# A knowledge-based $T^2$-statistic to perform pathway analysis for quantitative proteomic data

June En-Yu Lai

TIGP Bioinformatics Program, Academia Sinica, Taiwan
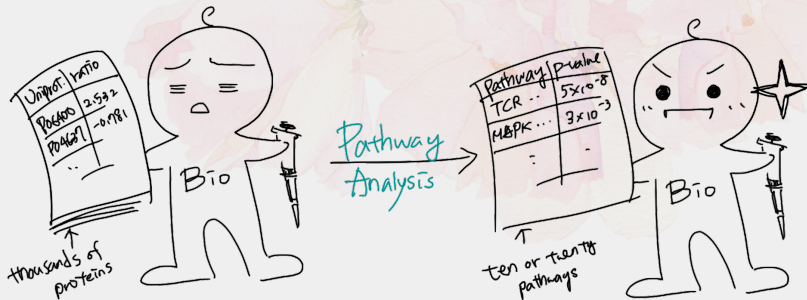
August 18, 2017

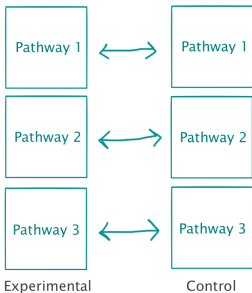# Functional Analysis for Quantitative Proteomic Data

- Proteomic data v.s. gene expression data
    - **Smaller sample size** (number of experiments).
    - Fewer identified entities.
    - The results are sensitive to experimental conditions and instruments.
- Functional analysis
    - **Pathway analysis** (PLoS Computational Biology, 2017).
    - Responsive subpathway locating (working manuscript).

# Null hypothesis: **Self-contained** v.s. Competitive

# The $T^2$-statistic for pathway analysis

The proposed $T^2$-statistic for a specific pathway $\mathcal{P}$ is then defined as,

$$T^2 = \mathbf{x}^T \mathbf{S}^{-1} \mathbf{x} \sim \chi_q^2$$

where

    $\mathbf{x}$  is the vector of expression ratios,

    $\mathbf{x}^T$  is the transpose of $\mathbf{x}$,

  $\mathbf{S}^{-1}$  is the inverse of **the covariance matrix $\mathbf{S}$**, and

    $q$  is the number of mapped proteins in $\mathcal{P}$.

# The $T^2$-statistic for pathway analysis

- We use the **confidence score** provided by protein-protein interaction databases to represent the strength of the covariance, and the **expression direction** provided by the testing dataset to indicate the sign of the covariance.

- On the basis of the confidence scores in $\mathfrak{I}$ and the protein expression ratios, each element $s_{ij}$ of $\boldsymbol{S}$ is determined by the following four rules:

$$s_{ij} = \begin{cases} 0.4 & \text{if } i = j. \\ c_{p_i p_j} & \text{if } i \neq j, \, c_{p_i p_j} \in \mathfrak{I}, \text{ and } x_i \cdot x_j \geq 0. \\ -c_{p_i p_j} & \text{if } i \neq j, \, c_{p_i p_j} \in \mathfrak{I}, \text{ and } x_i \cdot x_j < 0. \\ 0.0 & \text{if } i \neq j \text{ and } c_{p_i p_j} \notin \mathfrak{I}. \end{cases}$$

# The $T^2$-statistic for pathway analysis

$$T^2 = \mathbf{x}^T \mathbf{S}^{-1} \mathbf{x} \sim \chi_q^2$$

- If $\mathbf{S}$ is degenerate, we construct a Moore-Penrose pseudoinverse of $\mathbf{S}$ as a substitute, and $q$ becomes the rank of $\mathbf{S}$.
- The $p$-value of the pathway $\mathcal{P}$ is derived from the $\chi_q^2$ distribution.

# Pathway integration

# Performance evaluation — Really HARD!!!

- No accepted gold-standards.
- We tried to match the results reported by these methods to the biological ideas provided by the original publication (true positives).
- A good statistic should be able to test if a pathway is significant:
  - if a statistic reports very little significant pathways, it might have the problem of false negatives;
  - if a statistic reports a large number of significant pathways, it might have the problem of false positive.
- We supposed that the number of significant pathways is one of the attribute to evaluate these methods.

# General comparison

| Tools | Significance requirement | Ranking statistic |
|-------|--------------------------|-------------------|
| $T^2$ | raw $p$-value $\leq 0.05$ | number of mapped proteins |
| DPA | raw $p$-value $\leq 0.05$ | $p$-value |
| GSEA | FDR adjusted $p$-value (i.e. $q$-value) $\leq 0.25$ | NES |
| DAVID | EASE $\leq 0.1$ | $p$-value |
| IPA | Benjamini corrected $p$-value $\leq 0.05$ | $p$-value |

# General comparison: Consistency

# General comparison: Consistency

| Database | KEGG | | | | | Reactome | | | | | IPA |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | $T^2{\times}ST$ | $T^2{\times}HP$ | DPA | GSEA | DAVID | $T^2{\times}ST$ | $T^2{\times}HP$ | DPA | GSEA | DAVID | IPA |
| Uniprot | + | + | − | − | + | + | + | + | − | − | − |
| IPA | + | + | − | − | + | + | + | + | − | − | − |
| STRING | + | + | − | − | + | + | + | + | − | + | − |
| HitPredict | + | + | − | − | + | + | + | + | − | + | − |
| ST: low | + | + | − | − | + | + | + | + | − | − | − |
| ST: medium | + | + | − | − | + | + | + | + | − | − | − |
| ST: high | + | + | − | − | + | + | + | + | − | + | − |
| ST: highest | + | + | − | − | + | + | + | + | − | + | − |
| HP: low | + | + | − | − | + | + | + | + | − | − | − |
| HP: high | + | + | − | − | + | + | + | + | − | − | − |

