

# A Blind Date with Big Data

## Goal

- Measure correlations
- Rank relatedness

## Challenges

- Conditional dependencies
- Missing data



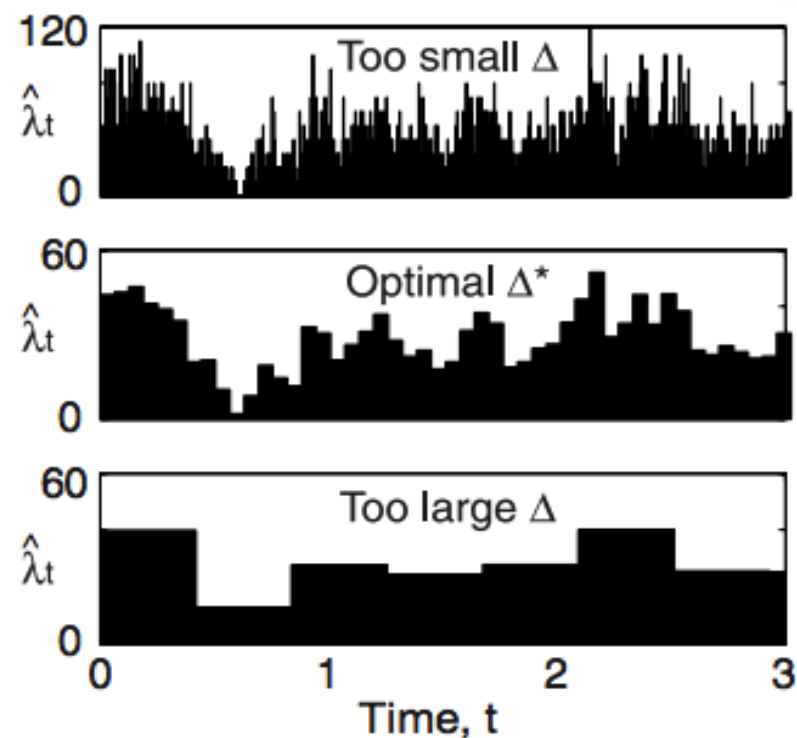


# Where We Are

## Correlation Calculations

- Mutual Information
- Optimal Histograms
- Efficient Algorithm
  - Ranks NHANES in hours

$$MI_{XY} = \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} \frac{n_{(i,j)}}{n_{(:,j)}} \log \frac{n_{(:,j)} * n_{(i,j)}}{n_{(i,:)} * n_{(:,j)}}$$



# Where We're Going

## Imputation

- Developing new method based on optimal bins

## Correlations

- Explore Kernel Density Estimation

## Conditionals

- Probabilistic Graphical Models

## Personal Goals

- Practice programming with variety of languages
- Study statistics, bioinformatics, information theory
- Learn about all projects of the Sabeti Lab
- Have fun!