

Matrix ASH: Modeling Genetic Effects Across Multiple Subgroups

Sarah Urbut, Gao Wang, Matthew Stephens

Department of Human Genetics and University of Chicago

Introduction

- Variation in gene expression is an important mechanism underlying susceptibility to complex disease.
- However, most studies to date have been conducted in a single immortalized peripheral cell type or single tissue framework
- The solution: GTEx! by 2016: 900 post-mortem donors, with approximately 30 tissues collected from each donor
- Our mission: *jointly analyze data on all tissues* to maximize power, and to identify and quantify the variability in effect sizes.

Objectives

- Combine information across tissues
- Capture distinct variation in effect sizes within and between subgroups: 'patterns of sharing'

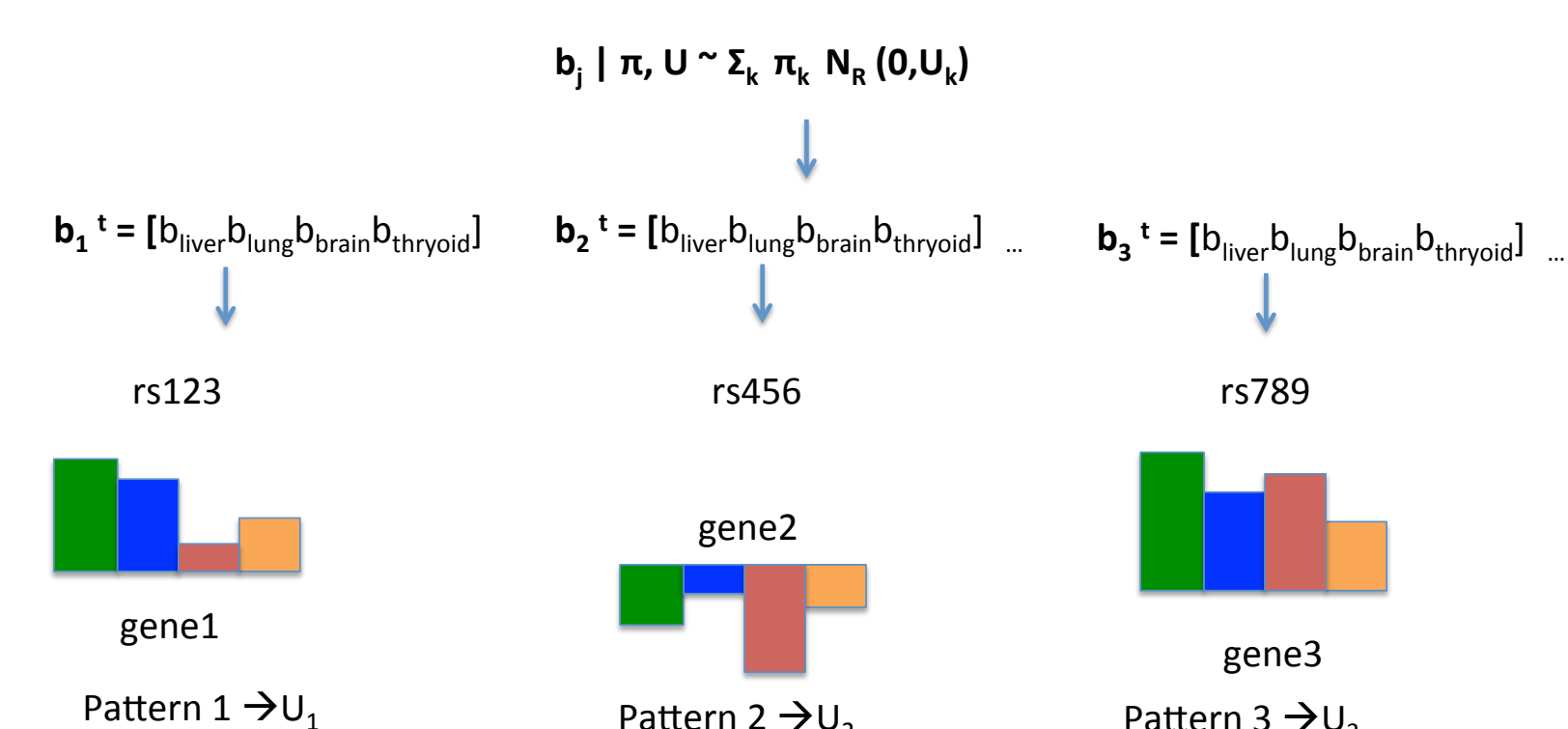


Figure 1: Figure caption

- Assume that each eQTL belongs to a group characterized by its effects across tissues.
- Within these groups, the tissues exhibit characteristic patterns of sharing
- Natural mixture model, in which we assume all the gene-snp pairs arise from a mixture of a finite number of Gaussian distributions
- Each component of the mixture is defined by the prior covariance matrix U_k from which the vector of standardized effect sizes \mathbf{b}_j is drawn.
- 'Learn' the relative proportions of each pattern of sharing from the data
- Key: Distinct data-sensitive diagonal and off-diagonal elements capture a wide array of patterns of sharing**

Mathematical Section

We assume the following mixture prior for the R dimensional vector of true effects, \mathbf{b}_j represents the genetic effect of SNP-gene pair j across $R = 44$ tissues:

$$\mathbf{b}_j | \pi, \mathbf{U} \sim \sum_{k,l} \pi_{k,l} N_{\mathbf{R}}(\mathbf{0}, \omega_l \mathbf{U}_k) \quad (1)$$

