

THE UNIVERSITY OF CHICAGO

FLEXIBLE STATISTICAL METHODS FOR JOINTLY MODELING EFFECTS

A DISSERTATION SUBMITTED TO
THE FACULTY OF THE DIVISION OF THE BIOLOGICAL SCIENCES AND
THE PRITZKER SCHOOL OF MEDICINE
IN CANDIDACY FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY

DEPARTMENT OF HUMAN GENETICS

BY

SARAH MARGARET URBUT

CHICAGO, ILLINOIS

JUNE 2017

Copyright © 2017 by Sarah Margaret Urbut
All Rights Reserved

TABLE OF CONTENTS

LIST OF FIGURES	iv
LIST OF TABLES	viii
0.1 Abstract	x
1 INTRODUCTION	1
1.1 State of the Field	4
1.2 Method Application	7
2 METHOD DEVELOPMENT	10
2.1 Multivariate adaptive shrinkage (<i>mash</i>)	10
2.2 Detailed Methods	13
3 APPLICATION OF MASH TO JOINTLY MAP CIS-EQTL IN MULTIPLE TISSUES	38
4 MASHCOMMONBASELINE: COMPARING CONDITIONS TO A COMMON CON- TROL	74
4.1 Defining the Model	76
4.2 Likelihood with <i>mashcommonbaseline</i>	78
4.3 <i>mashcommonbaseline</i> Simulation	80
4.4 Simulation Design	81
4.5 Simulation with Signal	82
4.6 Application of Deconvolution	83
4.7 Real Data Analysis: Method Application	86
4.8 Results: Patterns of Sharing	87
4.9 Improvements over Existing Methods	91
5 APPLICATION OF MASH TO GWAS OF MULTIPLE PHENOTYPES . .	94
5.1 Consortium	94
5.2 Data Analysis Procedure	96
5.3 Results	100
5.4 Too much significance?	103
6 CONCLUSION	108
REFERENCES	114

LIST OF FIGURES

- 2.1 **Overview of fitting procedure in mash, which estimates the multivariate distribution of effects present in the data.** The data (**right**) consist of a matrix of effect size estimates for a large number of units (rows) in multiple conditions (columns), together with their corresponding standard errors (here assumed to be 1 for each effect for simplicity). Colors (red/blue) indicate the sign of the effects (positive/negative), with shading intensity indicating size of effect. First, using the rows containing the strongest signals (**left**), we apply covariance estimation and dimension-reduction methods to estimate candidate “data-driven” covariance matrices (here U_2, \dots, U_9). To these we add several “canonical” covariance matrices, including the identity matrix, and matrices representing condition-specific effects. Each covariance matrix represents a “pattern” of effects that may occur in the data (summarized visually here by the first eigenvector, although each matrix is actually $R \times R$). We then scale each covariance matrix by a grid of scaling factors ω_l , varying from “very small” to “very large”, which allow for effect sizes to range from very small to very large. Finally, using the whole data set (**right**), we use maximum likelihood estimation to estimate weights (relative frequencies) $\pi_{k,l}$ for each (ω_l, U_k) combination; this corresponds to estimating how commonly each pattern–effect size combination occurs. 14
- 3.1 **Comparison of methods on simulated data.** Results are shown for two simulation scenarios: “shared structured” effects, where the non-zero effects are shared among the 44 conditions in complex structured ways similar to patterns in the GTEx data; and “shared unstructured” effects, where the non-zero effects are shared among the 44 conditions with effect sizes that are independent among conditions. Panel (a) shows accuracy of effect estimates. Panels (b) and (c) show ROC curves for detecting significant effects. The primary difference between (b) and (c) is that in (b) each effect is treated as a separate discovery in each condition (which requires condition-specific measures of significance), whereas in (c) each effect is treated as a single discovery across all conditions (which requires only a single measure of significance, as in traditional meta-analyses). In (b) we require the estimated sign (+/-) of each significant effect to be correct to be considered a “true positive”. Our new method (**mash**) outperforms other methods, particularly for “shared structured” effects, a scenario expected to be common in genomics applications. 47

3.2	Summary of primary patterns identified by mash in GTEx data. Shown are the heatmap of the correlation matrix, and barplots of the first 3 eigenvectors, of the covariance matrix U_k corresponding to the dominant mixture component identified by mash . This component accounts for 34% of all weight in the GTEx data. In all cases, tissues are color-coded as indicated in the heatmap legend. The first eigenvector reflects broad sharing among tissues, with all effects in the same direction; the second eigenvector captures differences between brain (and, to a less extent, testis and pituitary) vs other tissues; the third eigenvector primarily captures effects that are stronger in whole blood than elsewhere.	51
3.3	Examples illustrating of how mash uses patterns of sharing to inform effect estimates in the GTEx data. In panel a) each colored dot shows the raw effect estimate for a single tissue (color-coded as in Figure 3.2), with grey bar indicating ± 2 standard errors. These are the data input into mash . Panel b) shows the corresponding estimates output by mash (posterior mean, ± 2 posterior standard deviations). In each case mash combines information across all tissues, using the background information – patterns of sharing – it has learned from data on all eQTLs, to produce more precise estimates. Together, these three examples illustrate the flexibility of mash in combining information across different subsets of tissues for different eQTLs, depending on how their data match different patterns of sharing identified in the overall data. See main text for detailed discussion.	54
3.4	Histogram showing estimated number of tissues in which top eQTLs are “shared” by two different definitions, a) sign and b) magnitude. Sharing by sign means that eQTL have the same sign of effect; Sharing by magnitude means that they also have similar effect size (within a factor of 2). Left: All tissues; Center: non-brain tissues; Right: brain tissues.	59
3.5	Pairwise sharing by magnitude of eQTL among tissues. For each pair of tissues we consider the top eQTLs that are significant in at least one of the two tissues, and plot the proportion of these that are “shared in magnitude” – that is, have effect estimates that are the same sign and within a factor of 2 of one another. Pink triangles highlight groups of biologically-related tissues mentioned in the text as showing particularly high levels of sharing.	60

4.1 **Intuition behind mashcommonbaseline.** We illustrate the necessity of our approach by first simulating a matrix in (a) \hat{C} in which each row is simulated independently, i.e., $\hat{c}_j \sim N(0, I_8)$. It's corresponding correlation matrix (c) is diagonal. After removing the control condition gene expression measurement from every subsequent condition by premultiplying by the 7x8 contrast matrix \tilde{L} , we then observe the matrix of estimate deviations $\hat{\Delta}$ shown in (b). The correlation matrix of this matrix (d) has non-zero diagonal elements, indicating correlation between the elements of $\hat{\delta}_j$ 79

