

Mining Genetic Determinants of Human Disease with Graph-based Learning



TaeHyun Hwang and Rui Kuang

Department of Computer Science and Engineering, University of Minnesota Twin Cities

E-mail: thwang@cs.umn.edu, website: <http://compbio.cs.umn.edu/>



Understanding causative relations between genes and disease phenotypes has been a long standing challenge in genomics research. The recent 'omics projects offer the opportunity for a genome-wide study of the associations with the large scale human protein-protein interaction network, phenotype similarity network and their association network. Integrating the heterogeneous genomic and phenotypic data and exploring the modularity among the genomic and phenotypic objects becomes the central computational problem. The key challenge is two-fold: how to integrate heterogeneous data and how to utilize the modular structures. We propose to integrate the data as a large heterogeneous network and introduce an algorithm to explore the modules in the heterogeneous network. The proposed algorithm is a mathematically principled and yet intuitive label propagation method. Extensive empirical experiments support that the algorithm can improve the accuracy of disease phenotype-gene association prediction over currently available methods. The algorithm is also a general method that is applicable to similar network integration problems in other research areas.

BACKGROUND

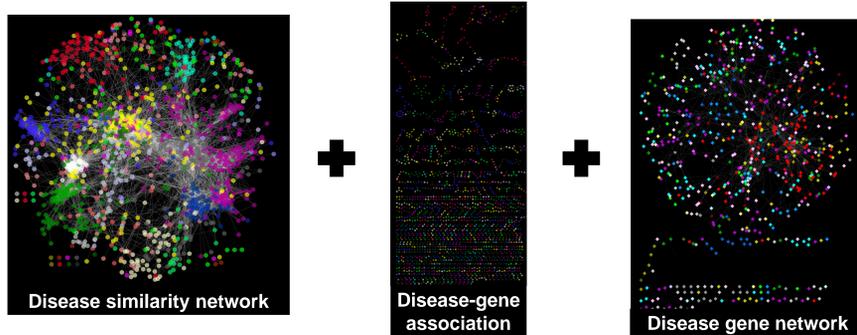
Motivation

- Many real world networks contain clusters, communities, cohesive groups or modules (such as functionally related genes, groups of people, and) among the objects in the network
- Several graph-based learning algorithms have been developed to utilize network structures to improve performance in different learning tasks on networks such as classification and ranking
- However, there is no general framework to explore network structures in the heterogeneous network containing different types of vertices and edges

METHODOLOGY

Problem Formulation

- Given:** Multiple network data. Each network has its own cluster structures containing different types of vertices and edges.

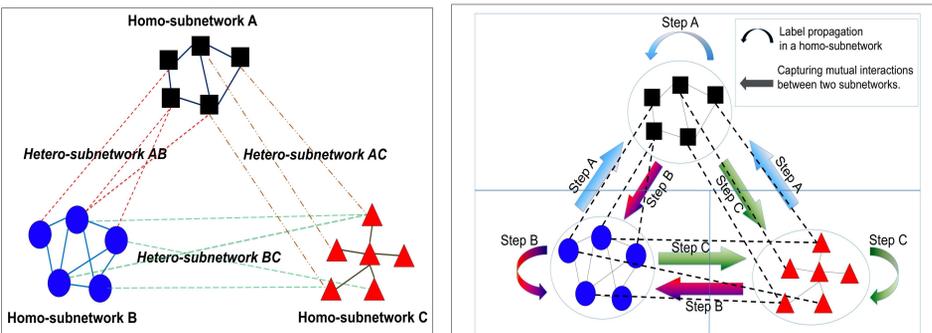


- Objective:**
 - Predict label / activation value to the unlabeled data for classification or ranking
 - Find community structure in the heterogeneous network

Approach

Mutual Interaction based Network Propagation (MINProp)

To handle label propagation on a complex heterogeneous network, MINProp sequentially performs network propagations on each individual homo and hetero-subnetwork



- A heterogeneous network**
This heterogeneous network contains three types of vertices, and accordingly three homo-subnetwork and three heterogeneous network

Running MINProp on a heterogeneous network with three homo-subnetworks:
Repeat until converge
Step A: Perform propagation in homo-subnetwork A with initialization from homo-subnetwork B and C.
Step B: Perform propagation in homo-subnetwork B with initialization from homo-subnetwork A and C.
Step C: Perform propagation in homo-subnetwork C with initialization from homo-subnetwork A and B.

