

Caring Without Sharing: GWAS in a Decentralized Setting

Armin Pourshafeie

Bustamante Lab

Joint work with:

Carlos D. Bustamante, Snehit Prabhu

Running Decentralized GWAS

- Why and what else?
- Methods
- Simulations/experiments



Introduction

Meta-studies

Limitations

QC

Methods

PCA

Association

QC

Results

PCA

Association

Introduction

- Goal: discover variants associated with a particular phenotype
- Discovering variants with small effect sizes requires large datasets
 - Data sharing can help
- Centralizing data is difficult (Hardware, policy, etc.)

Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

Association

QC

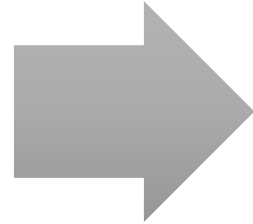
PCA

Association

Meta-studies

- Combine the results of previous studies on the same phenotype

Find compatible studies to combine



Choose a model (fixed effect vs. random effect) and combine the estimates

Introduction

Meta-studies

Limitations

QC

Methods

PCA

Association

Results

QC

PCA

Association

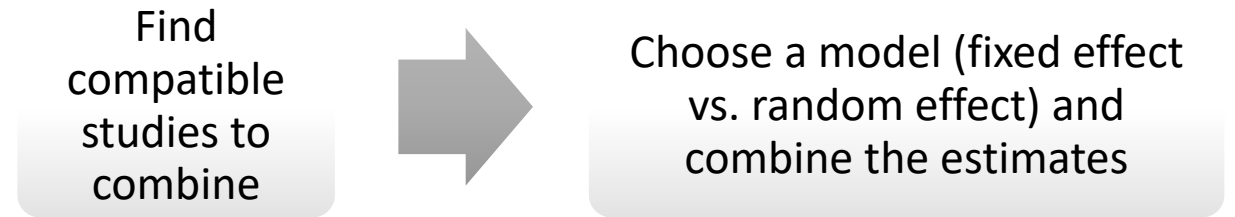
Meta-studies

- Pros:

- Familiar
- Readily available data
- Computationally efficiency
- Asymptotically statistical efficiency (Lin and Zeng 2010) (**asymptotic in each study**)
- Some level of privacy

- Cons:

- Unable to use small datasets
- Difficult/non-existent quality control
- Multiple regression is not possible
- Difficult to control for structure/duplicates



Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

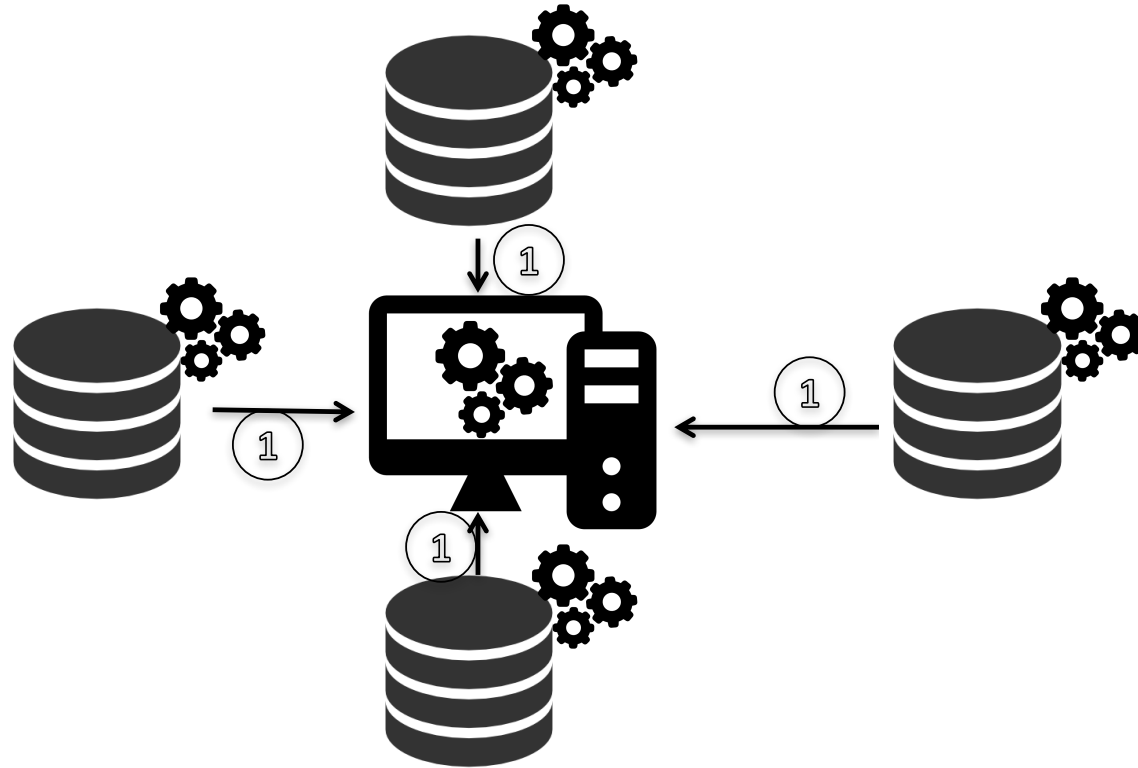
Association

QC

PCA

Association

Different paradigms



Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

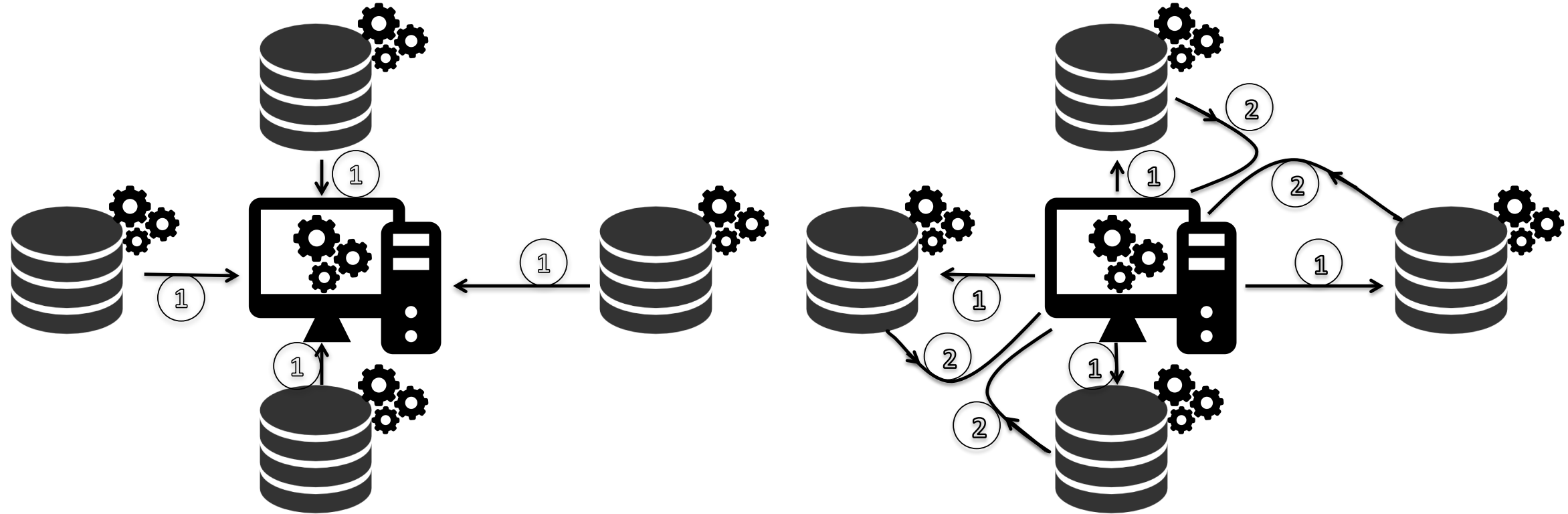
Association

QC

PCA

Association

Different paradigms



Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

Association

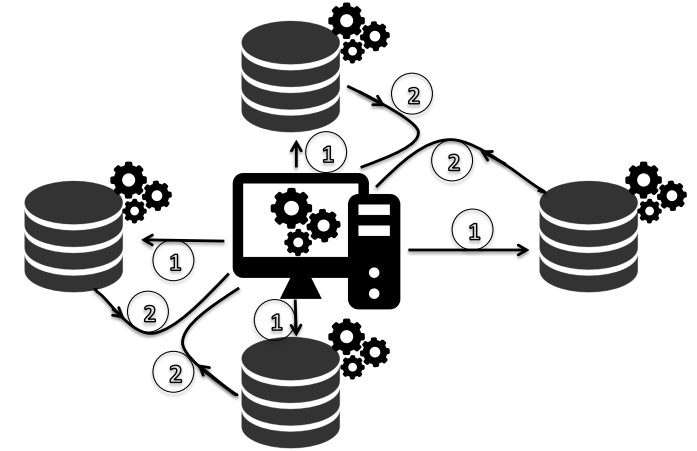
QC

PCA

Association

Decentralized GWAS

- Quality control
- Population structure control (PCA)