4.2 **Improvement in both Power and Accuracy.** The number of false positives versus true positive associations detected using `mashcommonbaseline` (accounting for residual correlation among errors) improves. We also demonstrate improved accuracy of the estimated deviations in the (RRMSE) as in equation 2.22. 83

4.3 **Inferred Patterns of Sharing in data from [1]:** In the top panel (a), `cormotif` identifies three separate clusters that are collapsed into the primary patterns of sharing in `mashcommonbaseline` This figure is adapted from ([1]). We display the two most common patterns of sharing (b) and (c), U_k as defined by the proportion of hierarchical weight received. At left is a heatmap of each scaled covariance matrix (i.e., $U_k/\max(\text{diag}(U_k))$) and at right it's first eigenvector. The two patterns of sharing which receive the majority of the weight both represent deviations which are broadly shared in both sign among subgroups, though the deviations appear strongest in Yersinia and Salmonella, with weaker correlations among BCG and RV+. 88

4.4 **Visualizing Quantitative Sharing using mashcommonbaseline:** In panel a and b, we display the distribution of sharing by magnitude (left) and sign (right) in the tuberculin data. In (b) we display pairwise sharing by magnitude (left) and sign (right) using `mashcommonbaseline` as computed in 3.5 and 3.4 For each pair of conditions, we consider the effects that are significant in at least one of the tissues, and estimate the proportion that have effect sizes that are within 2-fold of one another and the same sign (a) or of the same sign only (b) 90

4.5	Examples illustrating how <code>mashcommonbaseline</code> can capture more subtle patterns of sharing than restricting effects to a binary outcome in conditions. Similar to Figure 3.3, we plot the original noisy estimates (± 2 SD) at left and the posterior means (± 2 Posterior SD) at right in each panel. In the top panel, <code>mashcommonbaseline</code> recognizes differential expression in all groups, while in the bottom panel, <code>mashcommonbaseline</code> can still preserve subgroup specificity when it exists.	92
5.1	Pairwise Sharing of Effects Across Traits At left we demonstrate the learned patterns which receive the majority of the weight.	102
5.2	Pairwise Sharing of Effects Across Traits At left we demonstrate pairwise sharing by sign (left) and absolute value of magnitude (right) across 16 GWAS traits.	103
5.3	Number of SNPs Found Significant We display the number of significant (i.e., $lfsr \leq 0.05$) in each condition. At right, we display the number of 'disease-specific' SNPs in any condition - i.e., those SNPs that are significant in that condition alone	104
5.4	Pruned Significance. We demonstrate the number of remaining significant results after the pruning exercise described in section 5.4.	105
5.5	Aberrations Above, we demonstrate effects for whom univariate analyses concluded lack of significance while our joint analysis found significance. Top left, we examine the distribution of initial univariate p-values for the identified hits. We then demonstrate a metaplot of the initial univariate $\hat{\beta}$ and $\hat{\sigma}^2$ across subgroups for SNPs initially identified as insignificant in CD, BP and FN (clockwise, top right).	106

LIST OF TABLES

3.1	Errors in estimates of sharing for simulated data For each method we computed the mean absolute error of the estimated number of conditions that share effects (by either sign or magnitude) with the condition with the largest estimated effects. Here “raw” indicates the performance of the raw estimates being input into each method. Each number is a mean error across the 400 non-null effects from the “shared structured effects” scenario.	46
3.2	Summary of sharing among top eQTLs. Numbers show the proportion of effects meeting a given sharing criterion. “Shared by sign” requires that the effect has the same sign as the strongest effect among tissues. “Shared by Magnitude” requires that the effect is also within a factor of 2 of the strongest effect. Numbers in parentheses are obtained by a secondary <code>mash</code> analysis of subsets of tissues.	58
3.3	Supplemental Table 1: Comparison of accuracy of effect size estimates for each method. Results show the RRMSE for all effects ($\text{RRMSE}^{\text{all}}$), and for the subsets of effects that are truly non-null ($\beta \neq 0$; $\text{RRMSE}^{\text{Non-null}}$) and truly null ($\beta = 0$, $\text{RRMSE}^{\text{Null}}$). Values of $\text{RRMSE}^{\text{Null}} < 1$ indicate how shrinkage towards zero is helping improve the estimates of null effects. Values of $\text{RRMSE}^{\text{Non-null}} < 1$ indicate how pooling information across conditions can improve accuracy of estimates of non-null effects. (In the Independent simulations the shrinkage of all methods improves overall performance, despite hurting performance for the non-null effects, because most effects are null.)	65
3.4	Supplemental Table 2: Comparison of Identified Associations in Simulations. Results show the overlap of associations identified by each combination of methods. In both “Shared” scenarios <code>mash</code> captures the vast majority of the associations identified by the other methods.	65
3.5	Supplemental Table 3: Comparison of Identified Associations in GTEx Data. Results show the overlap of identified associations with each method. <code>mash</code> captures the vast majority of the associations identified by the other methods, in addition to other associations.	66
3.6	Comparison of Methods	67
5.1	Studies Considered We describe the consortium and the disease referenced by each abbreviation.	95
5.2	Associations Identified	101

- 5.3 **Summary of Significance By Method** Numbers show the proportion of effects meeting a given sharing criterion. “Total Effects” reports the total proportion of significant effects, while ‘In at Least One’ requires that the effect be significant in at least one subgroup across conditions. . . . 101

0.1 Abstract

New statistical methods for analyzing genomic datasets that measure many effects in many conditions are introduced (e.g. gene expression changes under many treatments). These new methods improve on existing methods by allowing for arbitrary correlations in effects among conditions. This flexible approach increases power, improves effect-size estimates, and facilitates more quantitative assessments of effect-size heterogeneity than simple “shared/condition-specific” assessments. We illustrate these features through three detailed analyses. The first is an assessment of locally-acting (“cis”) eQTLs in 44 human tissues. Our analysis identifies more eQTLs than existing approaches, consistent with improved power. More importantly, although eQTLs are often shared broadly among tissues, our more quantitative approach highlights that effect sizes can vary considerably among tissues: some shared eQTLs show stronger effects in a subset of biologically-related tissues (e.g. brain-related tissues), or in only a single tissue (e.g. testis). We then apply our method to a setting in which all conditions are compared to a common control, as well as to an analysis considering the genetic effect on multiple diseases simultaneously. Our methods are widely applicable, computationally tractable for many conditions, and available at <https://github.com/stephenslab/mashr>.

