

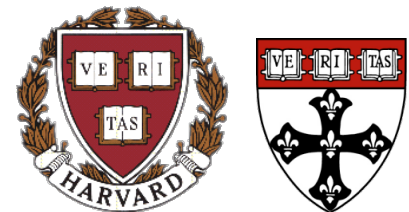
High-Sensitivity Pattern Discovery in High-Dimensional Heterogeneous Datasets

Gholamali (Ali) Rahnavard

Postdoctoral Associate/ The Huttenhower Lab

SAMSI Bioinformatics Workshop: Discovering Patterns in
Human Microbiome Data (HMD)

17 March 2015



Harvard T.H. Chan School of Public Health
Department of Biostatistics



Outline

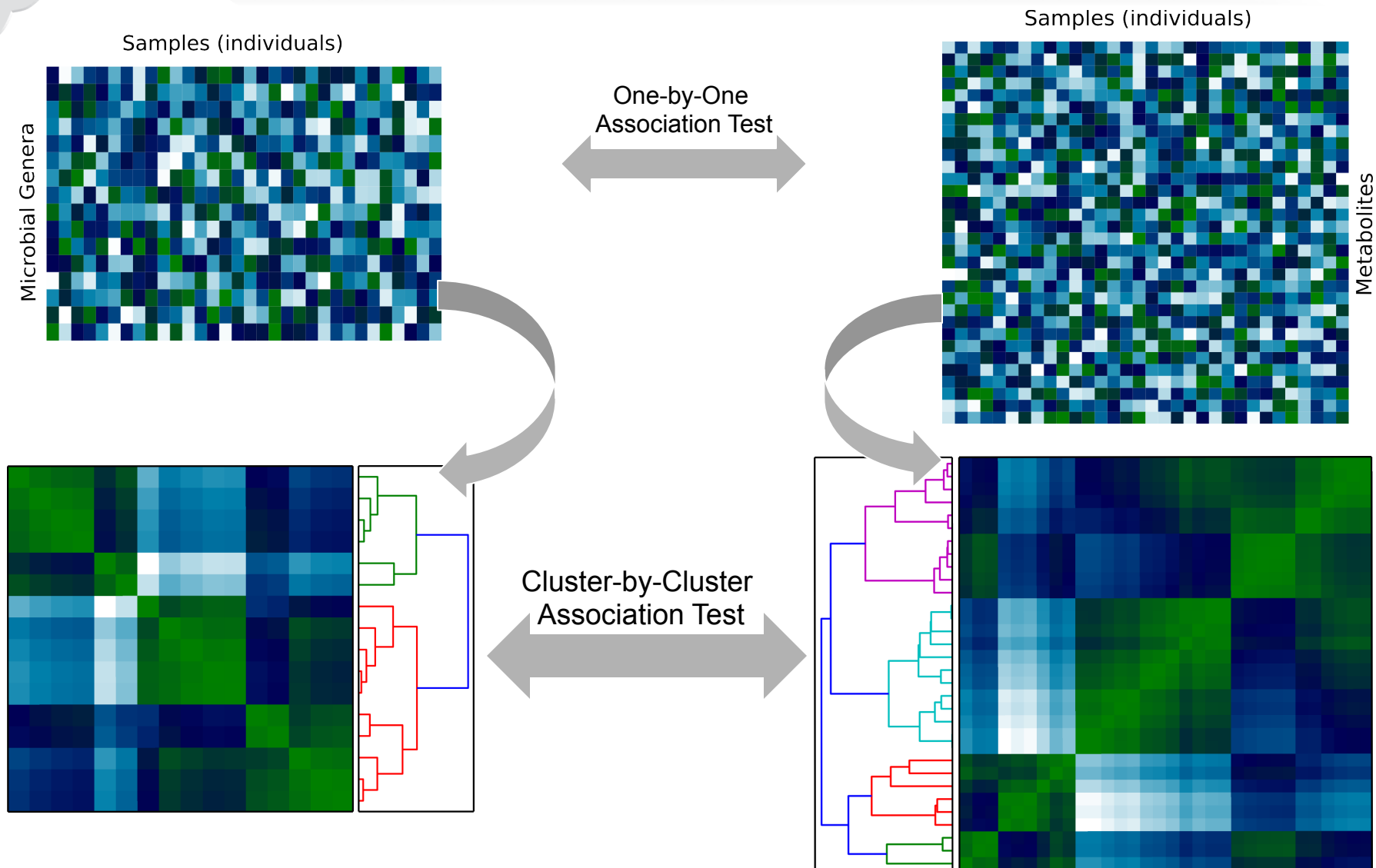
- Background
- Proposed Method
- Validation
- Evaluation
- Application
- Software

Fundamental Problem

How to discover quickly and confidently latent associations within high-dimensional, heterogeneous datasets in the presence of noise and collinearity?



