Reconstruction of Metabolic Networks from High Throughput Metabolic Data: In Silico Analysis of RBC metabolism

Ilya Nemenman¹, Michael Wall¹ Sean Escola², William Hlavacek¹

¹Los Alamos National Laboratory ²Columbia University Medical School



UNCLASSIFIED



Metabolic Networks: *(future)* Inference Problem from MassSpec / isotopically labeled data





Steady states because...

- Destructive measurements
- Uncorrelated errors
- Smaller errors (const. sample sz.)
- Less samples, but repeatable
- Only need topologies
- Only relative concentrations

Similar to mRNA arrays data





UNCLASSIFIED

2



Statistical dependency model



 $\frac{dA}{dt} = f(A, B, C, D, Enz)$ f(A, B, C, D, Enz) = 0 $P(A, B, C, D \mid Enz) = \delta(g)$ $P(A, B, C, D) = \langle \delta(g) \rangle = \exp[-\lambda_{ABCD}]$ $P(ABCD) \approx \exp[-\lambda_{AB} - \lambda_{AC} - \dots - \lambda_{CD}]$

Enz

Better model than for mRNA

- Direct coupling of nodes
- Simpler noise model
- Known modulators
- Interactions microscopically pairwise
- No directionality in steady state



UNCLASSIFIED

3



From activity to networks





ARACNE (Califano & Co)





Performance: Few false positives

- No false positives for tree networks
- No false positives under very general conditions for networks with only a few strong loops
- No false negatives under stronger conditions (many otherwise, but it's ok)

UNCLASSIFIED

Need to estimate MI reliably









MI estimation





B-cell dataset: cMYC network

~400 arrays (Dalla-Favera et al.)

No dynamics

~250 naturally occurring, ~150 perturbed

~25 phenotypes (normal, tumors, experimental perturbations)



- Protooncogene,
- 12% background binding,
- one of top 5% hubs
- significant MI with 2000 genes

Total interactions: 56 Pre-known: 22 New Ch-IP validated: 11/12

UNCLASSIFIED

8



Does good microarray performance guarantee good results for metabolites?

- Different noises
- Different nonlinearities
- Very dense
- ~1e7 ratios in kinetic rates/steady state concentrations
 - Interactions of low-abundance metabolites washed out
 - These are essential parts of environmental response pathways (intermediates)
- Need benchmark metabolic data sets (DREAM workshop, NYC, 09/06 http://dimacs.rutgers.edu/Workshops/ReverseEng/)



UNCLASSIFIED

9



Synthetic model

- 39 metabolites
- 44 individual reactions
- 107 pairwise interactions between distinct metabolites





Data sets

- Jamshidi et al. Mathematica code: generate ~1000 steady states with different values for Donnan ratio, glucose, intracellular Pi, Mg, and extracellular Na:
 - chemostat (ranges consistent with survival of RBCs in culture)
 - natural (ranges consistent with normal human blood work)
 - natural correlated (same with correlated parameters)
- Also got ~100hrs of time-dependent data with naturalistic evolution of control parameters
- The chemostat dataset:
 - Smallest mean concentration 4.5e-5
 - Largest mean concentration 1.2e+2
 - Smallest ratio std/mean 1.3e-14
 - Largest ration std/mean 4.5e-1



UNCLASSIFIED

11



Adding noise

Experimental noise simulated by adding additive noise and multiplicative noise

$$X = X_0 + A \cdot \operatorname{randn}() + B \cdot X_0 \cdot \operatorname{randn}()$$

- Remove nodes with std<noise
- Couple all neighbors of removed nodes for validation



UNCLASSIFIED

12



Example





Adjusting ARACNE parameters

- Best kernel: leave-one-out cross validation -- h = 0.1173 of the variable range (rank-order transform).
- To assure <1 falsely significant MI out of 39*38/2 = 741, select threshold corresponding to p-value = 1/741; I₀=0.019.



UNCLASSIFIED

14



Performance on RBC data for different noise levels

PRC for changing noise, I threshold, tolerance



$$p = \frac{N_{TP}}{N_{TP} + N_{FP}} = \frac{N_{TP}}{N_{P, found}}$$
$$r = \frac{N_{TP}}{N_{TP} + N_{FN}} = \frac{N_{TP}}{N_{P, tot}}$$

- Different DPI tolerances (0, 0.05, 0.1 for solid, dashed, dotted).
- Operation point for predetermined / threshold

Operated by the Los Alamos National Security, LLC for the DOE/NNSA



15

Problems

- Low abundance metabolites
- Bootstrapped data sets to increase *r*
- Start with constrained networks (by mass transfer)
- Regulated interactions (metabolic/transcriptional data sets needed)



Most interesting next thing

(detection of enzyme-coding genes)