CHAPTER 1

INTRODUCTION

Genomic studies often involve estimating and comparing many effects across multiple conditions or outcomes. Examples include studying changes in expression of many genes under multiple treatments [1]; or differences in histone methylation at many genomic locations in multiple cell lines [2]; or the effects of many genetic variants on risk of multiple diseases [3]; or the impact of many eQTLs in multiple cell-types or tissues [4–6].

While all these methods ask distinct biological questions, an initial goal is often to identify “significant” non-zero effects. When many conditions are considered, one might be interested in identifying effects as either ‘shared’ or specific to a particular condition. For example, in eQTL (expression Quantitative Trait Locus) studies, researchers are often interested in identifying tissue-specific effects, in the belief that they may have particular biological relevance. Another important goal is to compare effects, and to identify differences in effect among conditions. Such an analysis is of critical importance if one is to recognize that there exist degrees of continuous variation among effects that are considered ‘active’ in all conditions. As the power of a joint analysis increases, most effects might be called ‘shared’ in many conditions and the question of significance becomes less interesting than comparing among differences in effect magnitude and sign among conditions. All effects may then follow particular patterns, in which each pattern is characterized by a particular pattern of

sharing of effects among tissues. For instance, in eQTL analysis, some effects may tend to show large effects in certain tissues with more modest though significant effects in other tissues involved in a particular biological process, owing to their responsiveness to the biochemical mediators of the implicated process. Similarly, the same gene may show larger degrees of histone methylation in certain cell lines due to a shared need for silencing, while other cells may have smaller but non-trivial need.

The simplest, and perhaps most common, analysis strategy for such studies is to analyze the data in different conditions one at a time, and then compare the overlap of “significant” results in different conditions. Although appealingly simple, this “condition-by-condition” approach is unsatisfactory in several ways. For example, it can substantially under-represent sharing of effects among conditions, because many shared effects will be insignificant in some conditions just by chance. And when effects are shared among conditions it completely fails to exploit this, limiting its overall power [5].

To address these deficiencies of condition-by-condition analyses, several groups have developed methods for *joint* analysis of effects in multiple conditions (e.g. [2, 5–13]). The simplest of these methods build on traditional meta-analysis methodology [8, 9], and assume that the non-zero effects are present in every condition. Other methods are more flexible, allowing for condition-specific effects, for sharing of effects among subsets of conditions, and for heterogeneity in the shared effects [5, 6, 12]. Many of these methods also adapt themselves to the data at hand by learning patterns of sharing from the data, using a hierarchical model [5].

Nonetheless, existing methods remain limited in important ways. First, all of them make relatively restrictive assumptions about the correlations among non-zero effects. For example, [5] assumes correlations are non-negative, and that the non-zero effects are equally correlated among all conditions. In some applications correlations may be negative: for example, genetic variants that increase one trait may tend to decrease another. For tissues under parasympathetic and sympathetic control, such as the gut and skeletal muscle, an eQTL revealed under environmental stress may act to decrease and increase gene expression, respectively. For genetic variants implicated in hormone-related processes, we might imagine that tissues with opposing endocrine roles possess up-regulated and down-regulated gene expression. We find that some genetic effects associated positively with gene expression in the brain are negatively associated with gene expression in non-brain tissues, though sometimes this can be explained by LD (see discussion in “Linkage Disequilibrium”). And, often, some subsets of conditions will be more correlated than others: for example, in our eQTL application (below) effects in brain tissues are more correlated with one another than with effects in non-brain tissues. Second, the most flexible methods are computationally intractable for moderate numbers of conditions (e.g. 44 tissues in our eQTL application), and existing solutions to this problem substantially reduce flexibility. For example, eQTL-BMA in [5] solves the computational problem by restricting effects to be shared in all conditions, or specific to a single condition. Alternatively, `cormotif` [12] allows for all possible patterns of sharing in an elegant computationally-tractable way, but only under the more restrictive assumption that the non-zero effects are uncorrelated among conditions, which will often not hold

in practice. Third, and perhaps most importantly, existing methods typically focus only on *testing* for significant effects in each condition, and not on *estimating effect sizes*. As we illustrate here, estimating effect sizes can be essential to assessing heterogeneity of effects among conditions.

1.1 State of the Field

To give a more detailed understanding of available techniques, we briefly discuss three available methods in more detail. Current approaches to analyzing effects across multiple conditions typically resort to a configuration approach, focus primarily on measures of significance and fail to describe effect sizes. In a configuration in an eQTL study, for example, as described in [5], an effect is either active or inactive in a particular subgroup. With a limited number of subgroups (3 in their application), such an assumption is tractable as there are only $2^3 = 8$ possible configurations. However, [5] falls short in several important ways. First, it assumes that the distribution of effect sizes in all conditions in which the effect is active is equal within a configuration. Second, it assumes that the correlation among all non-zero effects is equivalent within a given configuration. Finally, it focuses on reporting a measure of significance for the effect and assigning the effect to a particular configuration, rather than making a statement about the effect size integrated over all possible configurations. Describing the effect size makes the identifiability of each pattern less important, because each vector of posterior effects integrates information across all patterns. Furthermore, reporting the effect recognizes that an effect can

be non-null in many conditions but with varying degrees of activity, giving way to the idea of quantitative heterogeneity inherent in biology.

For example, `corMotif`, a method designed to look for patterns of gene-expression ([12]), searches for a small number of latent probability vectors called correlation motifs to capture the major correlation patterns among multiple studies. Each vector corresponds to a set of probabilities for being differentially expressed in a particular subgroup. These probability vectors can then give way to a variety of configurations, making the method computationally tractable for large numbers of studies. However, as in [5], their method focuses on assigning genes to these classes and then ranking genes according to their probability of differential expression in a particular study [12]. This elucidates no information about effect sizes. Furthermore, they assume that the t-statistics and thus the effect sizes are uncorrelated among studies. For instance, an effect could be active in all conditions, but the effect size in one study tells us nothing about its effect in another.

Lastly, `metasoft` evaluates both random and fixed effect models to identify genetic associations. However, it reports only one estimated effect size rather than estimating effect sizes for each study, and again focuses on significance of a given gene. The authors focus again on testing for significance, rather than on describing the heterogeneity in effect sizes and signs among studies.

Finally, software implementations of existing methods are often tailored to a specific application, making it harder to apply to other settings: for example, `eQTL-BMA` [5] is primarily designed for eQTL applications, whereas `corMotif` [12] is tailored to

differential expression analyses. An exception here is the `metasoft` software [9, 14], which is designed to be *generic* in that it requires only effect estimates and their standard errors in multiple conditions, making it easily applicable to a wide range of settings.