Example: Multi-omics Integration



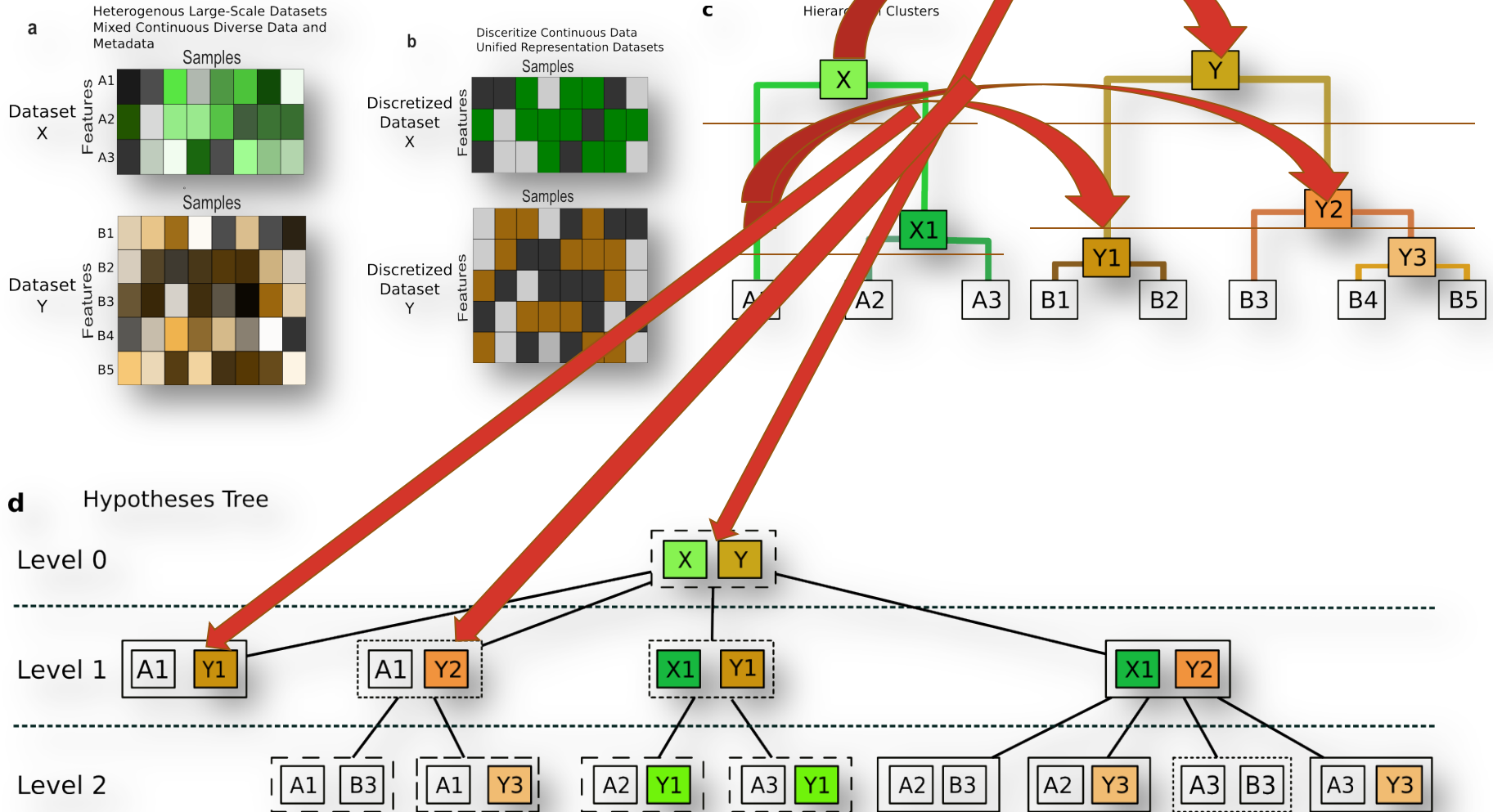


Proposed method

HIERARCHICAL ALL-AGAINST-ALL ASSOCIATION TESTING



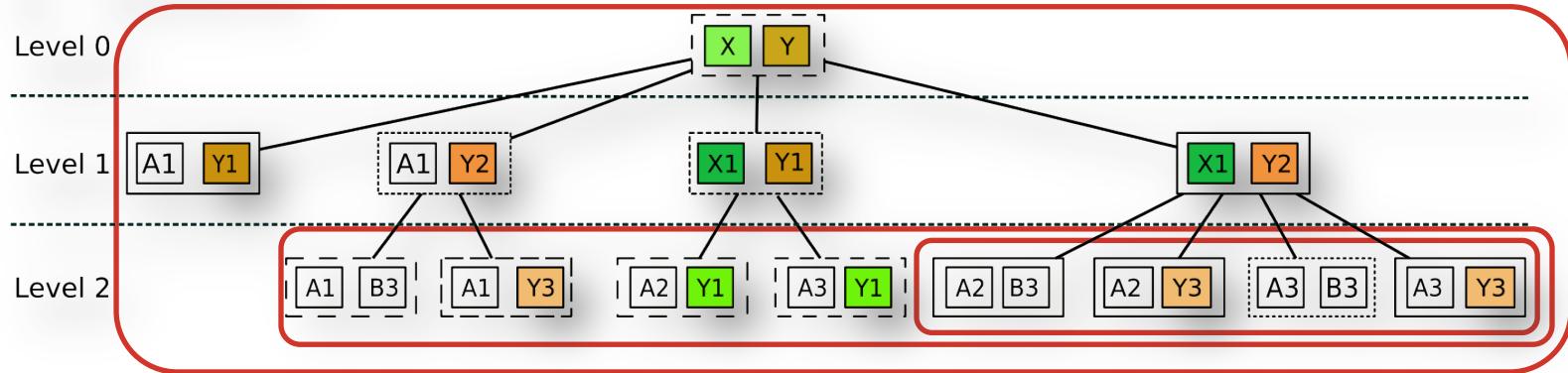
Discretizing, Clustering, and Coupling



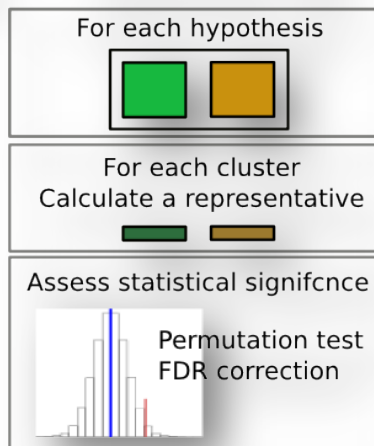


Hypothesis Testing, FDR Correction