- Imputation
- Association (logistic regression)



Introduction

Meta-studies

Limitations

QC

Methods

PCA

Association

Results

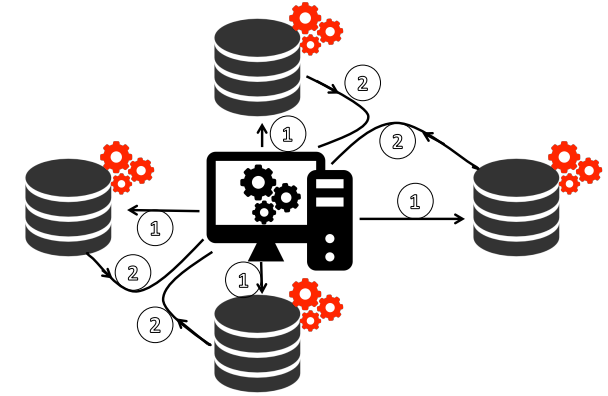
QC

PCA

Association

QC

- No communications (almost):
 - Calling quality, missing per individual
- Few communications
 - #(Missing, Homo-ref, Hetro, Homo-alt)
 - Missing-per loci, Allele Freq, Hardy-Weinberg
 - Relatedness:
 - Hashing (Dan He, et al. 2014)
- LD-pruning.
 - Hard :(
 - Pass in a matrix after thinning (very local pruning)



Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

Association

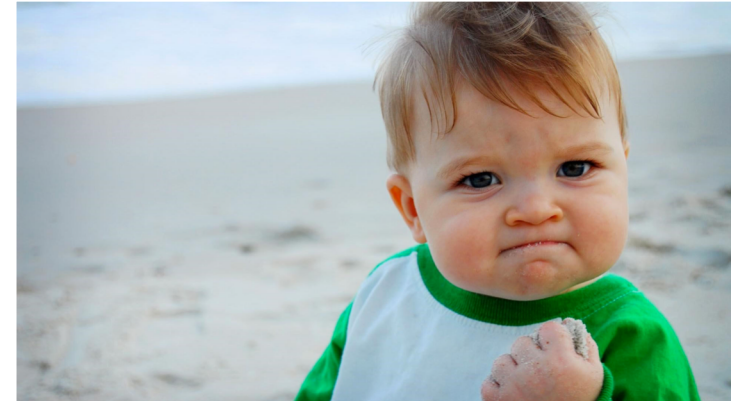
QC

PCA

Association

Population Structure Control (PCA)

- Easy:
 - Project everyone on dimensions discovered from a public dataset (1KG, Hapmap, etc.)
 - No need for LD pruning
 - Cheap, and fast



- Biased, not applicable to underrepresented populations

Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

Association

QC

PCA

Association


Population Structure Control (PCA)

- Hard:
 - ind-ind covariance matrix won't work

$$\bullet G_1 = \begin{bmatrix} ind_{1,1} \\ \vdots \\ ind_{1,N_1} \end{bmatrix}, G_2 = \begin{bmatrix} ind_{2,1} \\ \vdots \\ ind_{2,N_2} \end{bmatrix}, \dots, G_k = \begin{bmatrix} ind_{k,1} \\ \vdots \\ ind_{k,N_k} \end{bmatrix}$$

$$\bullet G^T G = \begin{bmatrix} block1 & & \\ & block2 & \\ & & block3 \end{bmatrix}$$

Missing



Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

Association

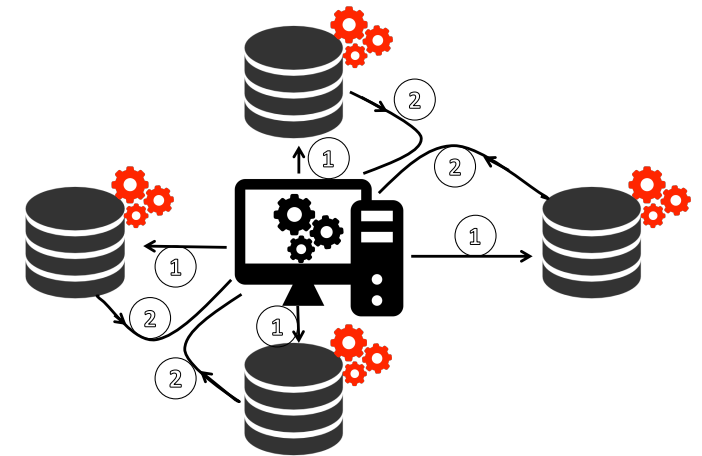
QC

PCA

Association

Population Structure Control (PCA)

- Use the LD-matrix (gene-gene) instead
 - Compute Gene-gene covariance matrix. All the genotypes of each individual is in a single dataset. The overall LD-matrix is simply the sum of these LD-matrices
- Pseudo-algorithm
 1. Compute the local LD-matrix
 2. Average the local LD-matrices at the center
 3. Perform eigen-decomposition
 4. Back solve for loadings at each silo



Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

Association

QC

PCA

Association

Population Structure Control (PCA)

- Pros:
 - This is impossible to do in meta studies
 - Can implement with differential privacy
- Cons:
 - The LD-matrix is very large
 - This method is inefficient with many small size silos

Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

Association

QC

PCA

Association

Association

- Notation:

- $\cdot^{(k)} := k^{th}$ silo

- $\ell^{(k)}(\beta) :=$ -log-likelihood function evaluated on silo k with parameter β

Centralized

Meta-study (FE)

$$\hat{\beta} = \operatorname{argmin}_{\theta} \sum_k \ell^{(k)}(\theta)$$

$$z^{(k)} = \operatorname{argmin} \ell^{(k)}(x)$$

$$\hat{\beta} = \sum w^{(k)} z^{(k)}$$

Assumption: $z^{(k)} = \beta + \varepsilon^{(k)}$

$$\hat{\beta} = \operatorname{argmin}_{\theta} \sum_k \min_x \ell^{(k)}(x^{(k)}) \text{ s.t. } x^{(k)} = \theta$$

$$L_{\rho}(\theta, \lambda, x) = \sum_k \ell^{(k)}(x^{(k)}) + \lambda^T (x^{(k)} - \theta) + \frac{\rho}{2} \|x^{(k)} - \theta\|_2^2$$

Lagrange Multiplier

Augmented Lagrangian
(Hestenes, Powell 1969)

Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

Association

QC

PCA

Association

Association

- $L_\rho(\theta, \lambda, x) = \sum_k \ell^{(k)}(x^{(k)}) + \lambda^T(x^{(k)} - \theta) + \frac{\rho}{2} \|x^{(k)} - \theta\|_2^2$
- Updates:
 - $z^{(k)} \leftarrow \operatorname{argmin}_x \ell^{(k)}(x^{(k)}) + \frac{\rho}{2} \|x^{(k)} - \theta + \lambda^{(k)}\|_2^2$ At each silo
 - $\theta \leftarrow \frac{1}{K} \sum_k z^{(k)}$ At the center
 - $\lambda \leftarrow \lambda^{(k)} + z^{(k)} - \theta$

See Boyd, Stephen, et al. "Distributed optimization and statistical learning via the alternating direction method of multipliers." *Foundations and Trends® in Machine Learning* 3.1 (2011): 1-122.

Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

Association

QC

PCA

Association

Results



- Simulated GWAS on POPRES¹
 - 2274 ind ~ 400k Loci
 - Simulated a case-control phenotype according to a logistic model
 - 50-50
 - 10 causal SNPS with effect size drawn from a gaussian + noise
- Two experiments:
 - 5 silos, random distribution ($n \approx 450$ per silo)
 - 2 silos, cases vs controls
- All regressions include 1 SNP + 5 PCs

1. Nelson, Matthew R., et al. *AJHG* (2008):

Introduction

Methods

Results

Meta-studies

Limitations

QC

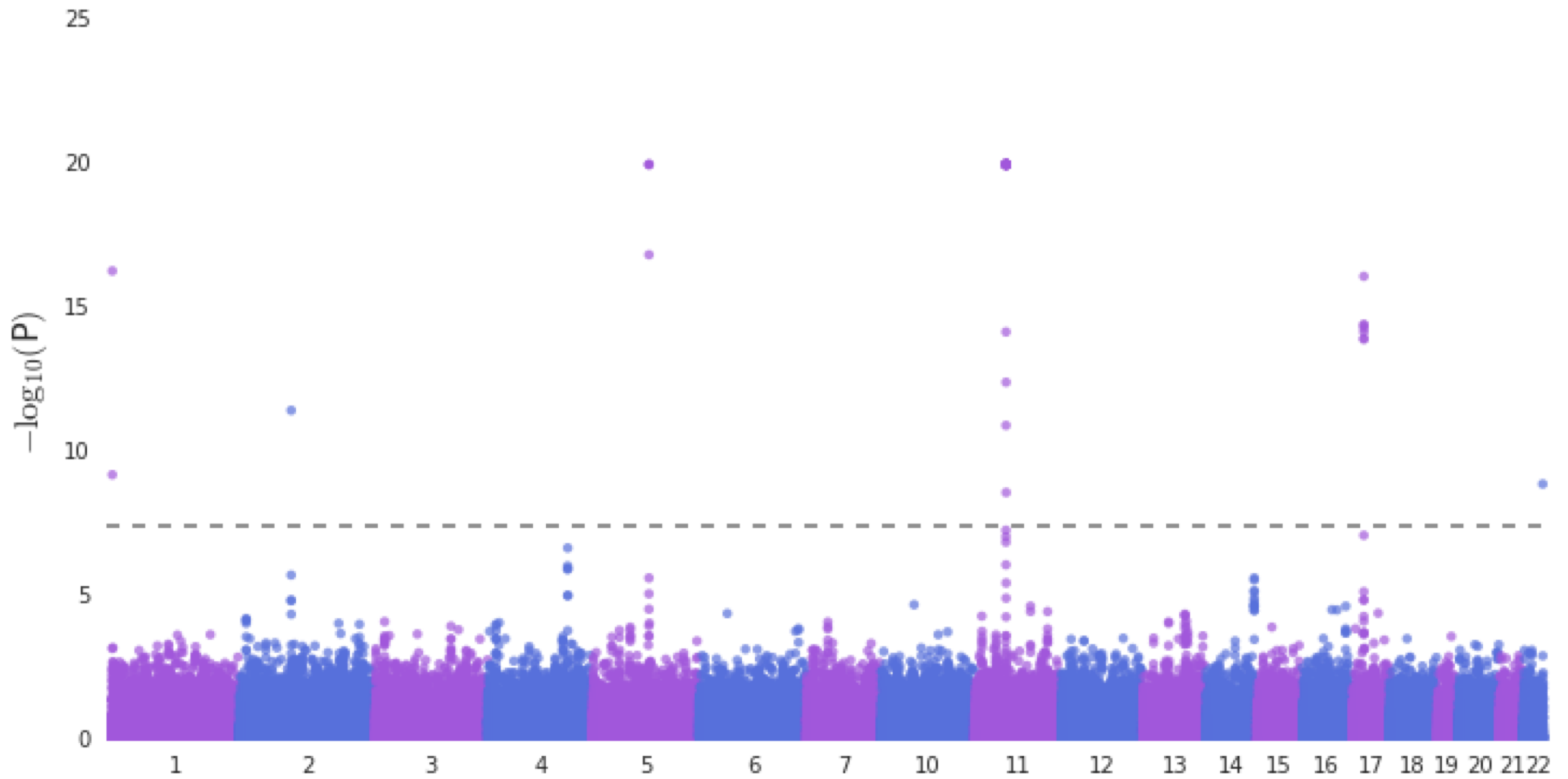
PCA

Association

QC

PCA

Association



Introduction

Methods

Results

Meta-studies

Limitations

QC

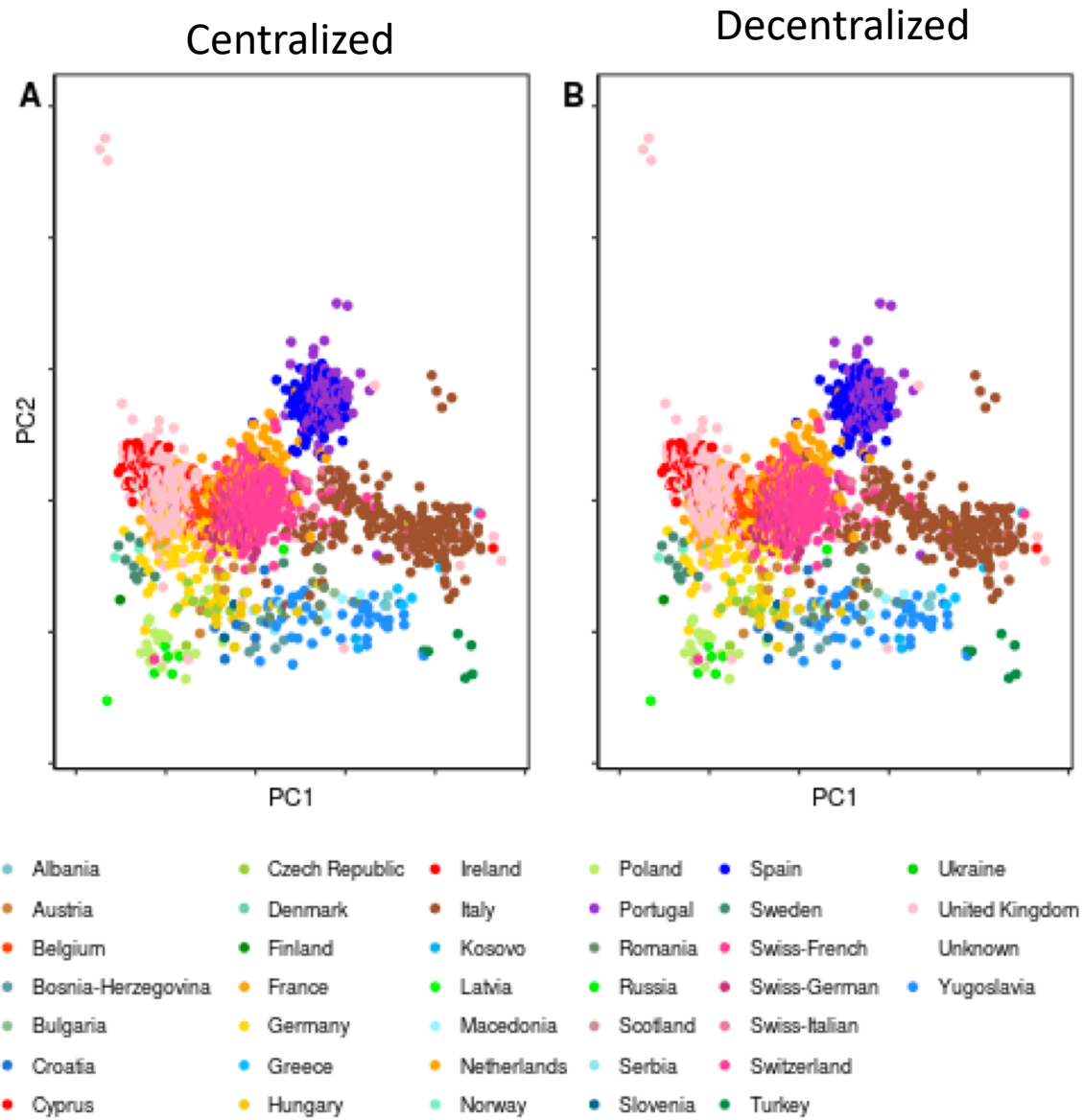
PCA

Association

QC

PCA

Association



Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

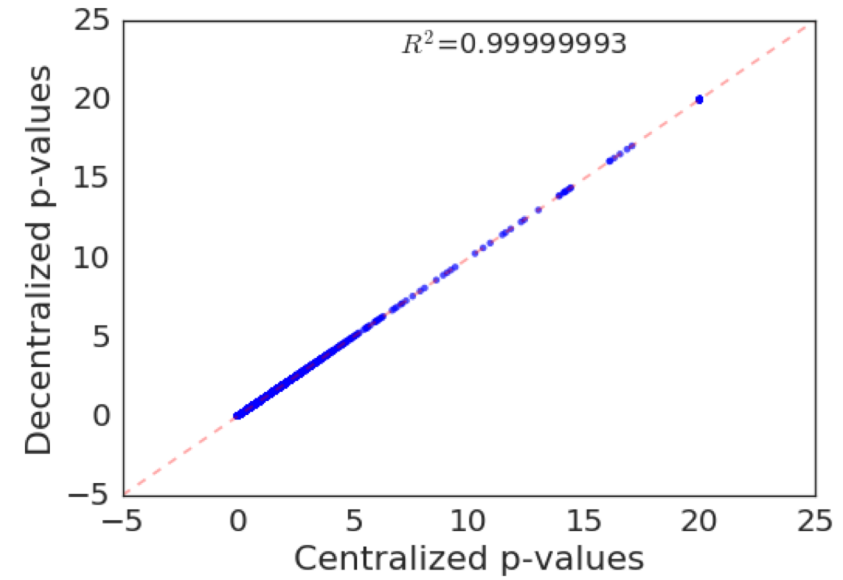
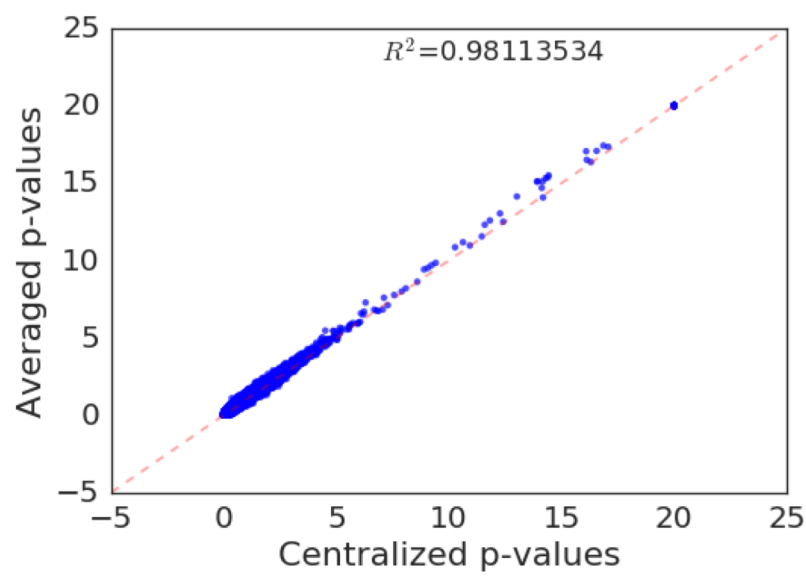
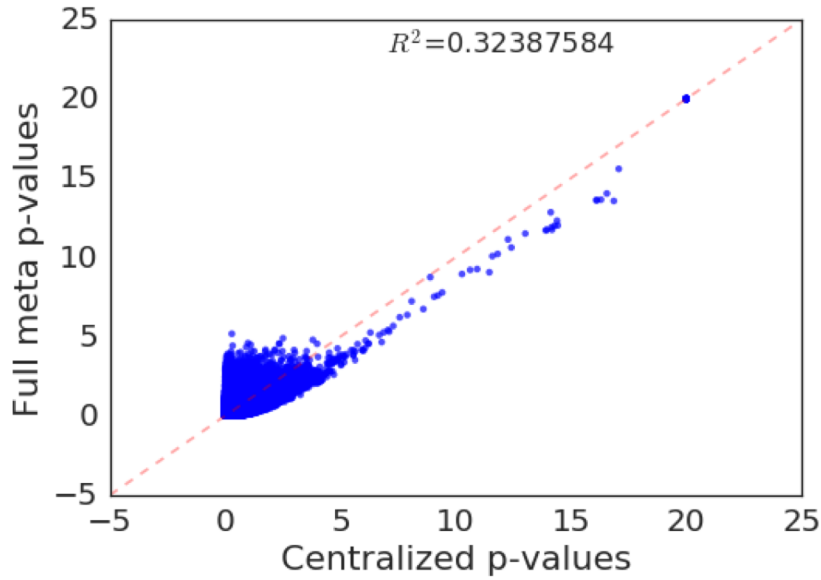
Association

QC

PCA

Association

Experiment 1: (iid distributed individuals, 5 Silos)