Here we introduce more flexible statistical methods that combine the most attractive features of existing approaches, while overcoming their major limitations. The methods, which we refer to as “multivariate adaptive shrinkage” (`mash`), build on recent work in [15] for testing and estimation of effects in a *single* condition, and extend them to *multiple* conditions. Key features of `mash` include: i) It is *flexible*, allowing for both shared and condition-specific effects, and capable of capturing stronger correlations in effects among some conditions than others; ii) It is *computationally tractable* for hundreds of thousands of tests in (at least) dozens of conditions; iii) it provides not only measures of significance, but also *estimates of effect sizes*, together with measures of uncertainty; iv) It is *adaptive*, meaning that its behaviour adapts to the patterns present in the particular data set being analyzed; and v) It is *generic*, requiring only a matrix containing the observed effects in each condition, and a matrix of their corresponding standard errors. (Indeed `mash` can work with just a matrix of Z scores, although that reduces the ability to estimate effect sizes.) Together these features make `mash` the most flexible and widely-applicable method available for estimating and testing multiple effects in multiple conditions.

As its name suggests, `mash` is built on the statistical concept of “shrinkage”. Here shrinkage refers to modifying estimates towards some value – often towards zero –

to improve accuracy. There are many good justifications for shrinkage, and it is widely viewed as a powerful statistical tool. However, it is seldom used in genomics applications. This may be due to the difficulty of deciding precisely *how much* to shrink. The “adaptive shrinkage” method in [15] solves this problem in univariate settings by *learning from the data* how much to shrink. Here we extend this to multivariate settings. Shrinkage in the multivariate setting is more complex than in the univariate setting, but also potentially more useful. In particular, the multivariate setting provides the opportunity not only to shrink estimates towards zero (which improves accuracy if most effects are small), but also to shrink effects in related conditions towards one another (which improves accuracy when effects are similar among conditions). This focus on multivariate shrinkage estimation, and more generally on joint estimation of effects across multiple conditions, distinguishes **mash** from existing approaches that focus primarily on testing for non-zero effects. Estimation is particularly useful in settings where, as in our eQTL application here, there is considerable sharing of effects among conditions, but where effect sizes also vary considerably.

1.2 Method Application

To demonstrate the potential for **mash** to provide novel insights we first apply it here to analyse (*cis*) eQTL effects in 16,069 genes across 44 human tissues. Compared with previous analyses of human eQTLs among multiple tissues [4–6], our analysis involves many more tissues, and provides more insight into sharing of effects by examining

variation in eQTL effect sizes among tissues. Focussing on the strongest “cis” eQTLs in each gene – which are the easiest to reliably assess – we find that the majority are shared among large numbers of tissues, in that their effects tend to be consistent in sign (positive or negative) across tissues. However, at the same time, effect sizes can vary considerably among tissues. Reassuringly, biologically-related tissues tend to show more correlated effects; for example, effects are often quite similar among the different brain tissues. Our analyses of variation in estimated effects among tissues suggest that assessments of “tissue-specific” vs “tissue-consistent” effects should pay attention to effect sizes, and not only to tests of significance. The contents of this chapter arise from work with Wang and Stephens *et al* ([16]).

We then describe using the `mash` framework to estimate deviations from a control in the situation in which all subgroups are compared to the same control. We call this method `mashcommonbaseline`, because there exists no baseline within each condition but each condition is compared to the same control. Such methods are common in gene expression analysis, for example, in which differential expression is assessed in comparison. Our goal in this application is to reduce the number of false positives identified in any one subgroup due to the inherent correlation induced when comparing all conditions to a common control.

Lastly, we describe the application of `mash` to GWAS data in which we analyze summary statistics aggregated across 16 biological traits and diseases, and explore patterns of sharing among diseases and effects. We hope that a joint analysis will help to explain some of the missing heritability endemic in current GWAS literature,

and allow us to make a better description of the quantitative heterogeneity evident among genetic effects on human disease.

CHAPTER 2

METHOD DEVELOPMENT

2.1 Multivariate adaptive shrinkage (**mash**)

Our method, **mash**, is designed to estimate the effects of many units in many conditions (n units in R conditions say). It takes as its input two $n \times R$ matrices, one containing “effect” estimates and the other containing their corresponding standard errors. For example, in the GTEx data analyzed here we consider the effects of hundreds of thousands of potential eQTLs (rows) in $R = 44$ tissues (columns). The method assumes that the true effects are centered on 0, and indeed allows that many effects – possibly the vast majority – may be at, or very near, zero. That is, the true effect matrix may be sparse. It also allows that some of the non-zero effects may be ‘shared’, being similar (though not necessarily identical) among conditions, while others may be ‘specific’ to only a subset of conditions. Although we illustrate **mash** on an eQTL application, it is sufficiently flexible to apply to most contexts involving many multivariate effects.

The **mash** method is an Empirical Bayes method with two steps: i) use all the observed data to learn typical patterns of sparsity, sharing and correlations among effects; ii) use these learned patterns to produce improved effect estimates, and corresponding measures of significance, for each unit in each condition. Step ii) is reasonably straightforward: it involves applying Bayes theorem to combine the

background information (learned patterns of sharing from Step i)) with the observed data for each effect (the estimates and standard errors in every condition). Step i) is the difficult part, and where the primary innovations of our work lie. Specifically, we introduce a flexible model that allows for sparsity of effects and correlations among non-zero effects, and introduce a novel and efficient two-step approach to fitting this model.

Our flexible model uses a mixture of multivariate normal distributions that allows for a range of effect sizes and patterns of correlation. Specifically, each R -vector of effects across conditions, \mathbf{b} , is assumed to come from a mixture distribution,

$$p(\mathbf{b}; \boldsymbol{\pi}, \mathbf{U}) = \sum_{k=1}^K \sum_{l=1}^L \pi_{k,l} N_R(\mathbf{b}; \mathbf{0}, \omega_l U_k), \quad (2.1)$$

where $N_R(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the multivariate normal density in R dimensions with mean $\boldsymbol{\mu}$ and variance covariance matrix $\boldsymbol{\Sigma}$; each U_k is a covariance matrix that captures some common “pattern” of (potentially-correlated) effects; each ω_l is a scalar scaling coefficient that corresponds to a different “size” of effect; and the mixture proportions $\pi_{k,l}$ determine the relative frequency of each pattern-size combination. The scaling coefficients ω_l take values on a fixed dense grid that spans “very small” to “very large”, to capture the full range of effects that could occur (the goal is that the grid is sufficiently large and dense that adding more values to it will not change results; see [15]).

