

Using Mutual Information to Quantify Toxicological Relationships in Quantitative High Throughput Screening Data

Sankar, Akshay¹; Shockley, Keith²

¹University of North Carolina – Chapel Hill, ²NIEHS Biostatistics and Computational Biology Branch

Abstract

BACKGROUND: Quantitative high throughput screening (qHTS) is an *in vitro* approach used to simultaneously test many chemicals over a broad range of concentrations to better assess chemical toxicity. Large qHTS datasets have several unidentified toxicological relationships. Mutual information (MI) is an information theory statistic used to describe the dependence between random variables. MI can describe both linear and non-linear relationships between chemicals in units of bits (0 or 1), where the value of MI represents the number of possible subgroups (2^{MI}) of paired responses between two chemicals. An MI of 0 bits describes a statistically independent relationship between chemicals.

OBJECTIVE: To develop an MI-based approach to describe toxicological relationships between biologically active chemicals in the BG1 ER agonist assay (phase II of Tox21).

METHODS: A total of 8,306 chemicals surveyed in the estrogen receptor agonist assay were filtered to 1,093 active chemicals. Interpolating splines were used to account for missing data and different concentration spacing. MI estimates were compared to Pearson correlation. Significance was based on permutation testing.

RESULTS: Observed MI values ranged from 0 bits to 3.9-bits across all 596,778 pairwise comparisons. In general, MI estimates captured linear and non-linear relationships between chemicals. A total of 80,333 and 19,458 significant chemical pairs ($\alpha = 5\%$) were found in the BG1 ER agonist assay using Pearson and MI, respectively.

CONCLUSIONS: MI can be used to quantify non-linear relationships between chemicals and find associations between chemicals not identified by Pearson correlation. Splines can be used to correct for missing data and offset concentrations ranges between chemical profiles. A scaled MI value, ranging from 0 to 1, may be useful to compare MI values of different chemical pairs.

Methods



- A total of 8,306 chemicals in the BG1 ER agonist assay were filtered to 1,093 active chemicals with homogeneous profiles.
- Interpolating splines were used to estimate response values along each individual profile and account for differences in concentration spacing.
- Response values were discretized into three standard deviation intervals of the negative controls. MI estimates were computed using discretized values according to Equation (1).
- The significance of the observed MI and Pearson correlation values were obtained via permutation testing after shuffling the BG1 ER agonist data set 30 times.
- Using MI values, the largest cliques, maximal and complete subgraphs, of significant pairwise comparisons ($\alpha = 5\%$) were analyzed with the Tox21 Enricher tool [http://hurlab.med.und.edu/tox21enricher/] to identify over-represented KEGG pathways.

Relevant Equations

$$H(X) = - \sum_{i=1}^n p(x_i) \log_2 p(x_i) \quad (1)$$

$$H(X|Y) = - \sum_{i=1}^n \sum_{j=1}^m p(x_i, y_j) \log_2 p(x_i|y_j) \quad (2)$$

$$I(X;Y) = H(X) - H(X|Y) = \sum_{i=1}^n \sum_{j=1}^m p(x_i, y_j) \log_2 \frac{p(x_i, y_j)}{p(x_i)p(y_j)} \quad (3)$$

$$I(X;Y) = I(Y;X) \geq 0 \quad (4)$$

- The mutual information, $I(X;Y)$, is the reduction in entropy (or uncertainty) of one variable (e.g., one chemical) given observations of the other variable (or a different chemical).
- Mutual information increases monotonically, is symmetric and is 0 when the random variables are statistically independent.
- An MI of 0 implies the random variables are statistically independent [$P(x,y) = P(x)P(y)$].

Conclusions

The MI-based approach:

- Captures about 16,000 associations between chemicals
- Accounts for missing data and differences in concentration spacing
- Is scale invariant and translation variant
- Nonparametric and symmetric

...with some caveats:

- MI fails to find significant associations between chemicals with uninformative (e.g. flat-line) concentration-response profiles
- Data discretization (3 standard deviations of negative controls) may not be optimum

Assessing the significance of MI estimates is difficult:

- Different pairwise comparisons might have different possible maximal MI values
- No probability distributions are available to reliably characterize the significance of small and large MI values
- P-value correction for the multiple testing problem was not addressed

Introduction

Quantitative high throughput screening (qHTS) produces concentration-response data for thousands of chemicals simultaneously and has applications in drug discovery and toxicity testing (Collins et al., 2008).