d Hypotheses Tree

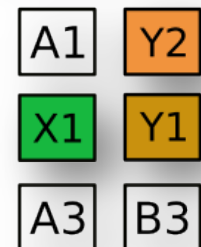


e Hypothesis Test



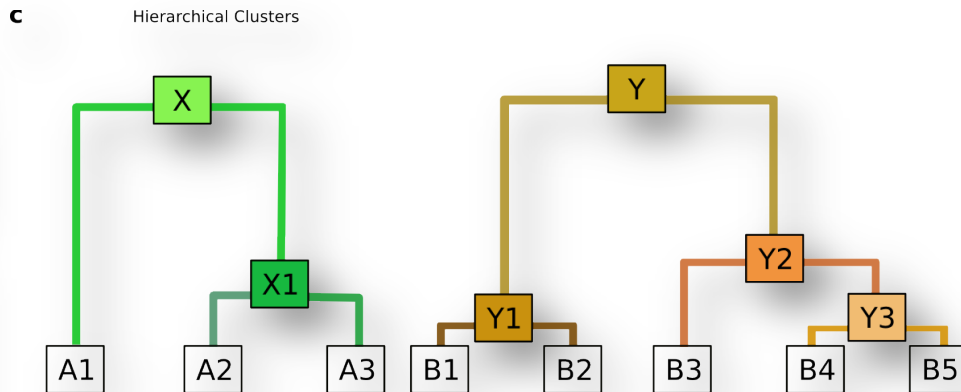
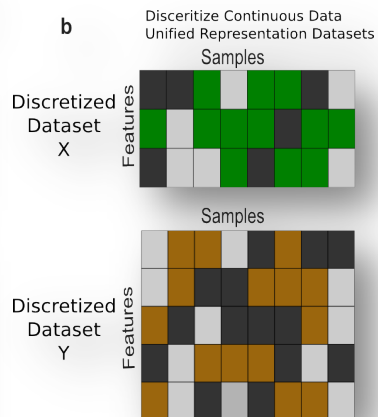
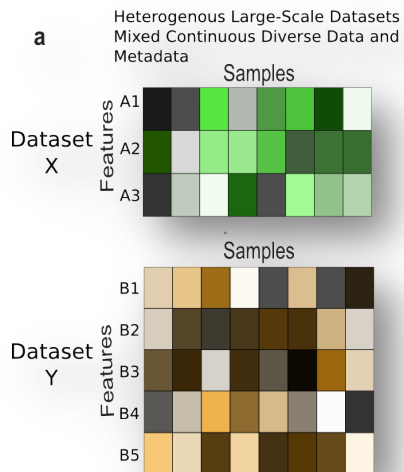
- Benjamini-Hochberg for FDR correction
 - Family of hypotheses
 - Level of hypotheses
 - All tested hypotheses