To fit this model, we use a novel two-step procedure illustrated in Figure 2.1:

i-a) Generate a large list of candidate covariance matrices $U_k = (U_1, \dots, U_K)$.

This list includes both “data-driven” estimates, and “canonical” matrices that have simple interpretations. The data-driven estimates are obtained by applying covariance estimation methods [17], and dimension reduction techniques (e.g. Principal components analysis, and sparse factor analysis [18]) to a subset of the effects matrix, specifically the rows of the effect matrix that have the largest (univariate) effects. The canonical matrices we use include the identity matrix (representing independent effects across conditions); a matrix of all 1s (representing effects that are equal in all conditions); and R matrices that represent effects that are specific to condition r ($r = 1, \dots, R$). See Detailed Methods for details.

i-b) Given this list, estimate $\boldsymbol{\pi}$ by maximum likelihood (using *all* observed effects, not only those used in Step i-a)).

The intuition is that Step i-a) can be relatively *ad hoc*, with the goal of producing a large list of matrices, only some of which may effectively capture key patterns in the data. Step i-b) is more formal, being based on the principle of maximum likelihood, and can rescue imperfections in Step i-a) by giving very low weight to covariance matrices that are not well supported by the data. Step i-b) is also the place where the overall sparsity of effects is taken account of: if most effects are zero, or very small, then this step will put most weight on very small effects (i.e. small scaling coefficients, ω). This modular approach has several attractive features. For example, Step i-b) is a convex optimization problem, and so can be solved efficiently

and reliably for large problems. And if researchers have ideas for additional ways to generate candidate matrices in Step i-a), these are easily plugged into the procedure.

The model (2.1) is quite flexible, and includes many existing methods for this problem as special cases (Detailed Methods). One potential drawback of flexible models is the possibility of “overfitting”. To address this we used a cross-validation procedure which trains the model on a random subset of the data (rows of the matrix) and then assesses its fit on the remaining data (“test data”). In practice we found overfitting not to be a major concern - that is, in general, we found that using more U_k typically improved, or at least did not harm, test set performance (e.g. Supplementary Figure 7). Thus, although `mash` is flexible, it is not *too* flexible. A still more flexible model could be obtained by estimating the means of the multivariate normal distributions in (2.1), rather than setting them to 0, but this would substantially increase the potential for overfitting.

2.2 Detailed Methods

Model and Fitting

Let b_{jr} ($j = 1, \dots, J; r = 1, \dots, R$) denote the true value of effect j in condition r . Further let \hat{b}_{jr} denote the (observed) estimate of this effect, and \hat{s}_{jr} the standard error of this estimate, so $z_{jr} := \hat{b}_{jr}/\hat{s}_{jr}$ is the usual z statistic for testing whether

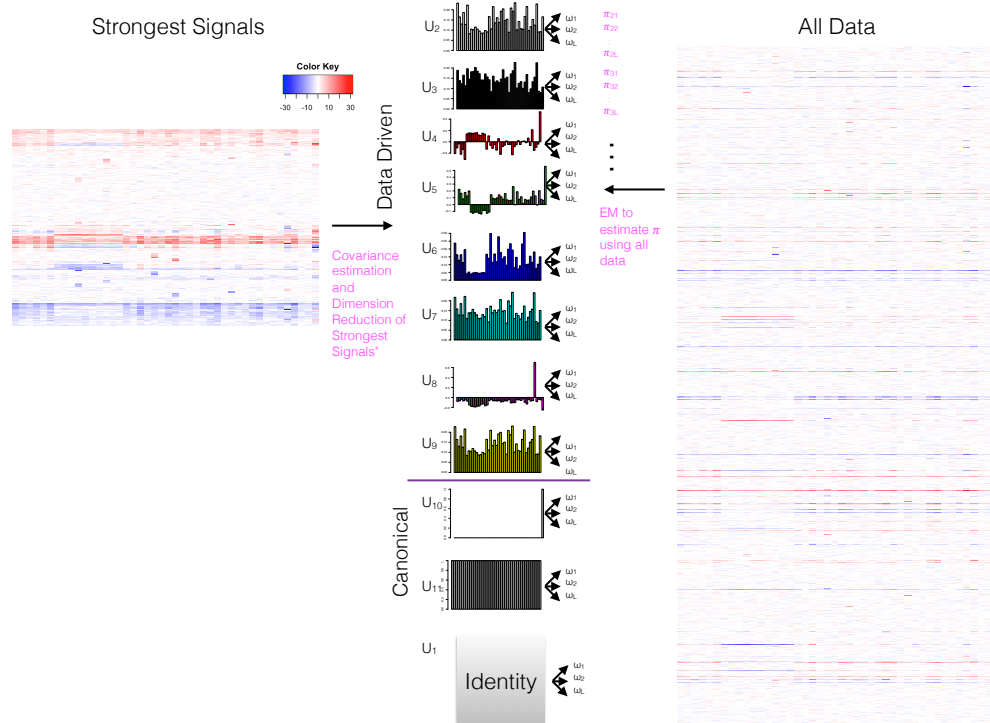


Figure 2.1: **Overview of fitting procedure in mash, which estimates the multivariate distribution of effects present in the data.** The data (**right**) consist of a matrix of effect size estimates for a large number of units (rows) in multiple conditions (columns), together with their corresponding standard errors (here assumed to be 1 for each effect for simplicity). Colors (red/blue) indicate the sign of the effects (positive/negative), with shading intensity indicating size of effect. First, using the rows containing the strongest signals (**left**), we apply covariance estimation and dimension-reduction methods to estimate candidate “data-driven” covariance matrices (here U_2, \dots, U_9). To these we add several “canonical” covariance matrices, including the identity matrix, and matrices representing condition-specific effects. Each covariance matrix represents a “pattern” of effects that may occur in the data (summarized visually here by the first eigenvector, although each matrix is actually $R \times R$). We then scale each covariance matrix by a grid of scaling factors ω_l , varying from “very small” to “very large”, which allow for effect sizes to range from very small to very large. Finally, using the whole data set (**right**), we use maximum likelihood estimation to estimate weights (relative frequencies) $\pi_{k,l}$ for each (ω_l, U_k) combination; this corresponds to estimating how commonly each pattern–effect size combination occurs.

b_{jr} is zero. Let B, \hat{B}, S and Z denote the corresponding $J \times R$ matrices, and let \mathbf{b}_j (respectively $\hat{\mathbf{b}}_j, \mathbf{z}_j$) denote the j th row of B (respectively \hat{B}, Z).

