# Static network inference (Barcelona variant)

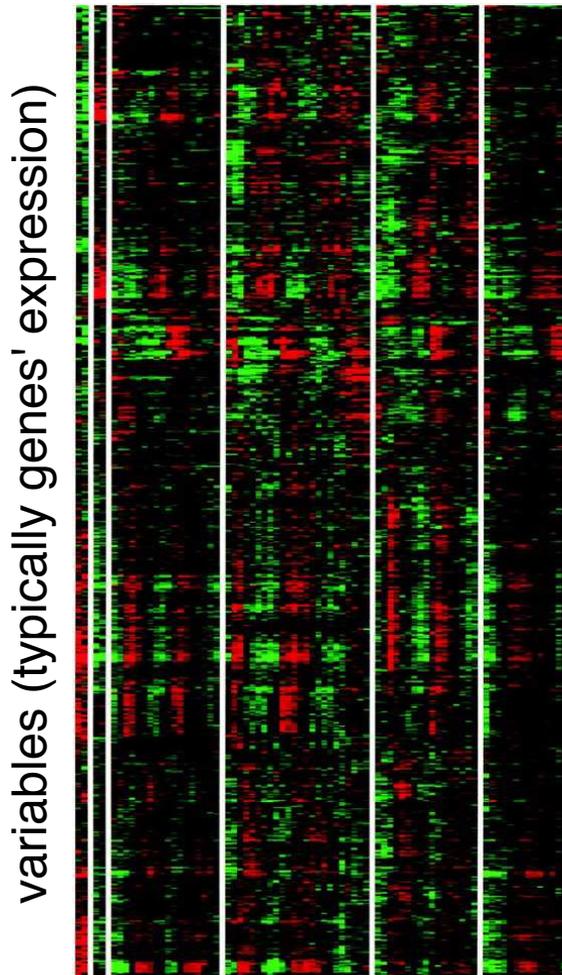BCH 5101: Analysis of -Omics Data

# Why network inference?

# Overview

- "Correlation" networks from expression – similar genes are connected

  - Correlation
  - Information theory & mutual information
  - Permutation testing
  - Relevance Networks
  - ARACNE

- Genetic interaction networks – when deleting two genes doesn't give what you'd expect

  - Synthetic interaction lethality
  - Avery & Wasserman's qualitative framework for epistasis analysis
  - A quantitative statistical model for epistasis analysis

# "Correlation" networks

Typically, we start with a data matrix measuring the expression of genes under different conditions.

conditions

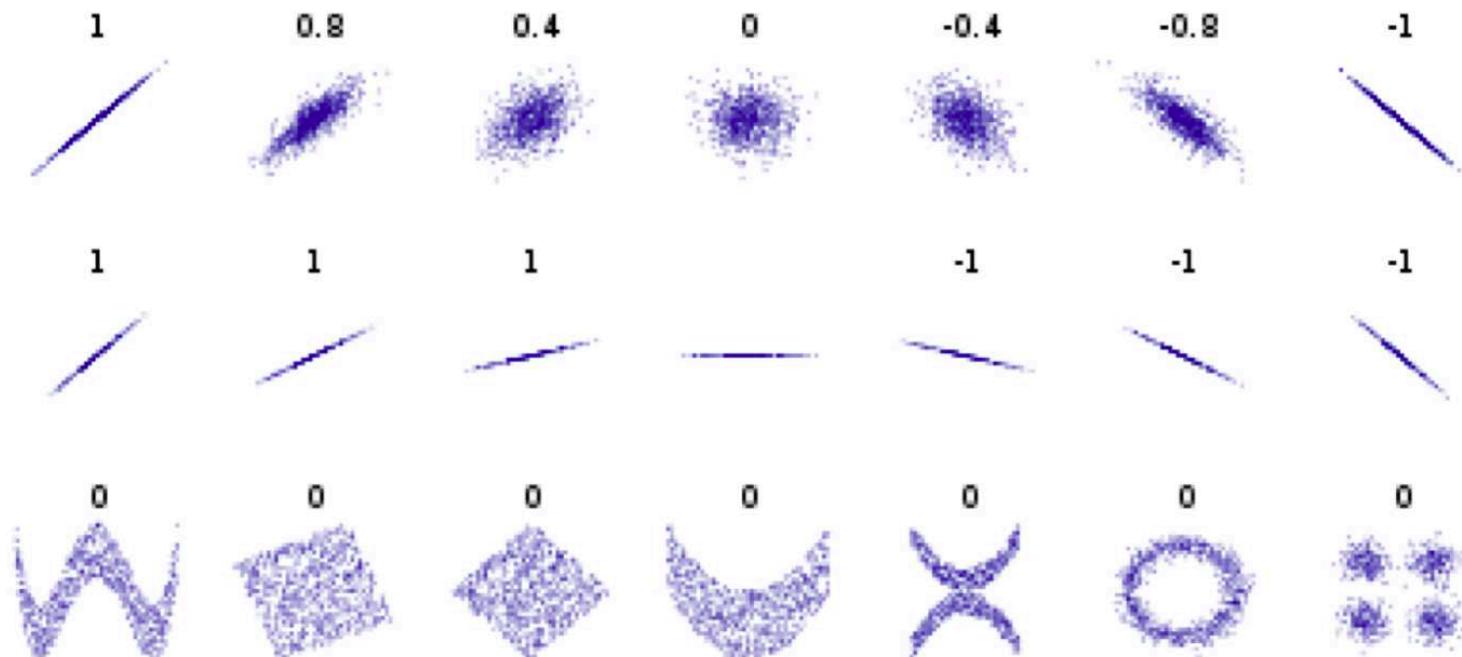variables (typically genes' expression)