f Associations

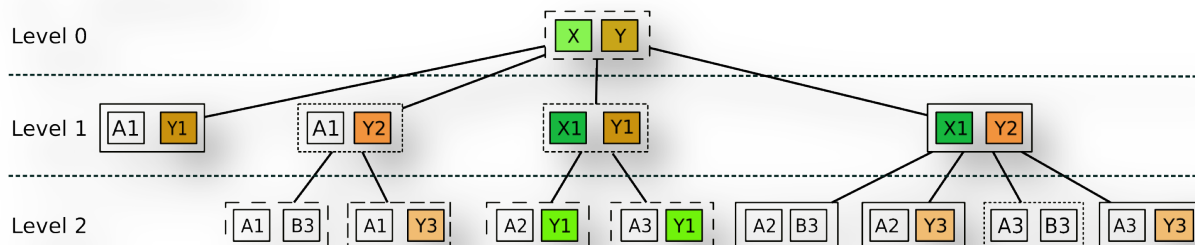




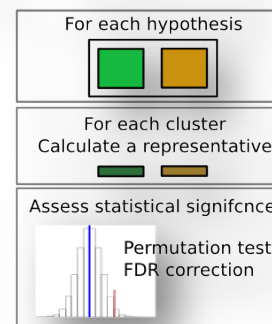
Method Overview



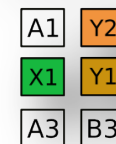
d Hypotheses Tree



e Hypothesis Test



f Associations





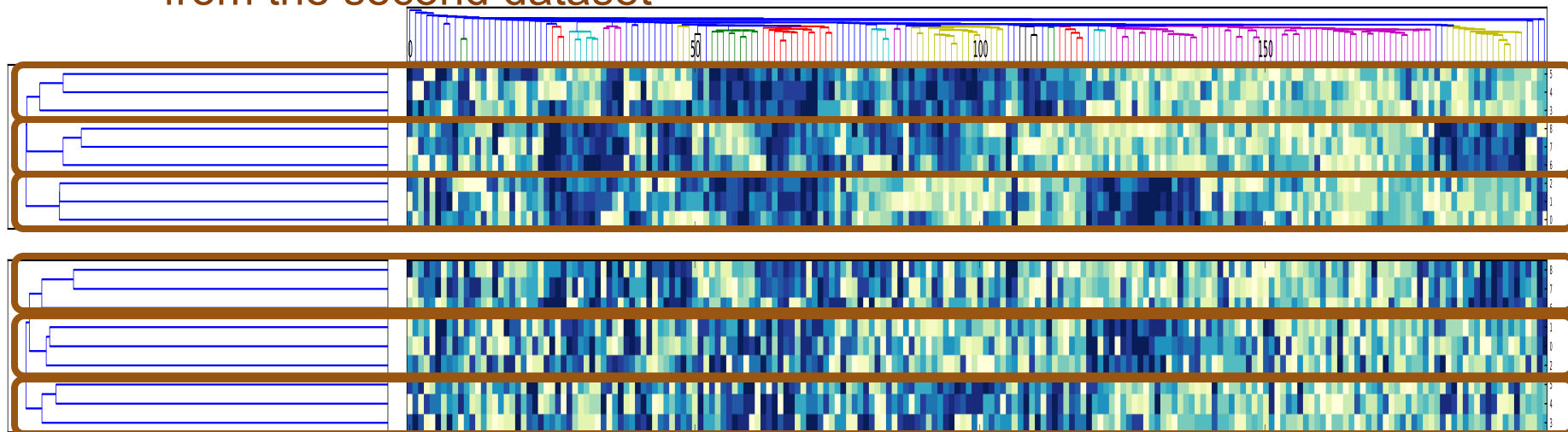
Validation

SYNTHETIC DATA



Simulated datasets

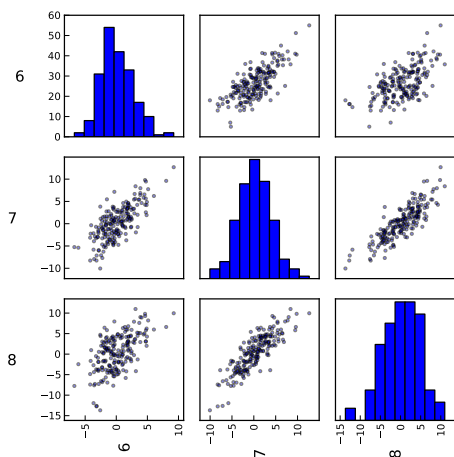
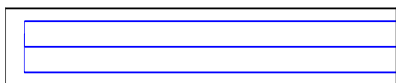
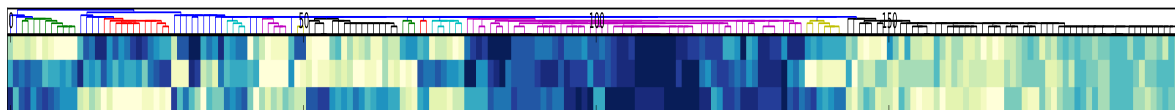
- Two synthetic datasets:
 - 12 features and 200 samples
 - 3 block of features with each datasets
 - Features within blocks are linearly correlated
 - The second dataset is spiked linearly from the first dataset
 - Each block from the first dataset is linearly correlated with one block from the second dataset



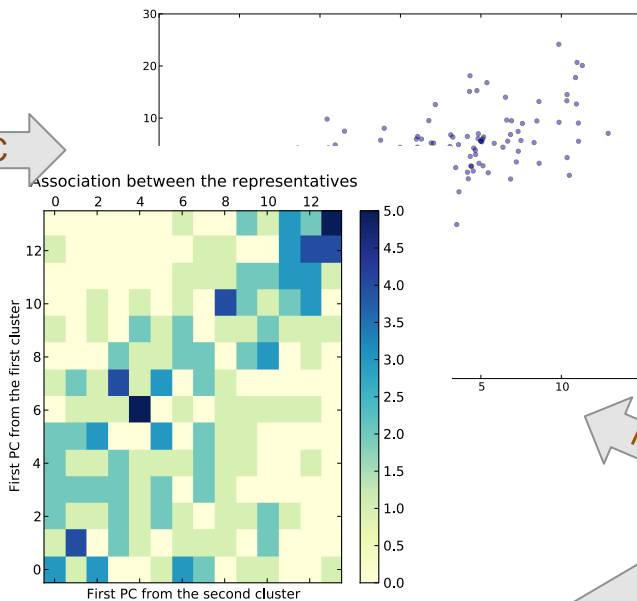


Associations

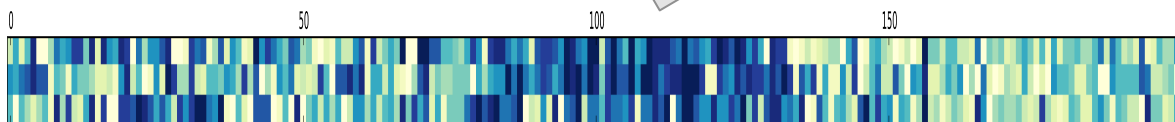
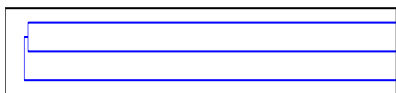
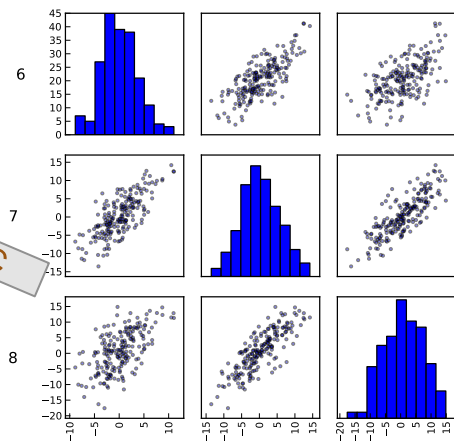
A cluster from the first dataset