We assume the vector $\hat{\mathbf{b}}_j$ is normally distributed about the true effects \mathbf{b}_j , with variance-covariance matrix V_j (defined below), and that the true effects follow (2.1).

That is,

$$p(\hat{\mathbf{b}}_j | \mathbf{b}_j, V_j) = N_R(\hat{\mathbf{b}}_j; \mathbf{b}_j, V_j), \quad (2.2)$$

$$p(\mathbf{b}_j | \boldsymbol{\pi}, \mathbf{U}) = \sum_{k=1}^K \sum_{l=1}^L \pi_{k,l} N_R(\mathbf{b}_j; \mathbf{0}, \omega_l U_k). \quad (2.3)$$

where $N_R(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the density of the R -dimensional multivariate normal (MVN) distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, and the scaling parameters $\omega_1, \dots, \omega_L$ are fixed on a dense grid (detailed below). Combining these two implies that the marginal distribution of $\hat{\mathbf{b}}_j$, integrating out \mathbf{b}_j , is

$$p(\hat{\mathbf{b}}_j | \boldsymbol{\pi}, \mathbf{U}, V_j) = \sum_{k=1}^K \sum_{l=1}^L \pi_{k,l} N_R(\hat{\mathbf{b}}_j; \mathbf{0}, \omega_l U_k + V_j). \quad (2.4)$$

This last equation comes from the fact that the sum of two MVNs is MVN.

Here the covariance matrix V_j is given by $V_j = S_j C S_j$ where C is a correlation matrix that accounts for correlations among the measurements in the R conditions, and S_j is the $R \times R$ diagonal matrix with diagonal elements $(\hat{s}_{j1}, \dots, \hat{s}_{jR})$. In settings where measurements in the R conditions are independent one would set $C = \mathbf{I}_R$, the $R \times R$ identity matrix, so $V_j = S_j^2$. However, in our GTEEx analysis the measurements

are correlated due to sample overlap (some individuals in common) among tissues; we estimate this correlation from the data (see Section “Estimating the correlation matrix C ”). The methods implemented here can be applied for any specified matrices V_j .

The two steps of `mash` are:

- i) Estimate $\mathbf{U}, \boldsymbol{\pi}$. This involves two substeps:
 - a) Create a list of both data-driven and canonical covariance matrices, $\hat{\mathbf{U}}$.
 - b) Given $\hat{\mathbf{U}}$, estimate $\boldsymbol{\pi}$ by maximum likelihood. (A key idea here is that if some matrices generated in a) do not help capture patterns in the data then they will receive little weight.) Let $\hat{\boldsymbol{\pi}}$ denote this estimate.
- ii) Compute, for each j , the posterior distribution $p(\mathbf{b}_j | \hat{\mathbf{b}}_j, \hat{\mathbf{U}}, \hat{\boldsymbol{\pi}}, V_j)$.

These steps are now detailed in turn.

Generate data-driven covariance matrices U_k

We first identify rows j of the matrix \hat{B} that likely have an effect in at least one condition. For example, in the GTEx data we chose the rows corresponding to the “top” SNP for each gene, which we define to be the SNP with the highest value of Z_j^{\max} where

$$Z_j^{\max} := \max_r \hat{\mathbf{b}}_{jr} / \hat{s}_{jr}. \quad (2.5)$$

(We used max here, rather than, say, the sum, to try to include effects that are very strong in a single condition and not only effects that are shared among conditions.) For the simulated data we ran the univariate adaptive shrinkage method `ash` on the data in each condition r separately, and computed $lfsr_{jr}$ for each effect j . We then chose the rows j for which at least one of the conditions showed a significant effect in this univariate analyses ($\min_r lfsr_{jr} < 0.05$).

Next we fit a mixture of MVN distributions to these strongest effects, using methods from [17]. Specifically results in [17] provide an EM algorithm for fitting a model very similar to (2.3) – (2.2) with the crucial difference that there is no scaling parameters on the covariances. That is,

$$p(\mathbf{b}_j | \boldsymbol{\pi}, \mathbf{U}) = \sum_k \pi_k N_R(\mathbf{b}_j; \mathbf{0}, U_k). \quad (2.6)$$

The absence of the scaling factors ω_l means that, compared with `mash`, the model (5.2) is less well suited to capture effects that have similar patterns (relative sizes across conditions) but vary in magnitude. However, by applying it here to only the largest effects we seek to sidestep this issue. Estimates of U_k from this EM algorithm are sensitive to initialization. Furthermore, we noticed an interesting feature of the EM algorithm: each iteration preserves the rank of the matrices U_k , so the ranks of the estimated matrices are the same as the ranks of the matrices used to initialize the algorithm. We exploited this fact by including low-rank matrices in our initialization to ensure that some of the estimated U_k are low-rank matrices. This helps stabilize the estimates since rank-penalization is one way to regularize

covariance matrix estimation.

To describe the initialization in detail, let \tilde{J} denote the number of “strongest effects” selected above, and let \tilde{Z} denote the column-centered $\tilde{J} \times R$ matrix of Z scores for these “strong effects”. To attempt to extract the main patterns in \tilde{Z} we perform dimension reduction on \tilde{Z} . Specifically we apply Principal Component Analysis (through Singular Value Decomposition, SVD) and Sparse Factor Analysis (SFA; [18]) to \tilde{Z} .

SVD yields a set of eigenvalues and eigenvectors of \tilde{Z} . Let λ_p, v_p denote the p th eigenvalue and corresponding (right) eigenvector. (So v_p is an R vector for $p = 1, \dots, R$.)

SFA yields a representation

$$\tilde{Z} = LF + E \tag{2.7}$$

where L is a sparse $J \times Q$ matrix of loadings, and F is a $Q \times R$ matrix of factors. Here we used $Q = 5$.

Given this we initialized the EM with $K = 3$ and

- $\tilde{U}_1 = \frac{1}{j} \tilde{Z}' \tilde{Z}$, the empirical covariance matrix of \tilde{Z} .
- $\tilde{U}_2 = \frac{1}{j} \sum_{p=1}^P \lambda_p^2 v_p v_p'$, which is a rank P approximation of the covariance matrix of \tilde{Z} . Here we used $P = 3$.
- $\tilde{U}_3 = \frac{1}{j} (LF)'(LF)$ which is a rank Q approximation of the covariance matrix

of \tilde{Z} .