For a given ω_l , we specify 4 'types' of $R \times R$ prior covariance matrices $\mathbf{U}_{k,l}$.

- $\mathbf{U}_{k=1,l} = \omega_l \mathbf{I}_R$
- $\mathbf{U}_{k=2,l} = \omega_l \frac{1}{J} \mathbf{Z}_{center}^t \mathbf{Z}_{center}$
- $\mathbf{U}_{k=3,l} = \omega_l \frac{1}{J} \mathbf{V}_{1..p} d_{1..p}^2 \mathbf{V}_{1..p}^t$ is the rank p eigenvector approximation of the tissue covariance matrices
- $\mathbf{U}_{k=4+Q-1,l} = \frac{1}{J} ((\Lambda \mathbf{F})^t \Lambda \mathbf{F})_q$ corresponding to the q_{th} sparse factor representation of the tissue covariance matrix
- $\mathbf{U}_{k=4+Q,l} = \frac{1}{J} (\Lambda \mathbf{F})^t \Lambda \mathbf{F}$ is the sparse factor representation of the tissue covariance matrix, estimated using all q factors.
- $\mathbf{U}_{k=5+Q:R+4+Q,l} = \frac{1}{J} ([100..]^t [100...])$
- $\mathbf{U}_{k=R+5+Q,l} = \frac{1}{J} ([111..]^t [111...])$

- Compute the mixture weights π_{kl} hierarchically
- For a given gene-snp pair, the Likelihood on \mathbf{b} :

$$\hat{\mathbf{b}} | \mathbf{b} \sim N_{\mathbf{R}}(\mathbf{b}, \hat{\mathbf{V}}) \quad (2)$$

- We know that for a single multivariate *Normal* the posterior on $\mathbf{b} | \mathbf{U}_0$ is simply:

$$\mathbf{b} | \hat{\mathbf{b}} \sim N_{\mathbf{R}}(\mu_1, \mathbf{U}_1)$$

$$p(\mathbf{b} | \hat{\mathbf{b}}, \hat{\mathbf{V}}, \hat{\pi}) = \sum_{k=1, l=1}^{K,L} \sim N_{\mathbf{R}}(\mu_{1kl}, \mathbf{U}_{1kl}) \tilde{\pi}_{k,l} \quad (3)$$

- Since the prior weights are computed from 'mostly null data' the prior weights will heavily weight the components with small posterior means (as determined by small prior variance in \mathbf{U}_k).
- $\tilde{\pi}_{k,l} = P(\text{Component} | \text{Data}) \propto P(\text{Data} | \text{Comp.}) \times P(\text{Comp.})$: Combine hierarchical and snp-specific information
- Allows pair to find its true match!

