chapelle, Bernhard Schölkopf, Alexander Zien, MIT press, Chapter 20*, 339-352.
- Shridhar, V., *et al.* (2001) Genetic analysis of early- versus late-stage ovarian tumors, *Cancer Res*, **61**, 5895-5904.
- Spellman, P.T., *et al.* (1998) Comprehensive identification of cell cycle-regulated genes of the yeast *Saccharomyces cerevisiae* by microarray hybridization, *Mol Biol Cell*, **9**, 3273-3297.
- TCGA Network (2008) Comprehensive genomic characterization defines human glioblastoma genes and core pathways, *Nature*, **455**, 1061-1068.
- Troyanskaya, O., *et al.* (2001) Missing value estimation methods for DNA microarrays, *Bioinformatics*, **17**, 520-525.
- Tsuda, K., Shin, H. and Schölkopf, B. (2005) Fast protein classification with multiple networks, *Bioinformatics*, **21 Suppl 2**, ii59-65.
- van 't Veer, L.J., *et al.* (2002) Gene expression profiling predicts clinical outcome of breast cancer, *Nature*, **415**, 530-536.
- Verhaak, R.G., *et al.* (2010) Integrated genomic analysis identifies clinically relevant subtypes of glioblastoma characterized by abnormalities in PDGFRA, IDH1, EGFR, and NF1, *Cancer Cell*, **17**, 98-110.
- Waldman, F.M., *et al.* (2000) Chromosomal alterations in ductal carcinomas in situ and their in situ recurrences, *J Natl Cancer Inst*, **92**, 313-320.
- Wu, C.C., *et al.* (2010) Prediction of human functional genetic networks from heterogeneous data using RVM-based ensemble learning, *Bioinformatics*, **26**, 807-813.
- Zhou, D., *et al.* (2004) Learning with local and global consistency, *Advances in Neural Information Processing Systems (NIPS)*, **16**, 321-328.