First PC



First PC



A cluster from the second dataset



Evaluation

PERFORMANCE

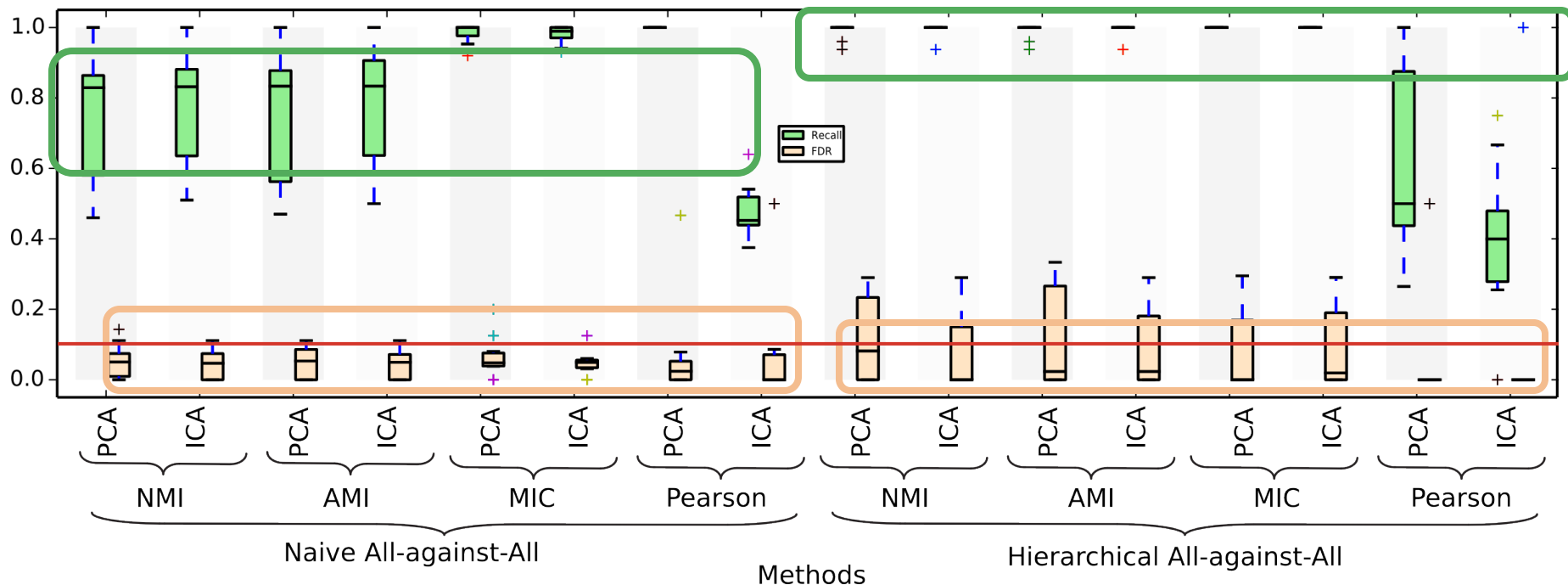


Simulated Datasets

- 10 different simulated datasets are used:
 - Datasets scaled from 4 features and 100 samples to 24 features and 500 samples.
 - These datasets have simulated clusters in their features with linear association within and between them.