In addition to the covariance matrices obtained from this EM algorithm, we added some more matrices based on the SFA results, specifically

- The 5 matrices $F'_q L'_q L_q F'_q$, which are each rank 1 matrices that reflect the effects captured by the q th factor in the SFA analysis ($q = 1, \dots, 5$).

The rationale here is that the factors in the factor analysis may directly reflect effect patterns in the data, and if so then these matrices will be a helpful addition. (We view such additions as a low-risk, because If they are not helpful then they will receive little weight when we estimate $\boldsymbol{\pi}$).

In total this procedure produces 8 data-driven covariance matrices for our GTEx analyses.

Generate canonical covariance matrices U_k

To these “data-driven” covariance matrices we add the following “canonical” matrices:

1. The matrix \mathbf{I}_R . This represents the situation where the effects in different conditions are independent, which may be unlikely in some applications (like the GTEx application here), but seems useful to include if only to exclude it.
2. The R rank-1 matrices $\mathbf{e}_r \mathbf{e}'_r$ where \mathbf{e}_r denotes the unit vector with 0s every-

where except for element r which is a 1. These represents effects that occur only in a single condition.

3. The rank-1 matrix $\mathbf{1}\mathbf{1}'$ where $\mathbf{1}$ denotes the R -vector of 1s. That is, the matrix of all 1s. This represents effects that are identical among all conditions.

The user can, if desired, add additional canonical matrices. For example, if R is moderate then one could consider adding the 2^R canonical matrices that correspond to shared (equal) effects in each of the 2^R subsets of conditions.

In total this procedure produces 46 canonical covariance matrices for our GTE_x analyses.

Standardize covariance matrices

Since (2.3) uses the same grid of scaling factors ω we standardize the matrices U_k obtained above so that they are similar in scale. Specifically, for each k , we divide every element of U_k by the maximum diagonal element of U_k (so that the maximum diagonal element of the rescaled matrix is one). These rescaled matrices provide the \hat{U} , completing step i)-a of `mash`.

Define grid of ω_l values

We choose a dense grid of ω_l ranging from “very small” to “very large”. [15] provides a specific way to select suitable limits $(\omega_{\min}, \omega_{\max})$ for this grid in the univariate case; we simply apply this method to each condition r in turn and take the smallest ω_{\min} and the largest of the ω_{\max} as the grid limits. The internal points of the grid are then obtained as in the univariate case [15], by setting $\omega_l = \omega_{\max}/m^{l-1}$, for $l = 1, \dots, L$, where $m > 1$ is a user-tunable parameter that affects the grid density and L is chosen to be just large enough so that $\omega_L < \omega_{\min}$. Our default choice of grid density is $m = \sqrt{2}$. In principle the grid should be made sufficiently dense that increasing its density would not change the answers obtained. In the GTE_x data we found results with $m = \sqrt{2}$ provided similar results to $m = 2$, supporting this choice.

Estimate $\boldsymbol{\pi}$ by maximum likelihood

Given $\hat{\boldsymbol{U}}, \boldsymbol{\omega}$, we estimate the mixture proportions $\boldsymbol{\pi}$ by maximum likelihood.

To simplify notation, let $\Sigma_{k,l} := \omega_l \hat{U}_k$, and replace the double index k, l with a single index p which ranges from 1 to $P := KL$. In this notation the prior (2.3) becomes

$$p(\mathbf{b}_j | \boldsymbol{\pi}, \boldsymbol{\Sigma}) = \sum_p \pi_p N_R(\mathbf{b}_j; \mathbf{0}, \Sigma_p), \quad (2.8)$$

and (2.4) becomes

$$p(\hat{\mathbf{b}}_j | \boldsymbol{\pi}, V, \boldsymbol{\Sigma}) = \sum_p^P \pi_p N_R(\hat{\mathbf{b}}_j; \mathbf{0}, \Sigma_p + V_j). \quad (2.9)$$

Assuming independence of rows of \hat{B} , the likelihood for $\boldsymbol{\pi}$ is given by

$$\begin{aligned} L(\boldsymbol{\pi}) &:= p(\hat{B} | \boldsymbol{\pi}, V, \boldsymbol{\Sigma}) \\ &= \prod_{j=1}^J p(\hat{\mathbf{b}}_j | \boldsymbol{\pi}, V, \boldsymbol{\Sigma}) \\ &= \prod_{j=1}^J \sum_p^P \pi_p N_R(\hat{\mathbf{b}}_j; \mathbf{0}, \Sigma_p + V_j). \end{aligned} \quad (2.10)$$

If the rows of \hat{B} are not independent then this may be interpreted as a “composite likelihood” [19]. By conditioning on V here, rather than treating it as part of the data, we are using a multivariate analogue of the approximation in [20].

Maximising this likelihood over $\boldsymbol{\pi}$ is a convex optimization problem, which here we solve using an EM algorithm [21], accelerated using SQUAREM [22]. This optimization problem is identical to the optimization over $\boldsymbol{\pi}$ in the univariate setting ($R = 1$) in [15], but involves a much larger number of components. If the matrix \hat{B} has many rows then to reduce computation time we can fit the model using a random subset of rows. For example, we used 20,000 rows in our GTE_x application. (It is important that this is a random subset, and not the \tilde{J} rows of strong effects used to generate the data-driven \hat{U}_k ; use of the strong effects in this step would be a mistake as it

would bias estimates of $\boldsymbol{\pi}$ towards large effect sizes.)

Posterior Calculations

To specify the posterior distributions, recall the following standard result for Bayesian analysis of an R -dimensional MVN. If $\mathbf{b} \sim N_R(\mathbf{0}, U)$, and $\hat{\mathbf{b}}|\mathbf{b} \sim N_R(\mathbf{b}, V)$ then

$$\mathbf{b}|\hat{\mathbf{b}} \sim N_R(\tilde{\boldsymbol{\mu}}, \tilde{U}), \quad (2.11)$$

where:

$$\tilde{U} = \tilde{U}(U, V) := (U^{-1} + V^{-1})^{-1}, \quad (2.12)$$

$$\tilde{\boldsymbol{\mu}} = \tilde{\boldsymbol{\mu}}(U, V, \hat{\mathbf{b}}) := \tilde{U}(U, V)V^{-1}\hat{\mathbf{b}}. \quad (2.13)$$