- Main idea: Make a big graph in which "similarly" expressed genes are connected.

- Could represent one TF regulating another, or co-regulated genes in a complex / pathway, or any number of other things . . . .

- The resulting graph can then be inspected / analyzed to extract biological meaning.

- What does "similar" mean?

- When are two variables similar enough?

# Pearson's linear correlation coefficient

Linear correlation between $(x_1, x_2, \ldots, x_n)$ and $(y_1, y_2, \ldots, y_n)$:

$$r(x,y) = \frac{\mathsf{Cov}(x,y)}{\sqrt{\mathsf{Var}(x)\mathsf{Var}(y)}} = \frac{\sum_{i=1}^{n} \frac{1}{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n} \frac{1}{n}(x_i - \bar{x})^2 \sum_{i=1}^{n} \frac{1}{n}(y_i - \bar{y})^2}}$$

where $\bar{x}$ and $\bar{y}$ are sample means.

# Lit. Example: Relevance Networks (Butte et al. PNAS 2000)

1. Remove variables with "low information content" – e.g., genes that are always on or always off.

2. For every pair of variables $(i, j)$ compute the Pearson's (linear) correlation coefficient across the conditions, $r_{ij}$, and then compute:

$$\hat{r}_{ij}^2 = \frac{r_{ij}}{|r_{ij}|} r_{ij}^2 = sgn(r_{ij}) r_{ij}^2$$

3. Choose a threshold, $\tau$, to determine statistically significant values of $\hat{r}_{ij}^2$

4. Connect nodes $i$ and $j$ with an undirected edge, if $\hat{r}_{ij}^2 > \tau$

# How to choose $\tau$? Permutation testing for significance

- Suppose you've got paired data on two variables, $x$ and $y$:
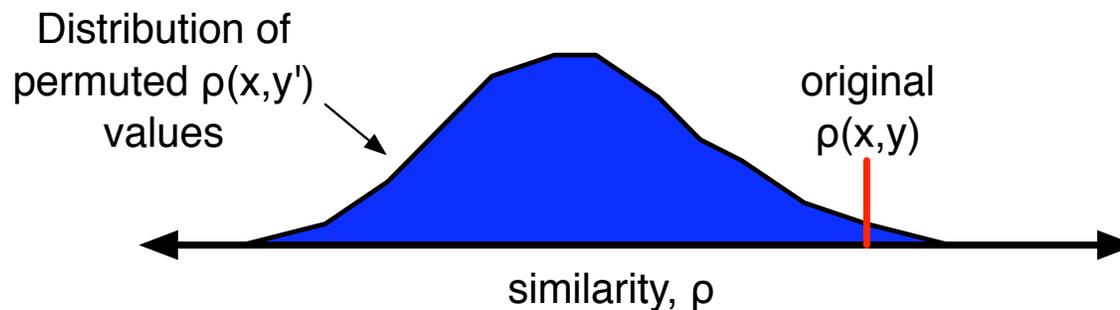
| $x$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | $x_7$ | $x_8$ |
|---|---|---|---|---|---|---|---|---|
| $y$ | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ | $y_8$ |