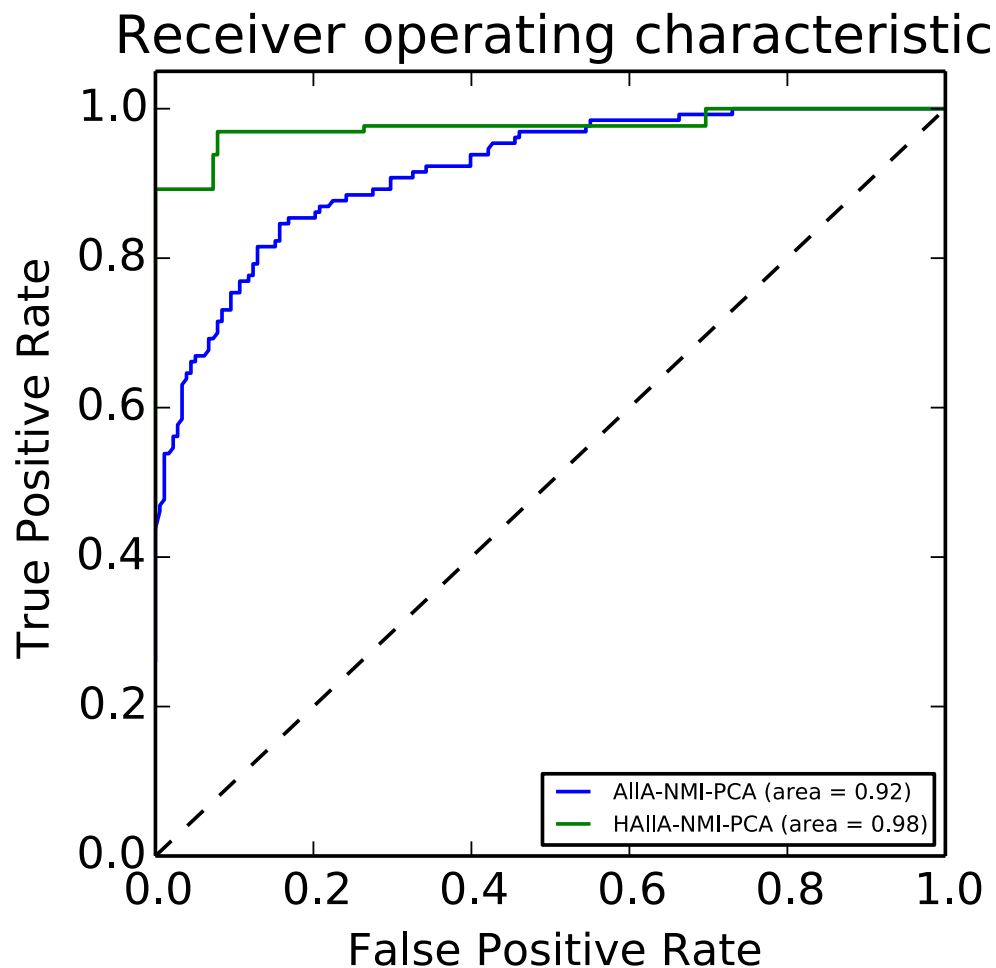
Naïve AIIA vs. HAIIA



- Higher statistical power in HAIIA
- Benjamini-Hochberg has controlled the FDR in the naïve All-against-All approach



HAIIA vs. AIIA: ROC Curve



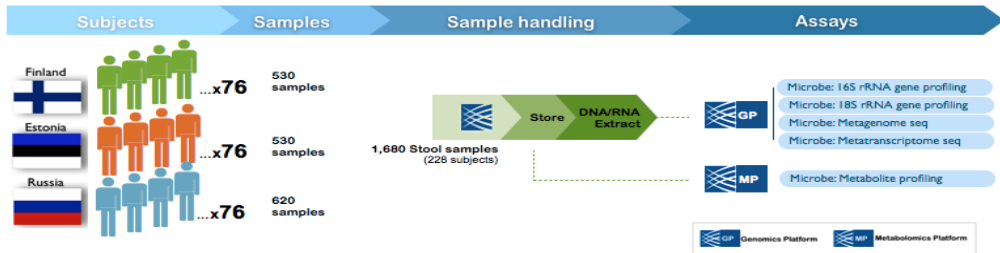
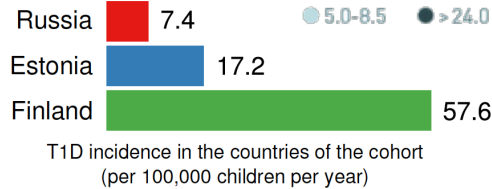
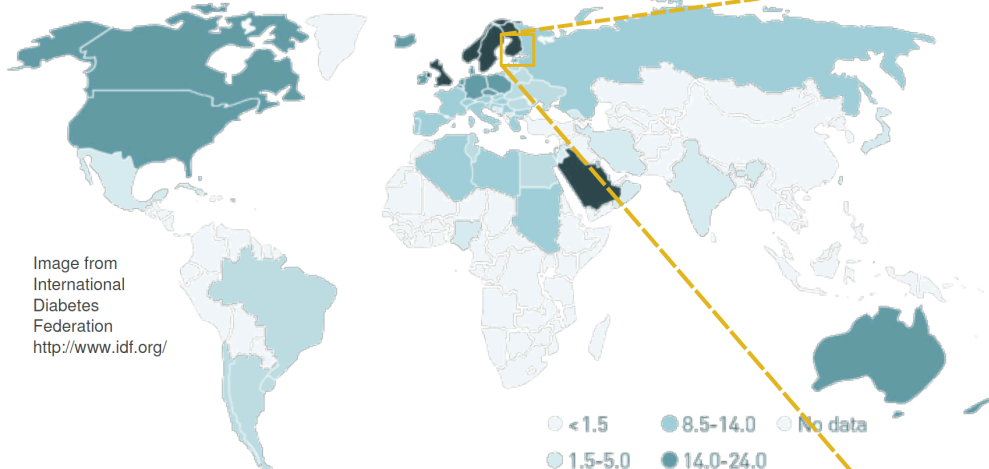


Application

GENERA VS. METABOLITES

The Multi-Country DIABIMMUNE Study

New cases of type 1 diabetes (0-14 years per 100,000 children per year), IDF 2011



Two datasets for 104 samples :