Introduction

Meta-studies

Limitations

Methods

QC

PCA

Association

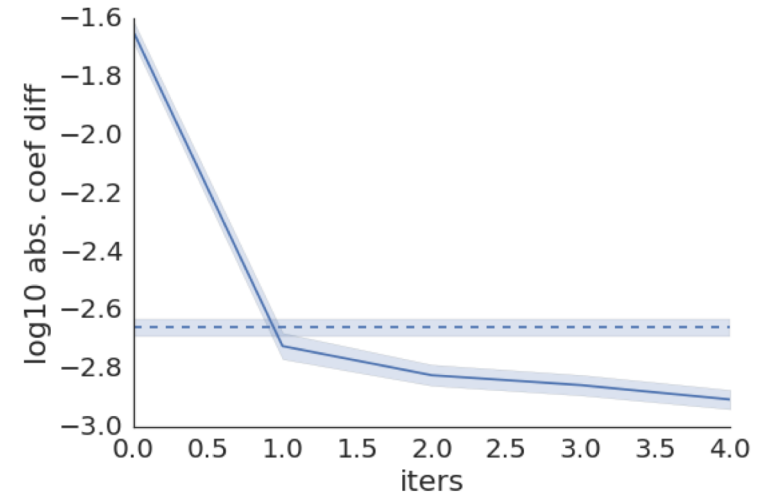
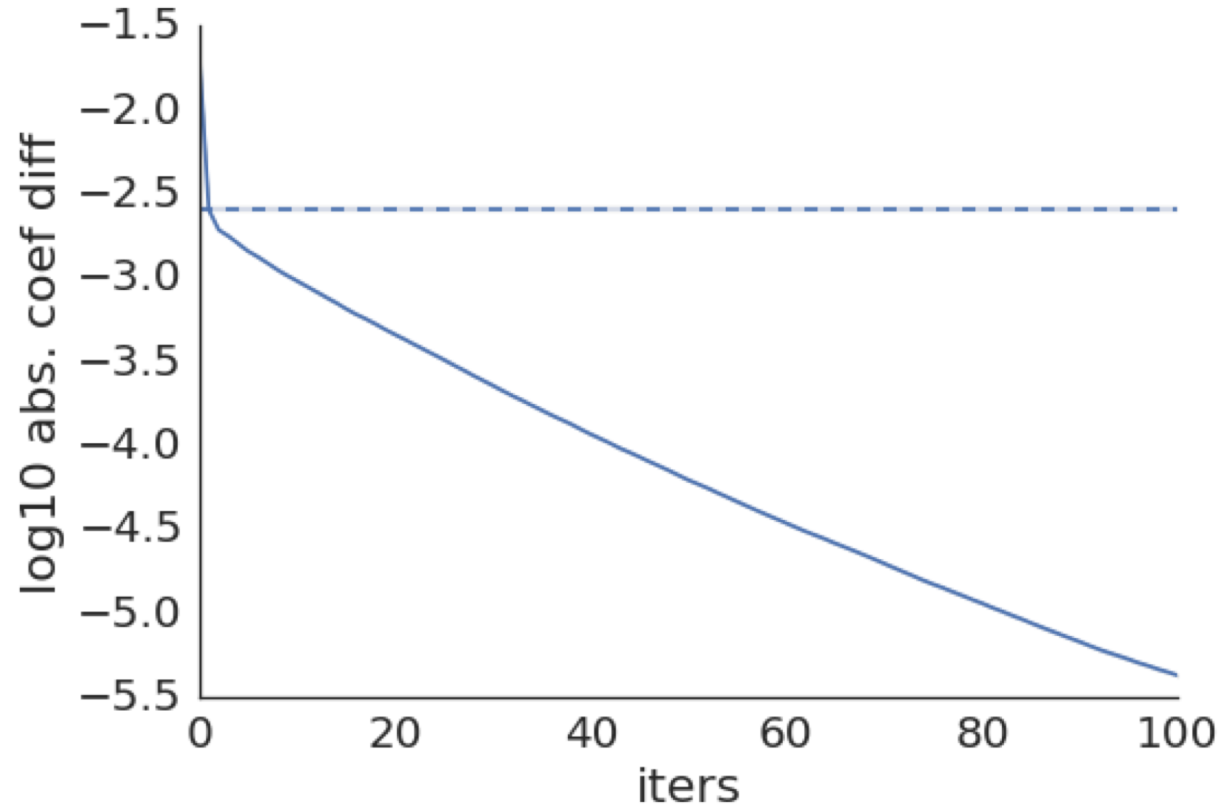
Results

QC

PCA

Association

Experiment 1: (iid distributed individuals, 5 Silos)



Introduction

Methods

Results

Meta-studies

Limitations

QC

PCA

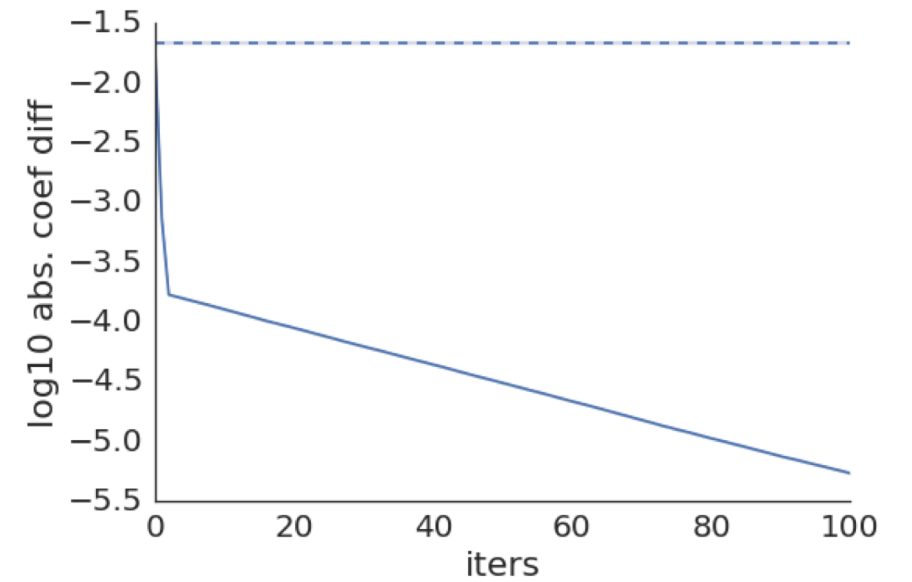
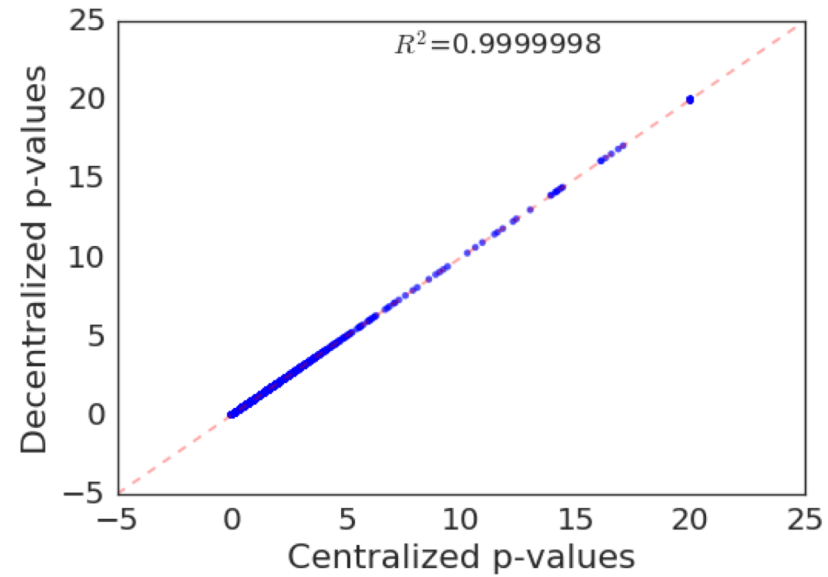
Association