MINProp workflow
Illustration of the MINProp algorithm. Label propagation initialized by the interactions with the other homo-subnetworks is sequentially performed on each individual homo-subnetwork

Regularization framework

$$\Omega(f) = \sum_{i=1}^k (f_i^T \Delta f_i + \mu_i \|f_i - y_i\|^2) + \frac{1}{2} \sum_{i=1}^{k-1} \sum_{j=i+1}^k \mu_{ij} [f_i^T f_j^T] \sum^{(i,j)} \begin{bmatrix} f_i \\ f_j \end{bmatrix}$$

- f : predicted label
- y : initial label
- Δ : graph laplacian of homo - subnetwork
- Σ : graph laplacian of hetero - subnetwork
- $\|f_i - y_i\|^2$: fitting term which keeps the final label values consistent with the initial labels
- $f_i^T \Delta f_i$: smoothness constrain on the homo - subnetwork
- k : number of subnetwork
- μ_i and μ_{ij} : positive constants

that enforces a consistent labeling of the strongly connected vertices in $V^{(i)}$

$$\sum_{i=1}^{k-1} \sum_{j=i+1}^k \mu_{ij} [f_i^T f_j^T] \sum^{(i,j)} \begin{bmatrix} f_i \\ f_j \end{bmatrix}$$

: smoothness term in the hetero - subnetwork

RESULTS

Task

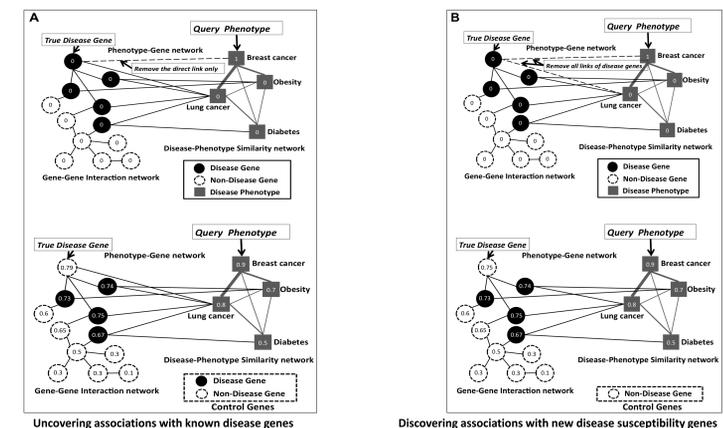
- Given query disease phenotype, we want to highly rank its disease causative genes. The ranking results can be used to discover novel disease genes.**

Datasets

Three heterogeneous network are used in this study.

- Disease similarity network**
: an undirected graph with 5080 disease vertices. Edges are weighted by the pairwise quantitative measurements of phenotypic overlap in text and clinical synopsis of OMIM
- Human protein-protein interaction network**
:ppi-network contains 34,364 binary-valued undirected interactions between 8919 human proteins (genes)
- Disease-gene association network**
:an undirected bipartite graph-with disease vertices and gene vertices. There are 1126 disease-gene associations in the network

Experiments



Performance of ranking disease genes in leave-one-out cross validation

Methods	Associations with known disease genes		Associations with new disease genes	
	Avg. AUC	Avg. AUC (win/draw/loss)	Avg. AUC	Avg. AUC (win/draw/loss)
MINProp vs. Random Walk	0.805	vs. 0.797 (738/75/313)	0.728	vs. 0.648 (1045/2/79)
MINProp vs. CIPHER-DN	0.863	vs. 0.738 (565/5/288)	0.821	vs. 0.738 (515/11/332)
MINProp vs. CIPHER-SP	0.805	vs. 0.734 (678/8/440)	0.728	vs. 0.729 (538/4/534)