- 16S rRNA gene sequencing to profile 20 genera
- 284 labeled metabolites



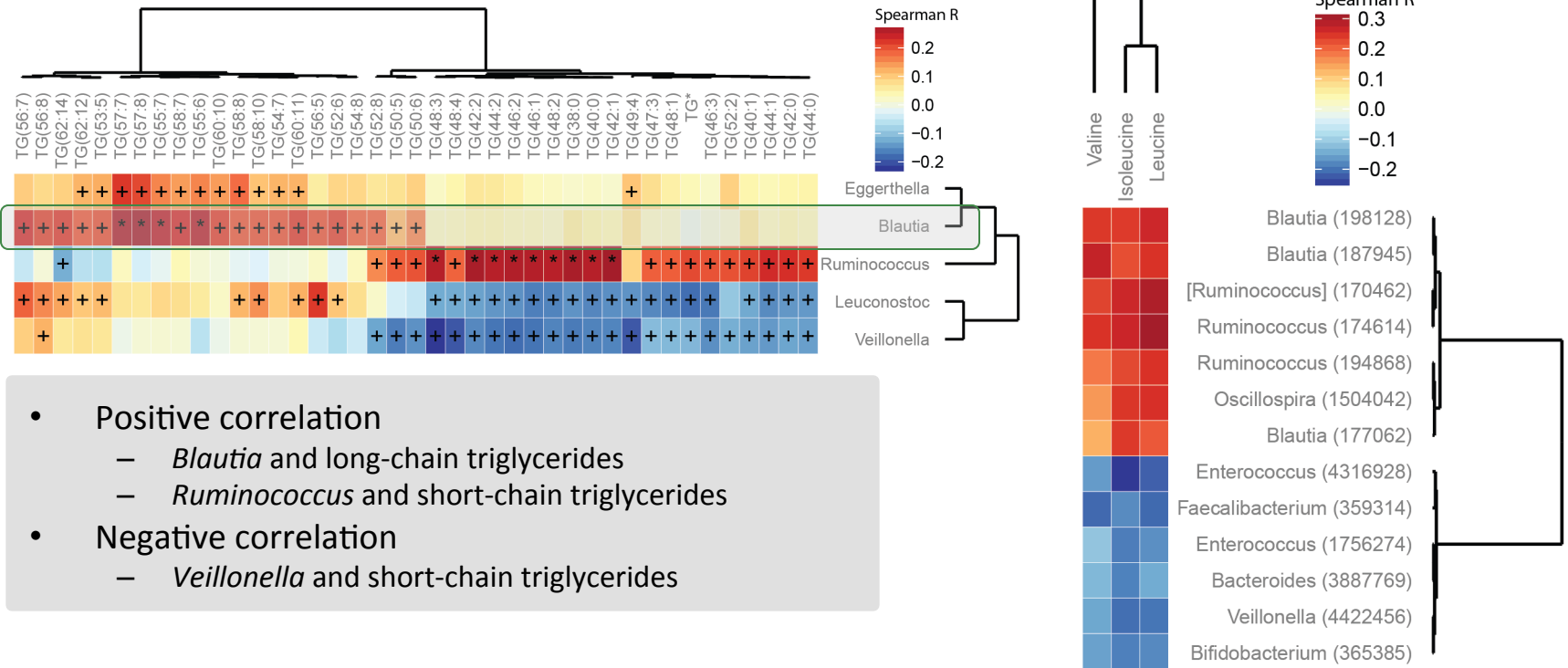
Aleksandar Kostic



Tommi Vatanen



Serum lipids and metabolites are associated with T1D-related microbial taxa

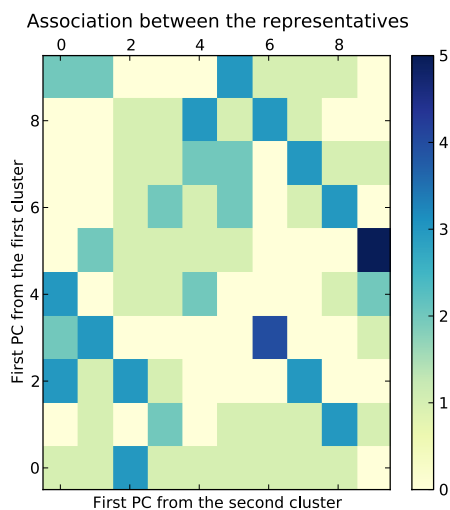
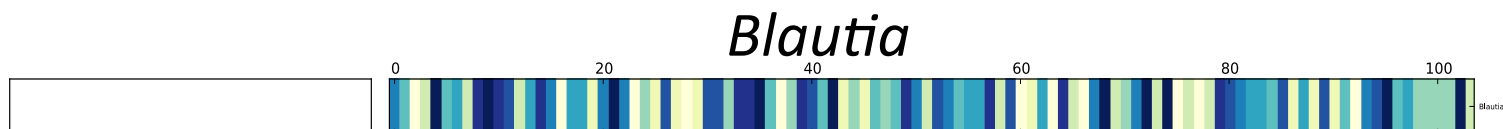


- Positive correlation
 - *Blautia* and long-chain triglycerides
 - *Ruminococcus* and short-chain triglycerides
- Negative correlation
 - *Veillonella* and short-chain triglycerides

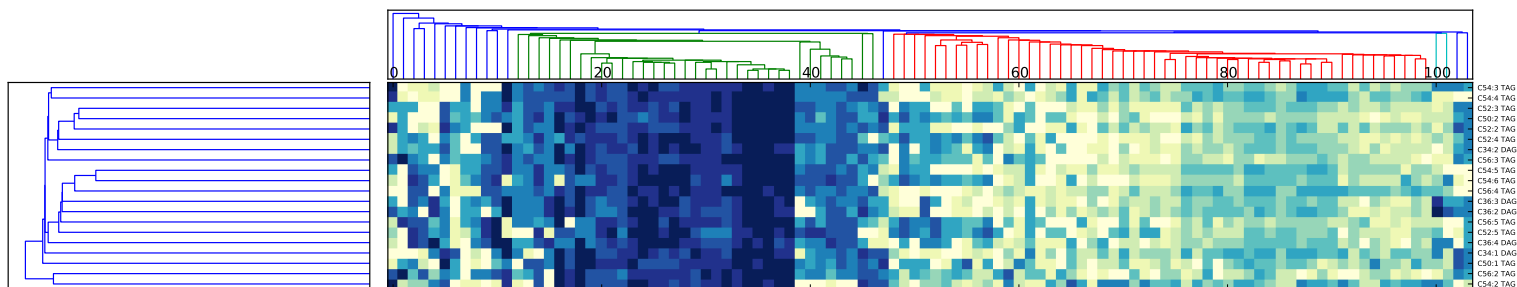
- Positive correlation
 - *Blautia* and branched-chain amino acids
 - *Ruminococcus* and branched-chain amino acids
- Negative correlation
 - *Veillonella* and branched-chain amino acids