- Suppose you've got any measure of similarity $\rho$, which assigns a score to such paired data, $\rho(x, y)$.
- $N$ times, randomly permute the $y$ values and recompute $\rho$. E.g.:

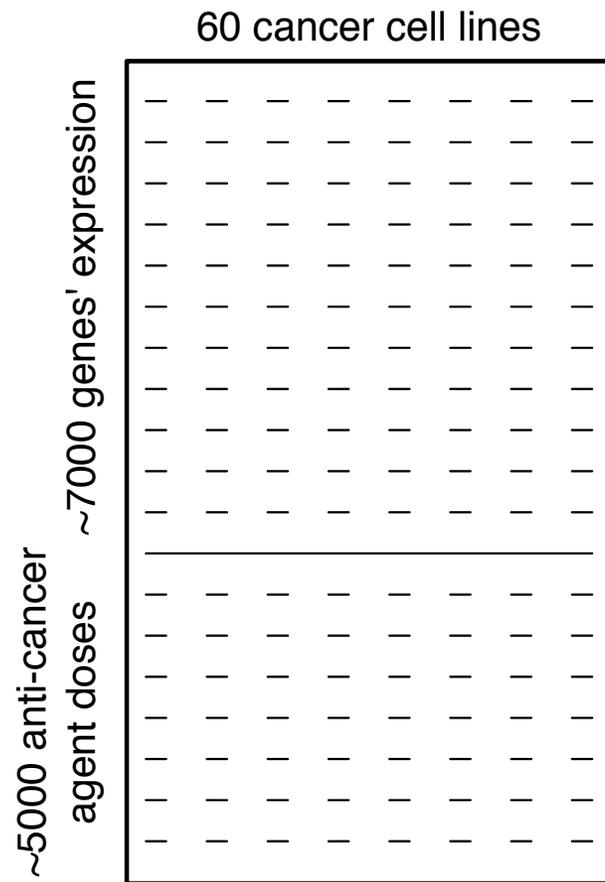| $x$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | $x_7$ | $x_8$ |
|---|---|---|---|---|---|---|---|---|
| $y'$ | $y_3$ | $y_7$ | $y_6$ | $y_1$ | $y_5$ | $y_4$ | $y_2$ | $y_8$ |

- The location of the original $\rho(x, y)$ with respect to the permuted $\rho$ values gives a p-value.
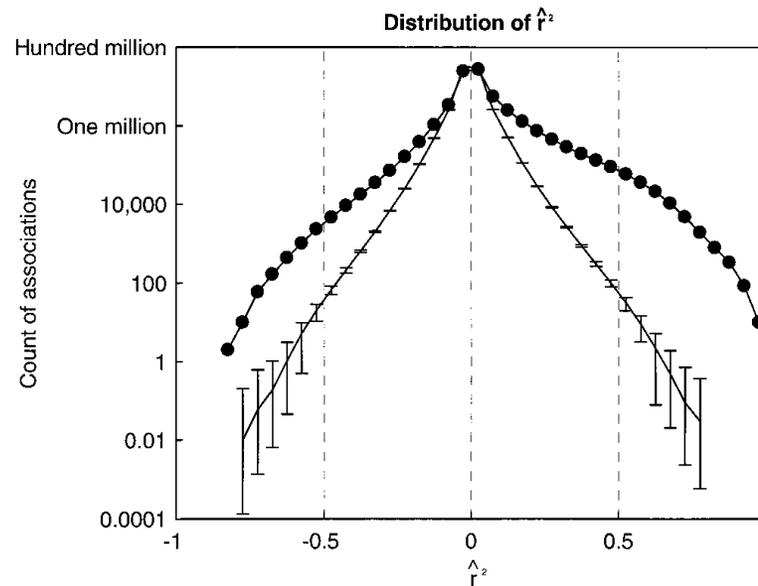


- Approach is distribution and similarity-measure neutral! Still need to choose a p-value threshold (or FDR)...

# Back to Butte et al.

Butte *et al.* took 60 cancer cell lines as conditions, and used as variables microarray expression data for 6,701 genes and susceptibility of those lines to 4,991 anti-cancer agents.



