Identifying Signaling Pathways

BMI/CS 776 www.biostat.wisc.edu/bmi776/ Spring 2017 Anthony Gitter gitter@biostat.wisc.edu

These slides, excluding third-party material, are licensed under <u>CC BY-NC 4.0</u> by Anthony Gitter, Mark Craven, Colin Dewey

Goals for lecture

- Challenges of integrating high-throughput assays
- Connecting relevant genes/proteins with interaction networks
- ResponseNet algorithm
- Classes of signaling pathway prediction methods

High-throughput screening

- Which genes are involved in which cellular processes?
- Hit: gene that affects the phenotype
- Phenotypes include:
 - Growth rate
 - Cell death
 - Cell size
 - Intensity of some reporter
 - Many others

Types of screens

- Genetic screening
 - Test genes individually or in parallel
 - Knockout, knockdown (RNA interference), overexpression, CRISPR/Cas genome editing
- Chemical screening
 - Which genes are affected by a stimulus?

Differentially expressed genes

- Compare mRNA transcript levels between control and treatment conditions
- Genes whose expression changes significantly are also involved in the cellular process
- Alternatively, differential protein abundance or phosphorylation



Assays reveal different parts of a cellular process



<u>KEGG</u>

Assays reveal different parts of a cellular process



Pathways connect the disjoint gene lists

- Can't rely on pathway databases
- High-quality, low coverage



- Instead learn condition-specific pathways computationally
- Combine data with generic physical interaction networks

Physical interactions

• Protein-protein interactions (PPI)



Appling Graz

- Metabolic
- Protein-DNA (transcription factor-gene)



Genes and proteins are different node types

Hairball networks

- Networks are highly connected
- Can't use naïve strategy to connect screen hits and differentially expressed genes



Yeger-Lotem2009

Identify connections within an interaction network



How to define a computational "pathway"

• Given:

- Partially directed network of known physical interactions (e.g. PPI, kinase-substrate, TF-gene)
- Scores on source nodes
- Scores on target nodes

• Do:

 Return directed paths in the network connecting sources to targets

ResponseNet optimization goals

- Connect screen hits and differentially expressed genes
- Recover sparse connections
- Identify intermediate proteins missed by the screens
- Prefer high-confidence interactions

Construct the interaction network







Weighting interactions

• Probability-like confidence of the interaction

Proteins

0	MP2K1_HUMAN	Homo sapiens	Temporarily not available for viewing in Netility.		
MK01_HUMAN Homo sapiens		Homo sapiens	Temporarily not available for viewing in Netility.		

Evidence

Source DB 🖨	Source ID 🛊	Interaction Type 🖨	PSI MI Code 🖨	PubMed ID 🖨	Detection Type 🖨	PSI MI Code 🖨
biogrid	857930	direct interaction	MI:0407	12788955	enzymatic study	MI:0415
ophid	17231	aggregation	MI:0191	11352917	confirmational text mining	MI:0024
ophid	17231	aggregation	MI:0191	15657099	deglycosylase assay	MI:1006
ophid	17234	aggregation	MI:0191	11352917	confirmational text mining	MI:0024
ophid	17234	aggregation	MI:0191	15657099	deglycosylase assay	MI:1006
biogrid	259225	direct interaction	MI:0407	12697810	t7 phage display	MI:0108
intact	EBI-8279991 🗗	phosphorylation reaction	MI:0217	23241949	biosensor	MI:0968

- Example evidence: edge score of 1.0
- 16 distinct publications supporting the edge

iRefWeb



Find the minimum cost flow



Prefer no flow on the low-weight edges if alternative paths exist



Formal minimum cost flow $\min_{f} \left(\sum_{i \in V', j \in V'} -\log(w_{ij}) * f_{ij} \right) - \left(\gamma * \sum_{i \in Gen} f_{Si} \right) \right)$

Subject to:

$$\sum_{j \in V'} f_{ij} - \sum_{j \in V'} f_{ji} = 0 \quad \forall i \in V' - \{S, T\}$$

Flow coming in to a node equals flow leaving the node

Formal minimum cost flow $\min_{f} \left(\sum_{i \in V', j \in V'} -\log(w_{ij}) * f_{ij} \right) - \left(\gamma * \sum_{i \in Gen} f_{Si} \right) \right)$

Subject to:

 $\sum f_{ij} - \sum f_{ji} = 0 \quad \forall i \in V' - \{S, T\}$ $i \in V'$

 $\sum f_{Si} - \sum f_{iT} = 0$ i∈Gen i∈Tra

Flow leaving the source equals flow entering the target

Formal minimum cost flow $\min_{f} \left(\sum_{i \in V', j \in V'} -\log(w_{ij}) * f_{ij} \right) - \left(\gamma * \sum_{i \in Gen} f_{Si} \right) \right)$

Subject to:

 $\sum_{j\in V'} f_{ij} - \sum_{j\in V'} f_{ji} = 0 \quad \forall i \in V' - \{S, T\}$

$$\sum_{i\in Gen} f_{Si} - \sum_{i\in Tra} f_{iT} = 0$$

Flow is non-negative and does not exceed $0 \le f_{ij} \le c_{ij} \quad \forall (i,j) \in E'$ edge capacity

Formal minimum cost flow $\min_{f} \left(\left(\sum_{i \in V', j \in V'} -\log(w_{ij}) * f_{ij} \right) - \left(\gamma * \sum_{i \in Gen} f_{Si} \right) \right)$

Subject to:

$$\sum_{j\in V'} f_{ij} - \sum_{j\in V'} f_{ji} = 0 \quad \forall i \in V' - \{S, T\}$$

$$\sum_{i\in Gen} f_{Si} - \sum_{i\in Tra} f_{iT} = 0$$

 $0 \leq f_{ij} \leq c_{ij} \quad \forall (i,j) \in E'$

Linear programming

- Optimization problem is a linear program
- Canonical form

maximize $\mathbf{c}^{\mathrm{T}}\mathbf{x}$ subject to $A\mathbf{x} \leq \mathbf{b}$ and $\mathbf{x} \geq \mathbf{0}$ Wikipedia

- Polynomial time complexity
- Many off-the-shelf solvers
- Practical Optimization: A Gentle Introduction
 - Introduction to linear programming
 - <u>Simplex method</u>
 - <u>Network flow</u>

ResponseNet pathways



- Identifies pathway members that are neither hits nor differentially expressed
- Ste5 recovered when STE5 deletion is the perturbation

ResponseNet summary

Advantages

- Computationally efficient
- Integrates multiple types of data
- Incorporates interaction confidence
- Identifies biologically plausible networks
- Disadvantages
 - Direction of flow is not biologically meaningful
 - Path length not considered
 - Requires sources and targets
 - Dependent on completeness and quality of input network

- Unlike PIQ, we don't have a complete gold standard available for evaluation
- Can simulate "gold standard" pathways from a network
- Compare relative performance of multiple methods on independent data
 - Top secret example



<u>Ritz2016</u>



 Natural language processing can also help semi-automated evaluation

• <u>Literome</u>

PMID: 14611643	that PKB mediates the of WNK1 at (details)
WNK1, the kinase mutated in an inherited high-blood-pressure syndrome, is a	
novel PKB (protein kinase B)/Akt substrate.	

<u>Chilibot</u>

 Our studies reveal a novel mechanism in which phosphorylation of STAT3 is mediated by a constitutively active JNK2 [MAPK9] isoform, JNK2 [MAPK9] α. <u>Ref: Oncogene, 2011, PMID: 20871632</u>

• <u>iHOP</u>

Akt1 😭, but not Akt2, phosphorylates palladin 🆙 at Ser507 in a domain that is critical for F-actin bundling. [2010]

Classes of pathway prediction algorithms



Classes of pathway prediction algorithms



Alternative pathway identification algorithms

- k-shortest paths
 - <u>Ruths2007</u>
 - <u>Shih2012</u>
- Random walks / network diffusion / circuits
 - <u>Tu2006</u>
 - eQTL electrical diagrams (<u>eQED</u>)
 - HotNet
- Integer programs
 - Signaling-regulatory Pathway INferencE (<u>SPINE</u>)
 - <u>Chasman2014</u>

Alternative pathway identification algorithms

- Path-based objectives
 - Physical Network Models (<u>PNM</u>)
 - Maximum Edge Orientation (<u>MEO</u>)
 - Signaling and Dynamic Regulatory Events Miner (<u>SDREM</u>)
- Steiner tree
 - Prize-collecting Steiner forest (<u>PCSF</u>)
 - Belief propagation approximation (<u>msgsteiner</u>)
 - <u>Omics Integrator</u> implementation
- Hybrid approaches
 - <u>PathLinker</u>: random walk + shortest paths
 - <u>ANAT</u>: shortest paths + Steiner tree

Recent developments in pathway discovery

- Multi-task learning: jointly model several related biological conditions
 - ResponseNet extension: <u>SAMNet</u>
 - Steiner forest extension: <u>Multi-PCSF</u>
 - SDREM extension: <u>MT-SDREM</u>
- Temporal data
 - ResponseNet extension: <u>TimeXNet</u>
 - Steiner forest extension
 - <u>Temporal Pathway Synthesizer</u> (unpublished)

Condition-specific genes/proteins used as input

- Genetic screen hits (as causes or effects)
- Differentially expressed genes
- Transcription factors inferred from gene expression
- Proteomic changes (protein abundance or posttranslational modifications)
- Kinases inferred from phosphorylation
- Genetic variants or DNA mutations
- Enzymes regulating metabolites
- Receptors or sensory proteins
- Protein interaction partners
- Pathway databases or other prior knowledge