Workshop on Data Analysis and Data Mining in Proteomics

Detecting False Positives and False Negatives in Protein Interactome using Network Topology

See-Kiong Ng Knowledge Discovery Department, I²R 12 May 2005





Are we there yet?

| | Coverage | Data quality |
|--------------------------|---|---|
| DNA genome sequence | 99% of genome sequence | 99.9% correct |
| mRNA profiling | 80-90% of transcripts represented | 90% of spots are good data |
| Protein interaction data | 10-30% of interactions catalogued | 50-70% of interactions are spurious |

False Positives

Unravelling the Protein-Protein Interactome

Part I:

The False Positive Challenge

How reliable are experimentally-derived protein interactions?

Some high-throughput protein interaction experiments, such as the popular Y2H method, have as much as <u>50%</u> false positives.

Sprinzak et al., JMB, 327:919-923, 2003

























Interaction Pathway Reliability

Conjecture:

"An interaction that is associated with an alternate path of reliable interactions is likely to be reliable."

- Biological functions are typically performed by highly interconnected networks of interactions
- Circular contigs are frequently observed in protein interaction and complex data
- Alternate paths are also observed in many real-world networks

Use alternative interaction paths as a measure to indicate functional linkage between the two proteins



Interaction Pathway Reliability

IRAP: Interaction Reliability by Alternate Pathway

A novel **global topology** measure for assessing interaction reliability computationally

$$I\!RAP^{\emptyset}(X\leftrightarrow Y) = \max_{\phi\in\Phi^{\varphi}(X,Y)} \prod_{\{U\leftrightarrow V\}\in\phi} \left(1-\frac{ig^{\emptyset}(U\leftrightarrow V)}{ig^{\varphi}_{\max}}\right)$$

where $ig_{\max}^{\mathcal{G}} = \max\{ig^{\mathcal{G}}(X \leftrightarrow Y) \mid (X \leftrightarrow Y) \in \mathcal{G}\}$ is the maximum *IG1* value in *G*, and $\Phi^{\mathcal{G}}(X, Y)$ is the set of all possible non-reproducible paths between *X* and *Y*, but excluding the direct path *X*« *Y*.







Discovering Biologically Interacting "Cross-Talkers"

- From our co-localization experiment, we also observed that there are 257 non-co-localized protein pairs with high IRAP (>0.95)
- A large proportion (53%) of these cross-talking pairs have the same MIPS functional annotation, suggesting that they could be *biologically interacting cross-talkers* such as those in signal tranduction pathways

| ProteinA | Cellular Localization | ProteinB | Cellular Localization | Functional Pathway | |
|----------|------------------------------|----------|------------------------------|---------------------|--|
| YDR299w | nucleolus-protein | YLR208w | cytoplasm-release of | Vesicular transport | |
| | transport | | transport vesicles from ER | (Golgi network) | |
| YOL018c | endosome, ER- | YMR117c | spindle pole body- | Cellular import | |
| | syntaxin SNARE | | spindle pole component | | |
| YDL154w | nucleus-recombination | YBR133c | cytoplasm- neg. | Meiosis | |
| | | | regulator of kinase | and budding | |
| YGL192w | nucleus-put. Adenosine | YBR057c | cytoplasm-meiosis | Development of | |
| | methyltransferase | | potentially in premeiosis | asco-basido | |
| | for sporulation | | DNA synth | -zygo spore | |
| YDR299w | nucleolous- protein | YPL085w | cytoplasm,ER-veiscle coat | both in vesicular | |
| | transport | | protein interacts cytoplasm, | transport | |
| | | | with sec23p | | |
| YEL013w | vacuole-phosphorylated | YFL039c | cytoskeleton-actin | Protein targeting | |
| | protein which interacts with | | | and budding | |
| | Atg13p for cyto to vacuole | | | | |
| | targeting vacuole targeting | | | | |

Examples of interactions with high IRAP values (≥ 0.95) between non-co-localized proteins ("cross-talkers") involved in the same cellular pathway

False Negatives

Unravelling the Protein-Protein Interactome

Part II:

The False Negative Challenge

Approach

- Use the Alternative Path Model to detect false negatives in an experimental PPI network:
 - Compute an IRAP value for each *unlinked* protein pair in the network
 - If it has a high IRAP value, then it may be a false negative

Experiment

Dataset

- 10,199 non-redundant interactions between 4,336 yeast proteins from MIPS with date Jan. 18, 2005
- "Verified true" interactions in PPI network
 - Ito. core set (833 interactions)
- Apply IRAP on dataset with the 833 true interactions hidden from the program
- IRAP program re-discovered 730 interactions



Results

| Subcellular Localization | Protein Pairs | | | | | | | | | | |
|-----------------------------|---------------|---------|-----------|---------|---------|---------|---------|----------|---------|---------|-----|
| | | YDL1S0W | YL R287 C | YPR173C | YOR275C | YDR259C | YDR311W | YGL 153W | YOR180C | YOR327C | YDL |
| 725 | YDL130W | | 0.97403 | | | 0.97403 | | | | | |
| | YL R287C | | | | | | | | | | |
| 72.6/7.60 | YPR173C | | | | 0.97419 | | | | | | |
| | YOR276C | | | | | | | | | •• | |
| 750 | YDR259C | | | | | | 0.97403 | 0.96154 | | 0.98701 | 0.9 |
| | YDR311W | | | | | | | 0.96138 | | 0.96138 | |
| 760 | YGL153W | | | | | | | | 0.98701 | 0.97419 | 0.5 |
| | YOR180C | | | | | | | | | | |
| 750 | YOR327C | | | | | | | | | | 0.5 |
| | YDL226C | | | | | | | | | | |



Conclusions

- There are latent local & global network "motifs" that indicate the likelihood of protein interactions
- These network "motifs" can be exploited in computation elimination of *false positives* and *false negatives* from high-throughput Y2H experiments and possibly other highly erroneous interaction data
- IRAP is so far the most effective topologically-based computational measure for assessing the reliability (false positives) of proteinprotein interactions detected by high-throughput methods
- IRAP can also discover new interactions (false negatives) not detected in the experimental PPI network



References

J. Chen, W. Hsu, M.L. Lee, S.-K. Ng (2004) "Systematic Assessment of High-Throughput Experimental Data for Reliable Protein Interactions using Network Topology", in Proceedings of the 16th IEEE International Conference on Tools with Artificial Intelligence (ICTAI 2004), November 15-17, Florida, pages 368-372.

J. Chen, W. Hsu, M.L. Lee, and S.-K. Ng (2005) "Discovering Reliable Protein Interactions from High-Throughput Experimental Data using Network Topology", accepted for publication in *Artificial Intelligence in Medicine*.