QC

PCA

Association

Experiment 2: (Cases vs. Controls)



Introduction

Meta-studies

Limitations

QC

Methods

PCA

Association

QC

Results

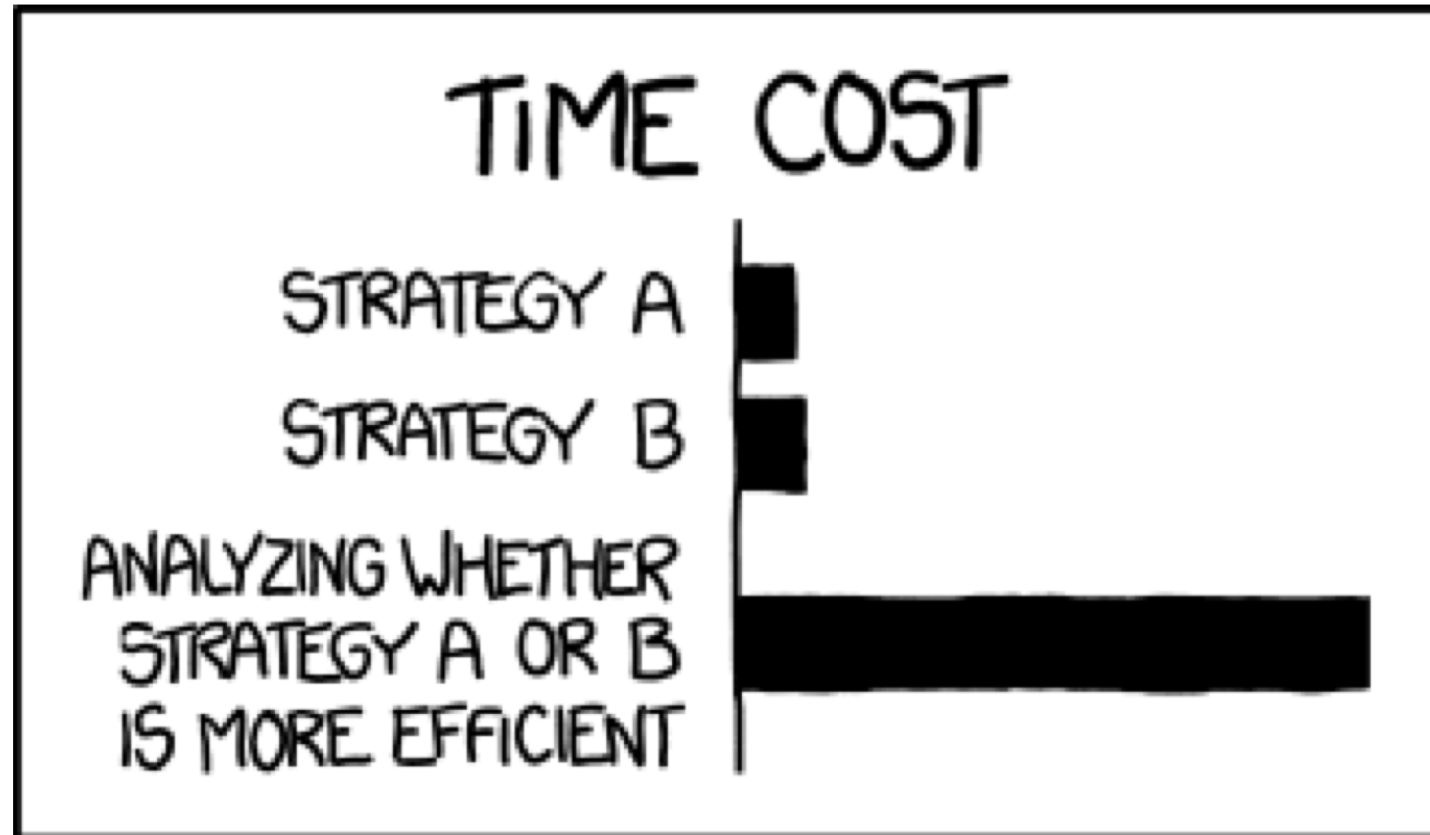
PCA

Association

Acknowledgments

- Collaborators
 - Carlos Bustamante
 - Snehit Prabhu
- Funding:
 - NHGRI SGTP
- **Thank you!**

Questions?



THE REASON I AM SO INEFFICIENT