- 544 genes and 93 anti-cancer agents discarded as being low information
- 68,345,586 pairwise correlations computed
- $\tau = 0.8$ was deemed statistically significant

# Network found



Found 1222 links among
834 genes and anti-cancer agents
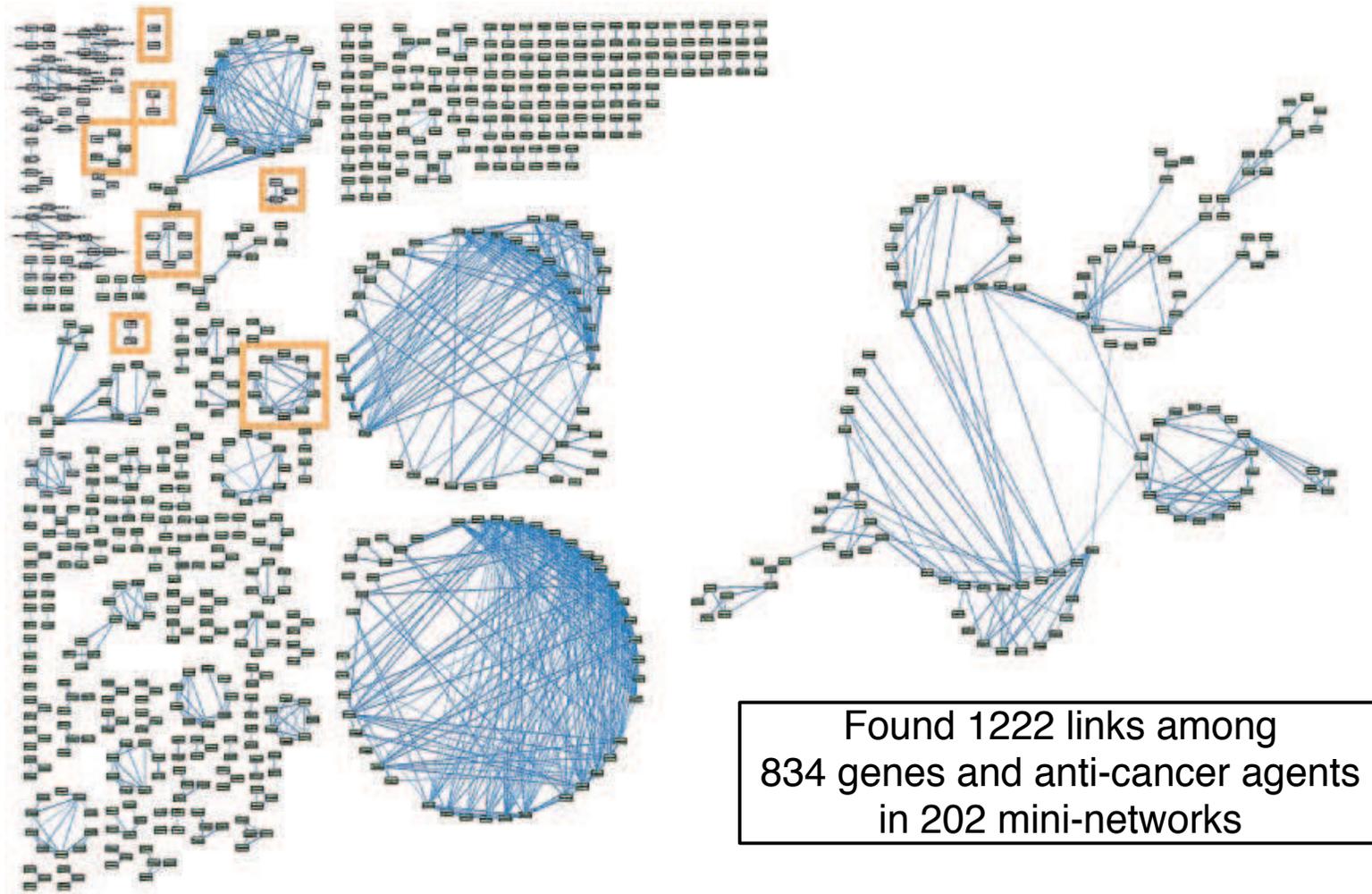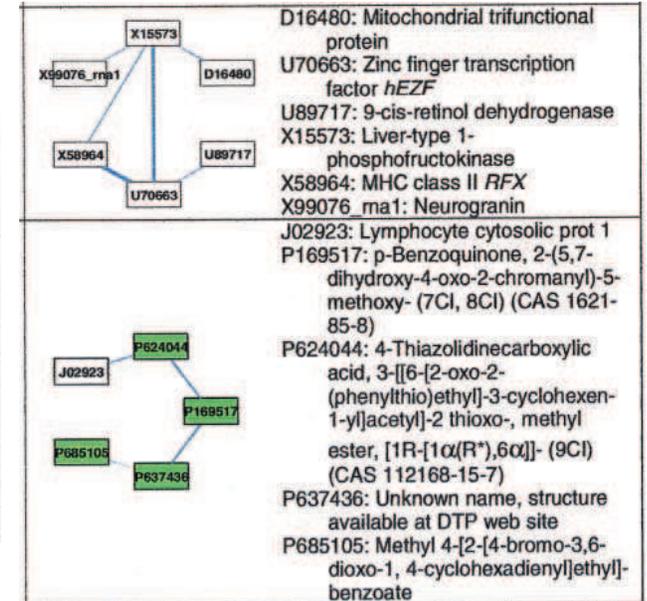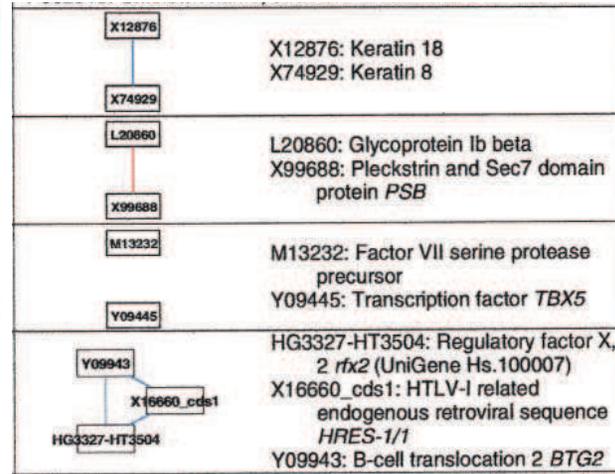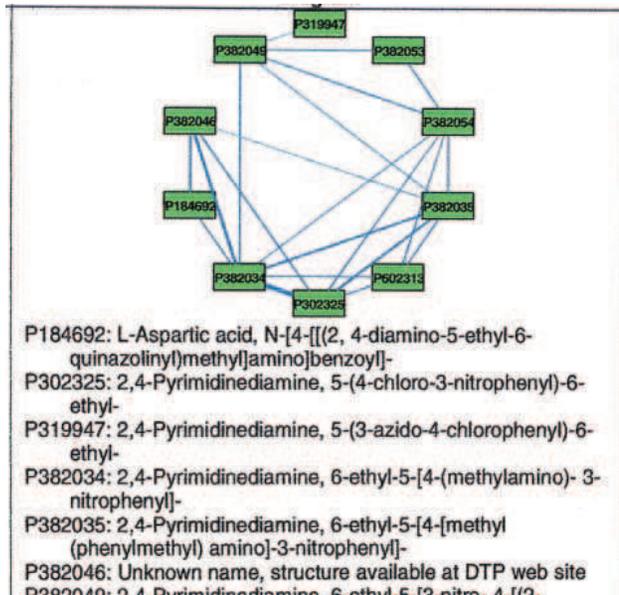in 202 mini-networks

Fig. 2. Relevance networks constructed from the joined databases of baseline gene expression in 60 cancer cell lines and measures of susceptibility of the same cell lines to anticancer agents. The pairs of features (anticancer agents in green boxes, genes in white boxes) with $r^2$ at or greater than $\pm 0.80$ were drawn with line thickness proportional to $r^2$. Features without an association at $\pm 0.80$ were removed. Associations with negative $r^2$ are in red. Seven networks are highlighted in orange and are in Table 1. Large versions of all figures and descriptions for each accession number may be found at http://www.chip.org/genomics.

# Some subnetworks



- Found many links between genes and between anti-cancer agents

- Only one gene was linked to an anti-cancer agent – possibly meaningful

- Their $\tau$ is very conservative (IMHO)

- Should they compute a different $\tau$ for each pair of variables?

$\Rightarrow$ Analysis of subnetworks leads to generation of hypothesis; the Relevance Networks algorithm has been used in many subsequent studies.