

14 PREDICTING PROTEIN INTERACTIONS

- * Prediction challenges
 - + Effect of point mutations
 - + Structure of complexes
 - + All interacting proteins
- * Use structural data to predict complexes.
- * Simply want to see if two proteins will interact or not (docking problem)

1. Docking
which surface of A interacts with which surface of B?

- * ↑ partners,
- + evaluate all possible relative position + orientation
- + Allow for structural changes
- + Measure energy of interaction

2. Problems: Slow, False Positives.

- * no taking into account relative strength of interaction w/ all other proteins.
(may not be the best interaction)

3. Limit Search Space

- * filter BEFORE structural comparisons,
- + avoid expensive compute until confident

4. PRISM, PrePPI

4.1 PRISM idea:

- + limited # of protein "architectures"
- + architectural motifs define interactions
- + find spatially similar regions of protein complements of known interface.

[Predict "hot spots" Deep learning]

4.1.2 Two Parts

- * rigid body structural comparisons of target protein to known PPI interface
- * flexible refinement using docking energy
- * valuation metrics: structural similarity & conservation of binding "hot spots"

4.1.3 Algorithm

- * identify interface of template (dist. cutoff)

- * align surface of query to half interfaces

- * ignore rest of structure, test

structure match

" " at hotspots

sequence match " "

- * flexible refinement

4.1.4 Use PDB for template interaction surfaces

Protein Data Bank

4.2 PrePPI

to the Query Q_A, Q_B

- * find homologous proteins of known structure M_A, M_B .

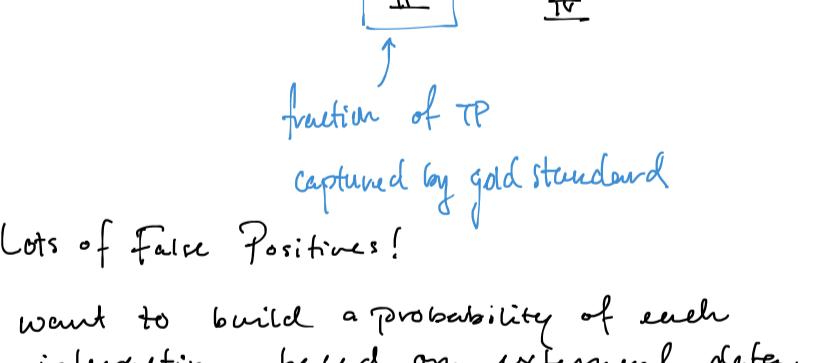
* find structural neighbors $\{N_A, N_B\}$ [$\sim 1500/\text{structure}$]

- * find known interaction among neighbors.

- * align sequences of M_A, M_B to N_A, N_B based on structure

4.2.1 Aligning Sequences.

Template complex



Evaluate on 5 metrics

1. SIM (structural similarity)

2. SIE COV of interaction pairs that can be assigned.

3. Predicts, Pinup, cons-PPISP

* Drawback: doesn't perform flexible structure analysis.

5. Outline

- * High throughput measuring PPI

- * Estimate interaction probabilities.

- * BayesNet prediction of PPI

5.1 "Protein Complexes take the bait"

- * Tag proteins & Mass Spectroscopy

- * Run gel & analyze.

input → protein → tag

False Positives.

- non-specific binding proteins

- proteins binding to tag

False Negatives

- near tag interactions.

- low concentration levels

5.2 TAP-tag (in vivo pull-down) (affinity capture)

[original]

TAP → protein → CFP → protein → NH₂

fusion protein created

in genome.

5.3 Yeast Two Hybrid

DNA binding domain → DBD Bait → AD → reporter gene

transcribe only if AD (activation domain) attached

Cognate Binding Site

detected by

two types of problems:

- * prediction (factors → observation)

- * inference (observation → factors)

7. Estimated INT PPI using other kinds of data

7.1 Gene expression

- * Correlated expression of genes.

7.2 Co-evolution

- * multiple instances of appearance/disappearance of same gene pairs

detected by

two types of problems:

- * prediction (factors → observation)

- * inference (observation → factors)

7.3 Bayesian Model

$$Pr[\text{real PPI} | \text{Data}] = \frac{Pr[\text{data} | \text{true PPI}] Pr[\text{true PPI}]}{Pr[\text{data}]}$$

$$\text{likelihood ratio} = \frac{Pr[\text{true PPI} | \text{data}]}{Pr[\text{false PPI} | \text{data}]} \quad (\text{cancel out } Pr[\text{data}])$$

$$RF = \text{ranking function} = \log \left[\frac{Pr[\text{data} | \text{true PPI}]}{Pr[\text{data} | \text{false PPI}]} \right] \quad (\text{i.e ignore priors})$$

↑ (A, B) pairs.

Protein A, Protein B

→ compute RF_{AB} based on experiments 1, 2, ..., n

7.4 Bayesian Networks.

* in vivo: gene regulation, signaling, prediction

* consists of graph & set of probabilities. ("Learned" from data)

7.5 Predicting interactions.

in vivo pull-down { Gravin, Hb } → protein → tag

Y2Hybrid { Uef1, Ito } → protein → tag

Bayes

not good b/c 2^{n-1} possible relations.

Fully connected

P1, P2 ful

membrane protein

highly expressed

X1, X2, ..., Xn

detected by

two types of problems:

- * prediction (factors → observation)

- * inference (observation → factors)

7.6 Error Rates

Gold Standard

Experiment

FN

FP & TP

Assume:

$$\frac{I}{II} = \frac{III}{IV}$$

fraction of TP captured by gold standard

7.7 Lots of False Positives!

* want to build a probability of each interaction based on external data.

* Bayesian Model

$$Pr[\text{real PPI} | \text{Data}] = \frac{Pr[\text{data} | \text{true PPI}] Pr[\text{true PPI}]}{Pr[\text{data}]}$$

$$\text{likelihood ratio} = \frac{Pr[\text{true PPI} | \text{data}]}{Pr[\text{false PPI} | \text{data}]} \quad (\text{cancel out } Pr[\text{data}])$$

$$RF = \text{ranking function} = \log \left[\frac{Pr[\text{data} | \text{true PPI}]}{Pr[\text{data} | \text{false PPI}]} \right] \quad (\text{i.e ignore priors})$$

↑ (A, B) pairs.

Protein A, Protein B

→ compute RF_{AB} based on experiments 1, 2, ..., n

7.8 Bayesian Networks.

* in vivo: gene regulation, signaling, prediction

* consists of graph & set of probabilities. ("Learned" from data)

7.9 Predicting interactions.

in vivo pull-down { Gravin, Hb }

Y2Hybrid { Uef1, Ito }

Bayes

not good b/c 2^{n-1} possible relations.

Fully connected

P1, P2 ful

membrane protein

highly expressed

X1, X2, ..., Xn

detected by

two types of problems:

- * prediction (factors → observation)

- * inference (observation → factors)

7.10 Estimated INT PPI using other kinds of data

7.11 Gene expression

* Correlated expression of genes.

7.12 Co-evolution

* multiple instances of appearance/disappearance of same gene pairs

of same gene pairs