Results

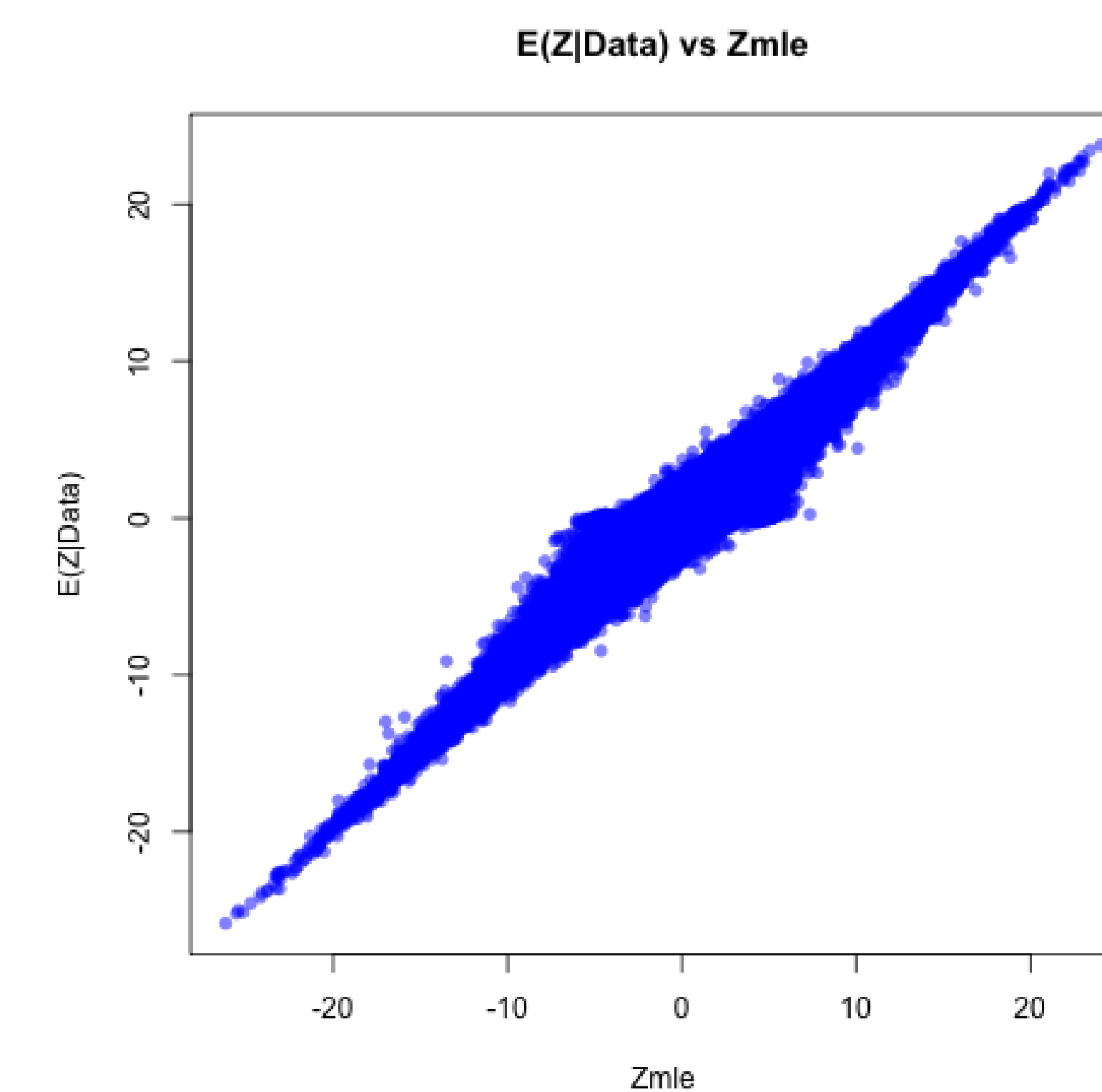


Figure 2: E(Z|Data) vs MLE: Using information gleaned from our high 'prior' probability of observing a particular effect due to the effects other tissues leads us to shrink some summary statistics more than others

Identify strong tissue specific QTLs as gene-snp pairs in which *LFSR* is low in only one tissue.