Blautia and long-chain triglycerides



long-chain triglycerides





Software

FEATURES



Extensibility

Testing Approaches

- Naïve AIA
- Hierarchical AIA

Similarity Metrics

NMI
MIC
Pearson
AMI

Decomposition Methods

PCA
ICA
KPCA
PLS
CCA



Packaging

- Python package
 - Software
 - Unit tests
 - Evaluation Modules
- Online References
 - Hierarchical All-against-All(HAIIA)
 - <http://huttenhower.sph.harvard.edu/halla>
 - The Huttenhower Lab
 - <http://huttenhower.sph.harvard.edu>



Installation and Running

- Download
- `python setup.py install`
- Command line
 - `halla -X dataset1.txt -Y dataset2.txt -o output --plotting-results`



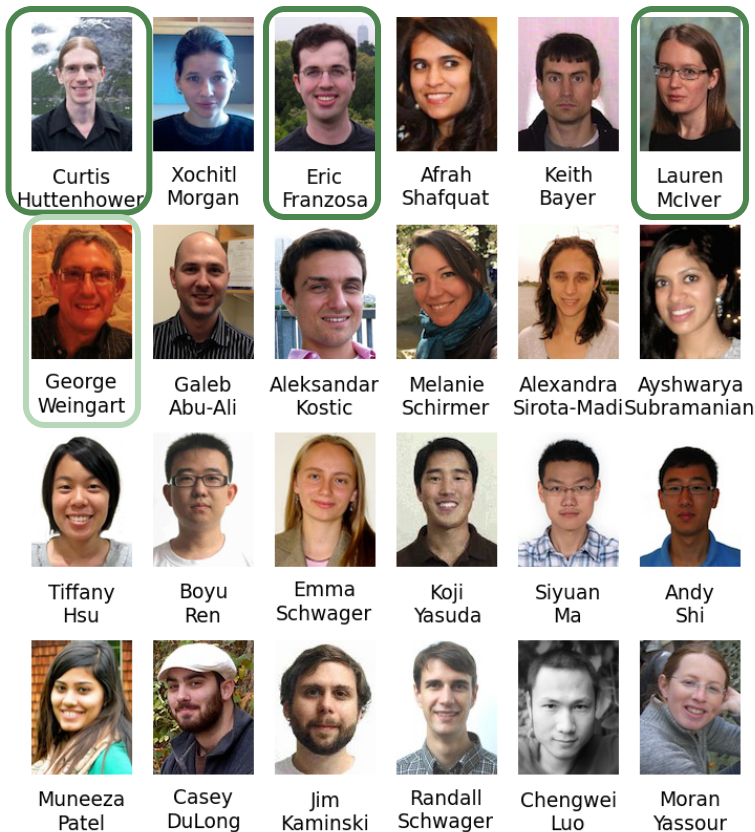
Summary Advantages

- Works with heterogeneous data (continuous+ categorical)
- Able to detect individually weak, but collectively strong signals
- Interpretable approach to data dimension reduction
- Assigns q-values to associations with high sensitivity
- Extensible to the use of different clustering and association methods



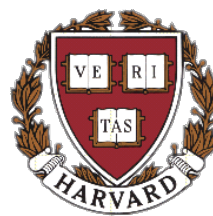
Acknowledgments

The Huttenhower Lab



Alumni

- Levi Waldron
- Yo Sup Moon

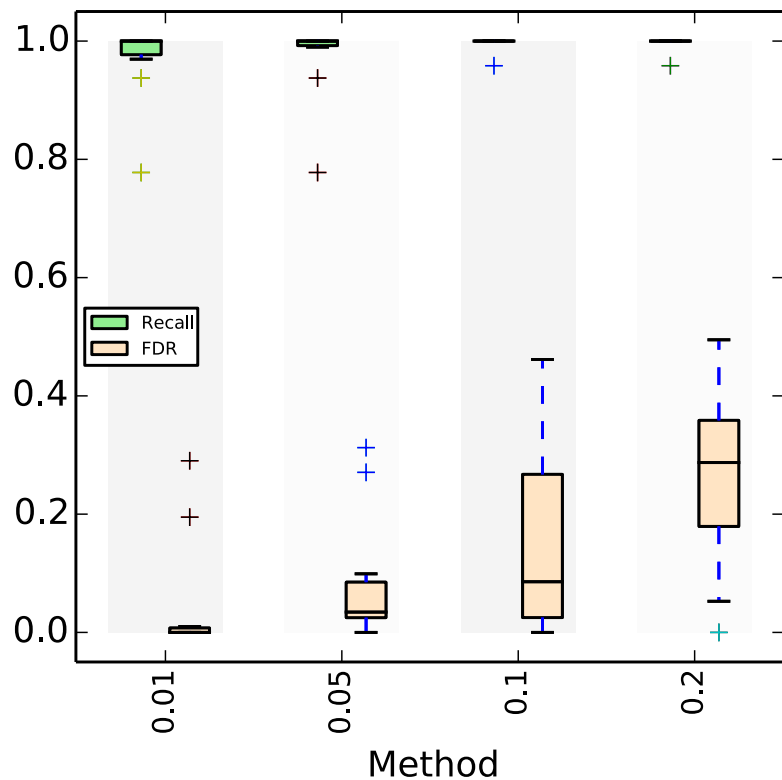


Harvard T.H. Chan School of Public Health
Department of Biostatistics



FDR Controlling

Benjamini-Hochberg-Family FDR correction



Benjamini-Hochberg-Level FDR correction