In Tox21, about 10,000 chemicals are tested at 15 concentration levels to more accurately assess toxicological response. qHTS datasets will be generated for 200-300 different assays (Shockley, 2015).

Describing toxicological relationships in qHTS data may help toxicologists prioritize chemicals for further study or, ultimately, predict toxicity.

Entropy (H) represents how much information (or uncertainty) is conveyed by the probability distribution of the chemical profile (e.g., chemical X or chemical Y). Mutual information (I) is a measure of the amount of information shared between chemical profiles (Shannon, 1948).

Mutual information is a nonparametric measure of dependence that can describe linear and non-linear toxicological relationships.

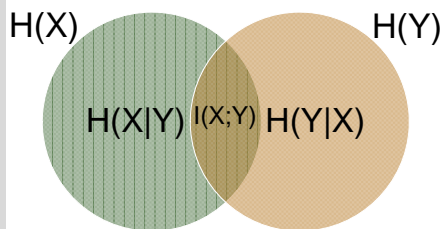


Figure 1: MI vs. Pearson Correlation

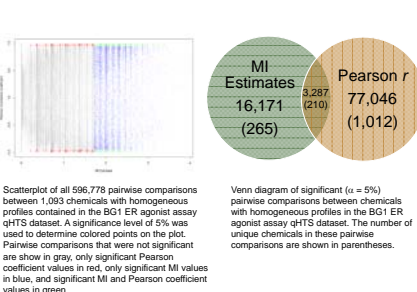
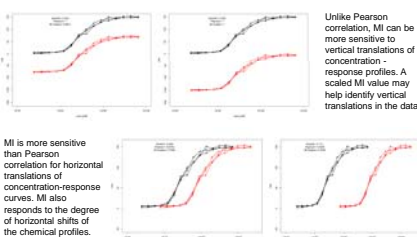


Figure 3: Shifted Profiles



Objective

The underlying concentration-response profiles in qHTS data vary widely across different chemicals. Thus, the form of toxicological relationships between chemicals in qHTS data are not known: they may be linear, exponential, quadratic etc.

We sought to develop a nonparametric approach to:

- identify relationships between chemicals (linear or nonlinear)
- determine the significance of these associations
- compare our new measure with Pearson correlation

Figure 2: Experimental qHTS Examples

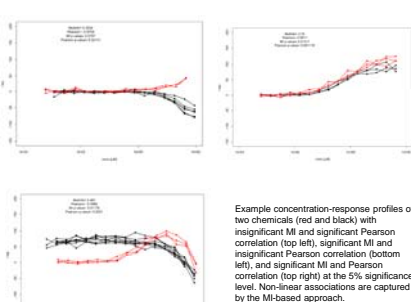
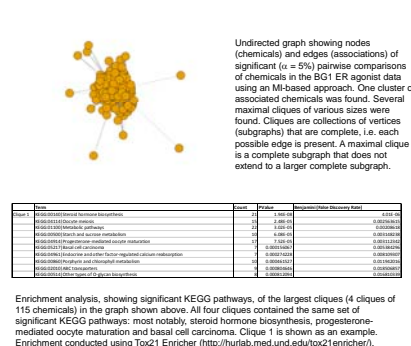


Figure 4: Tox21 Enrichment Analysis



Future Directions

- Q1: Can we use Markov models to relax the independent and identically distribution (iid) assumption and more accurately measure MI?
- Q2: Can we optimize discretization procedures to maximize the mutual information between two chemical profiles?
- Q3: Can we develop methods to reliably compare and assess the significance of MI observations?
- Q4: How can MI be extended to describe toxicological relationships between three or more chemicals?

References

- Collins FS, Gray GG and Bucher JR. (2008). Transforming environmental health protection. Science. 319: 906-907.
- Shockley KR. (2015). Quantitative high-throughput screening data analysis: challenges and recent advances. Drug Discovery Today 20: 296-300.
- Shannon CE. (1948). A mathematical theory of communication. Bell System Technical Journal. 27: 379-423, 623-656.

Acknowledgments

Summer Internship Program (NIEHS)

Branch of Computational Biology and Biostatistics (NIEHS)

Tox21 (NTP) for BG1 ER agonist data: <https://ntp.od.nih.gov/tox21>