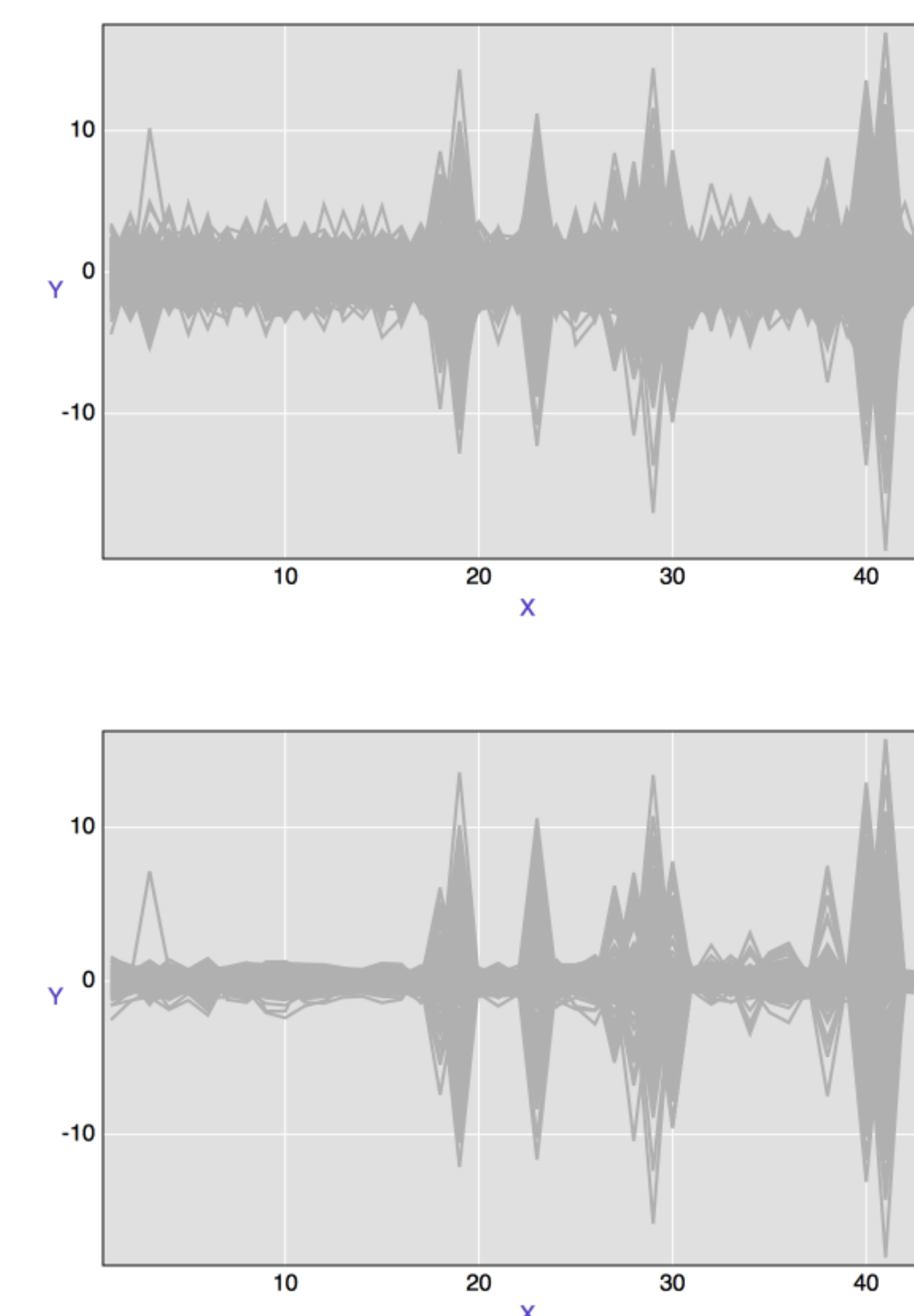
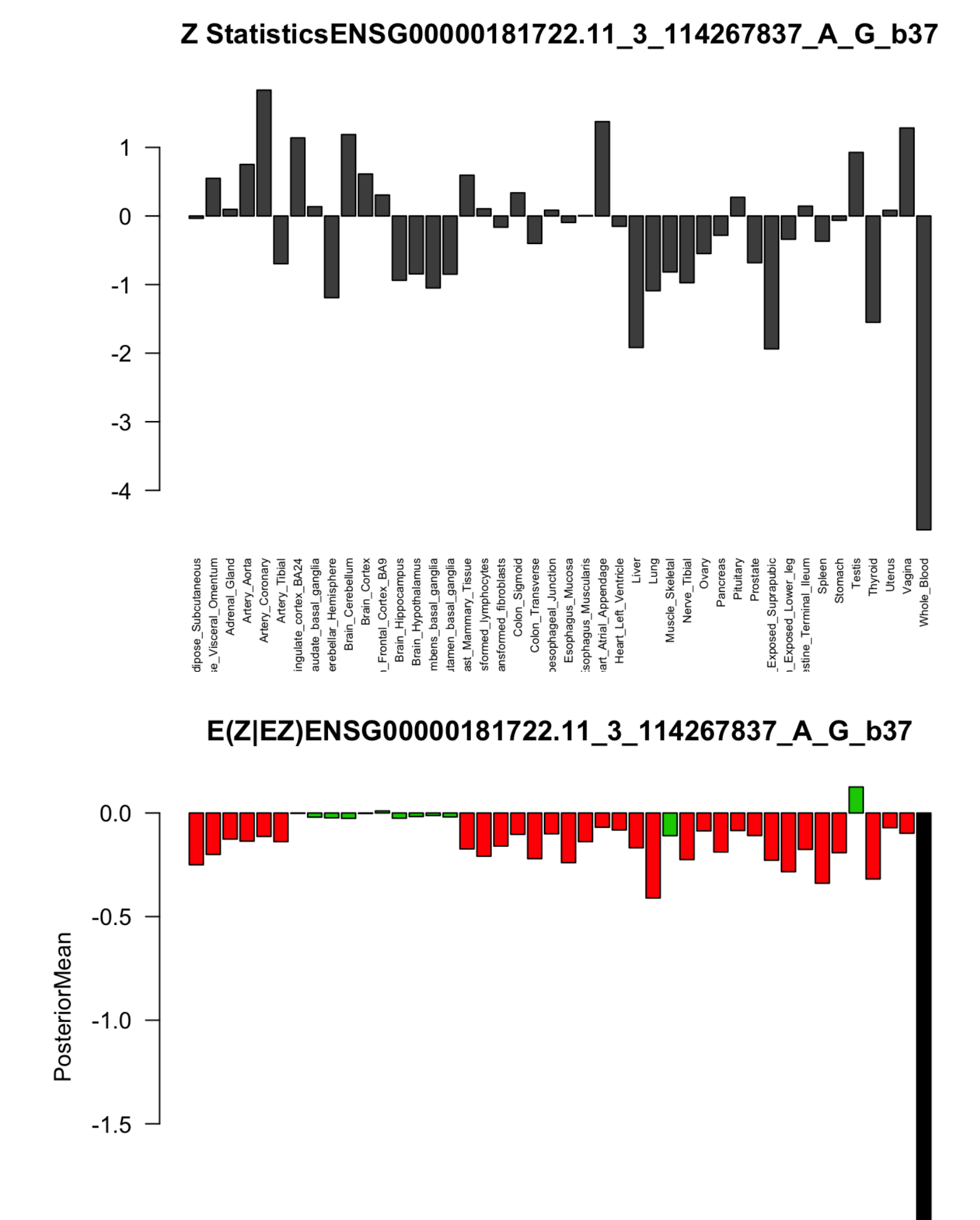


Figure 3: Small statistics are 'denoised' without losing strong tissue-specific effects

Conclusion



- Small positive effects in several tissues are pushed towards the overall negative effect
- Tissue specific effect in whole blood is allowed to remain: high prior belief in tissue-specific effects here.

... but we don't shrink strong effects - large effects preserved (see extremes of Figure 2)

Additional Information

- Learn about the overall heterogeneity of the data set
- For each gene snp pair, $Pr(\text{Inconsistent Signs})$ in 100 simulations from (3) according to $\tilde{\pi}$