# General comparison: Target pathways

| Dataset | TCR | | | PKA | | Myogenesis | | CML | | MAPK |
|---|---|---|---|---|---|---|---|---|---|---|
| Treatment | $\alpha$-CD3$\epsilon$ | | | PGE2 | | Serum-free | | Dasatinib | | U0126 |
| | 5 min | 15 min | 60 min | 1 min | 60 min | 24 hr | 72 hr | 5 nM | 50 nM | 10 $\mu$M |
| KEGG pathway | T cell receptor signaling pathway | | | cAMP signaling pathway | | ECM-receptor interaction | | Chronic myeloid leukemia | | MAPK signaling pathway |
| $T^2 \times$ ST  *p*-value | 1/13  < 0.0001 | 1/57  < 0.0001 | 2/36  0.0019 | 16/45  < 0.0001 | 17/47  < 0.0001 | 15/54  < 0.0001 | 16/75  < 0.0001 | 20/111  < 0.0001 | 23/119  < 0.0001 | 3/118  < 0.0001 |
| $T^2 \times$ HP  *p*-value | 1/14  < 0.0001 | 1/59  0.0002 | 2/17  0.0011 | 17/49  < 0.0001 | 18/51  < 0.0001 | 16/56  < 0.0001 | 15/70  < 0.0001 | 20/113  < 0.0001 | 25/121  < 0.0001 | 3/117  < 0.0001 |
| DPA  *p*-value | - | 8/68  0.0007 | 20/71  0.0024 | 11/15  0.0414 | 7/19  0.0178 | 18/73  0.0004 | - | - | - | 2/17  0.0007 |
| GSEA  *q*-value | 68/69  0.1786 | 7/60  0.1039 | - | - | - | 22/31  0.1416 | - | - | - | - |
| DAVID  *p*-value | 1/47  < 0.0001 | 1/53  < 0.0001 | 1/35  < 0.0001 | - | - | 40/73  0.0041 | 40/73  0.0041 | 8/54  0.0004 | 8/54  0.0004 | 8/119  0.0002 |
| IPA  *p*-value | 1/225  < 0.0001 | 2/232  < 0.0001 | 89/185  0.0022 | 60/79  0.0186 | 60/79  0.0191 | - | - | 52/129  0.0013 | 52/129  0.0013 | - |

# Case Study: TCR downstream phosphoproteome

| Dataset | Pathway Title | $T^2 \times$ ST | $T^2 \times$ HT | DPA | GSEA | DAVID | IPA |
|---------|---------------|-----------------|-----------------|-----|------|-------|-----|
| 5 min | T cell receptor signaling pathway | 1/13 | 1/14 | - | 68/69 | 1/47 | 1/225 |
| 15 min | Ras signaling pathway | 2/57 | 2/59 | - | - | 10/53 | - |
| | Regulation of actin cytoskeleton | 3/57 | 3/59 | - | - | 9/53 | 15/232 |
| | MAPK signaling pathway | 4/57 | 4/59 | 65/68 | - | 41/53 | 4/232 |
| | PI3K-Akt signaling pathway | 24/57 | 25/59 | - | 6/60 | - | 169/232 |
| 60 min | mTOR signaling pathway | 25/36 | - | 67/71 | - | 24/35 | 68/185 |

# Accurate and inaccurate estimation: Toy example



|                                          | Correlated Data (generated by $S$)           | Independent Data (generated by ï¡ð$I$)        |
| ---------------------------------------- | -------------------------------------------- | -------------------------------------------- |
| Correlated Null (normalized by $S$)      | M2: accurate estimation                      | M4: inaccurate due to false positive PPI scores |
| Independent Null (normalized by $I$)     | M3: inaccurate due to incomplete knowledge   | M1: accurate estimation                      |

# Accurate and inaccurate estimation: Real cases

| Dataset | Experiment | Pathway Title | original *p*-value | 30% permuted | 60% permuted | 30% purged | 60% purged |
|---------|-----------|---------------|---------|---------|---------|---------|---------|
| TCR | 5 min | T cell receptor signaling pathway | < 0.0001 | 100% | 100% | 100% | 100% |
| | 15 min | Ras signaling pathway | 0.0017 | 100% | 100% | 100% | 100% |
| | | Regulation of actin cytoskeleton | < 0.0001 | 100% | 100% | 100% | 100% |
| | | MAPK signaling pathway | < 0.0001 | 100% | 100% | 100% | 100% |
| | | PI3K-Akt signaling pathway | < 0.0001 | 100% | 100% | 100% | 100% |
| | 60 min | mTOR signaling pathway | 0.0002 | 97% | 98% | 100% | 100% |
| PKA | 1 min | Regulation of actin cytoskeleton | < 0.0001 | 100% | 100% | 100% | 100% |
| | | PI3K-Akt signaling pathway | < 0.0001 | 92% | 81% | 86% | 84% |
| | | MAPK signaling pathway | < 0.0001 | 100% | 100% | 100% | 100% |
| | | Rap1 signaling pathway | < 0.0001 | 100% | 100% | 100% | 100% |
| | | cAMP signaling pathway | < 0.0001 | 100% | 100% | 100% | 100% |
| | | Glycolysis / Gluconeogenesis | 1 | 0% | 0% | 0% | 0% |
| | 60 min | Cell cycle | < 0.0001 | 100% | 100% | 100% | 100% |
| | | mTOR signaling pathway | 0.0024 | 100% | 100% | 100% | 100% |
| | | Base excision repair | < 0.0001 | 100% | 100% | 100% | 100% |

Thank You♥