Exploring modularity of genes

CC	MINProp vs. RW	MINProp vs. CIPHER-DN	MINProp vs. CIPHER-SP	Hybrid vs. CIPHER-SP
	Avg. AUC	Avg. AUC	Avg. AUC	Avg. AUC
[0.1, 1]	0.875 vs. 0.776	0.889 vs. 0.855	0.875 vs. 0.813	0.886 vs. 0.813
[0.01, 0.1]	0.902 vs. 0.799	0.906 vs. 0.799	0.902 vs. 0.801	0.911 vs. 0.801
[0, 0.01]	0.653 vs. 0.586	0.770 vs. 0.688	0.654 vs. 0.693	0.692 vs. 0.693
Total	0.728 vs. 0.648	0.821 vs. 0.738	0.728 vs. 0.729	0.756 vs. 0.729

MINProp is capable to discover novel disease genes

MM#	Phenotype Name	HGNC Symbol	Ranking (U)	Ranking (U)	Ranking (U)	Status
			MINProp	CIPHER-SP	Random Walk	
17300	PHENOCROMOCYTOMA	PHO	1105 (0.138)	1105 (0.137)	1105 (0.137)	known
60774	MENINGIOMA, FAMILIAL	PTEN	1307 (0.147)	1307 (0.147)	1307 (0.147)	known
60980	RUO-CHABH SYNDROME	PTEN	1307 (0.147)	1307 (0.147)	1307 (0.147)	known
166710	OSTEOPOROSIS	COL1A1	4086 (0.907)	4086 (0.907)	4086 (0.907)	known
60426	LEUKEMIA, ACUTE MYELOID	JAK2	154 (0.044)	154 (0.044)	154 (0.044)	known
20200	ADRENOCORTICAL CARCINOMA	TP53	239 (0.140)	239 (0.140)	239 (0.140)	known
30009	NEPHROPSYCHIASIS	TP53	239 (0.140)	239 (0.140)	239 (0.140)	known
60187	STROKE, ISCHEMIC	PCSK9	448 (0.056)	448 (0.056)	448 (0.056)	new
18478	THYROID CARCINOMA	BRCA1	1182 (0.131)	1182 (0.131)	1182 (0.131)	known
60413	THALASSEMIA	HBA2	443 (0.056)	443 (0.056)	443 (0.056)	new
17780	PROLAPSE SUSCEPTIBILITY 1	ELRN	239 (0.140)	239 (0.140)	239 (0.140)	new
27300	TESTICULAR TUMORS	KIT	347 (0.029)	347 (0.029)	347 (0.029)	new

Future Works

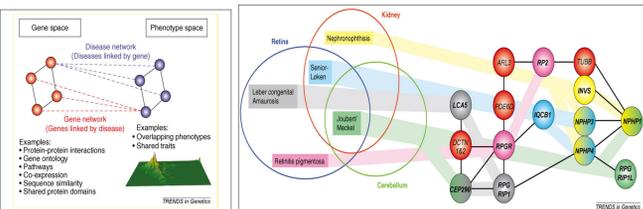
- Integrate other types of genomic and clinical data into our current model to improve the performance.

Network Data Integration perspective:

- Heterogeneity**
: The heterogeneous networks may contain several subnetworks consisted of different types of vertices and edges
: Some vertices are connected to many others, some to very few, and interaction strengths and dynamics may vary widely
- Bias**
: The heterogeneous networks may have subnetworks containing unbalanced sizes, different noisy levels and different edge-weight scales
- Scalability**
: The heterogeneous networks may have more than tens of thousands vertices and edges

Biological Network Data Integration perspective:

- Network View**
: Available genomic data needs to be integrated with "network" view of systems biology
- Exploring modular structure of complex biological systems**
: Need to develop an efficient technique for exploring communities on a large scale
- Interpretability**
: Results from the technique should be interpretable and biologically meaningful.



"Network" view of human disease and gene association.
Each node indicates disease and gene. Phenotypically similar disease can be grouped together in the disease network, and genes have similar functions can be grouped together in gene networks. Disease genes for phenotypically similar diseases may reside in similar biological modules and subnetworks in the protein interaction networks.