This result is easily extended to the case where the prior on \mathbf{b} is a mixture of MVNs (2.3). In this case the posterior distribution is simply a mixture of MVNs:

$$p(\mathbf{b}_j|\hat{\mathbf{b}}_j, \hat{V}_j, \hat{\boldsymbol{\pi}}) = \sum_p^P \tilde{\pi}_{jp} N_R(\mathbf{b}_j; \tilde{\boldsymbol{\mu}}_{jp}, \tilde{U}_{jp}) \quad (2.14)$$

where $\tilde{\boldsymbol{\mu}}_{jp} = \tilde{\boldsymbol{\mu}}(\Sigma_p, V_j, \hat{\mathbf{b}}_j)$ (equation (2.13)), $\tilde{U}_{jp} = \tilde{U}(\Sigma_p, V_j)$ (equation (2.12)), and

$$\tilde{\pi}_{jp} = \frac{\hat{\pi}_p N_R(\hat{\mathbf{b}}_j; \mathbf{0}, \Sigma_p + V_j)}{\sum_{p=1}^P \hat{\pi}_p N_R(\hat{\mathbf{b}}_j; \mathbf{0}, \Sigma_p + V_j)}. \quad (2.15)$$

From this is is straightforward to compute the posterior mean

$$\hat{\mathbf{b}}_j := E(\mathbf{b}_j | \hat{\mathbf{b}}_j, \hat{V}_j, \hat{\boldsymbol{\pi}}) = \sum_p^P \tilde{\pi}_{jp} \tilde{\boldsymbol{\mu}}_{jp}, \quad (2.16)$$

and posterior variance

$$\text{Var}(b_{jr} | \hat{\mathbf{b}}_j, \hat{V}_j, \hat{\boldsymbol{\pi}}) = \sum_{p=1}^P \tilde{\pi}_p (\tilde{U}_{jp,rr} + \tilde{\boldsymbol{\mu}}_{jp,r}^2) - [\sum_p^P \tilde{\pi}_{jp} \tilde{\boldsymbol{\mu}}_{jp,r}]^2, \quad (2.17)$$

as well as the local false sign rate.

Local False Sign Rate

To measure “significance” of an estimated effect β_{jr} we use the “local false sign rate” [15]:

$$lfsr_{jr} := \min[\Pr(\beta_{jr} \geq 0 | D), \Pr(\beta_{jr} \leq 0 | D)] \quad (2.18)$$

where D denotes all the available data. More intuitively, $lfsr_{jr}$ is the probability that we would get the sign of the effect β_{jr} incorrect if we were to use our best guess of the sign (positive or negative). Thus a small $lfsr$ indicates high confidence in the sign of an effect. The $lfsr$ is more conservative than its analogue, the local false discovery rate [23], because requiring confidence in the sign of an effect is more stringent than requiring confidence that it be non-zero. More importantly the $lfsr$ is more robust to modeling assumptions than the lfd [15], a particularly important issue in multivariate analyses where modeling assumptions inevitably play a larger

role.

Bayes Factors testing Global Null

Although not our primary focus, it is straightforward to use the fitted model to compute Bayes Factors for the alternative model ($\mathbf{b} \neq 0$) vs the null model $\mathbf{b} = 0$. Specifically

$$\text{BF}_j = p(\hat{\mathbf{b}}_j | \hat{\mathbf{U}}, \hat{\boldsymbol{\pi}}, V_j) / p(\hat{\mathbf{b}}_j | \mathbf{b} = 0, V_j) \quad (2.19)$$

where the numerator is given by (2.4) and the denominator by (2.2) with $\mathbf{b} = 0$.

The EZ model, and applying mash to Z scores

The model (2.3) assumes \mathbf{b}_j are independent of their standard errors V_j . We refer to this as the “exchangeable effects” (EE) model [24]. An alternative assumption is to allow that the effects may scale with standard error, so that effects with larger standard error tend to be larger. That is:

$$S_j^{-1} \mathbf{b}_j | \boldsymbol{\pi}, \mathbf{U}, \boldsymbol{\omega}, \mathbf{V}_j \sim \mathbf{g}(), \quad (2.20)$$

where $g()$ represents the mixture of multivariate normal distributions in (2.3). We refer to (2.20) as the “Exchangeable Z” (EZ) model, because the left of this equation is the vector of Z scores for effect j .

As described in [15], this EZ model can be fit by applying exactly the same code as the EE model to the Z statistics, with the standard errors of the Z statistics set to be 1. (That is, set $\hat{B} = Z$ and $\hat{s}_{jr} = 1$.) One advantage of this model is that it can be fit using only the Z scores, and does not require access to both the estimates and their standard errors. The $lfsr$ can also be computed using only the Z scores. However, the posterior mean estimates that arise from this model are estimates of $S_j^{-1}\mathbf{b}_j$; transforming these to estimates of effect sizes \mathbf{b}_j requires knowledge of S_j .

We analyzed the GTEx data using both EE and EZ models. Results were qualitatively similar in terms of patterns of sharing, but the EZ model performed better in cross-validation tests of model fit (see below), and so we report results from that model.

Estimating the correlation matrix C

To estimate the correlation matrix C we exploit the fact that C is the correlation matrix of the Z scores \mathbf{z}_j under the null ($\mathbf{b}_j = 0$). Specifically we estimate C using the empirical correlation matrix of the z scores for the effects j that are most consistent with the null, $\mathcal{N} = \{j : \max_r(|z_{jr}|) < 2\}$:

$$\hat{C} = \frac{1}{|\mathcal{N}|} \sum_{j \in \mathcal{N}} \mathbf{z}_j \mathbf{z}_j' \tag{2.21}$$

For the GTEx data the measurements in different tissues are not very highly corre-

lated: all elements of the estimated C were < 0.25 and 95% were < 0.1 . However, in cross-validation tests (below) this estimated C produced better model fit than ignoring correlations ($C = I_R$).

Cross-validation of model fit

To compare the performance of different strategies for selecting the covariance matrices U_k we use a cross-validation-based approach to assess model fit. In brief, this involves first dividing the data matrix into two groups by selecting half the rows to form the “training data”, with the remaining rows forming the “test data”. We then apply `mash`, as above, to the training data: use the strongest effects to select candidate U_k , and then learn the weights $\pi_{k,l}$ from all the training data (or a random subset if the data are large; we used 20,000 effects in our analysis). This provides an estimate of the distribution of effects \hat{g} . We assess the “fit” of this estimated g by how well it predicts the test data. That is, by computing $p(\hat{B}|\hat{\boldsymbol{\pi}}, V, \hat{U})$, given by (2.2), for the test data.

This strategy facilitates experimentation with ways to estimate \hat{U} . In particular, if new ways to generate \hat{U} are suggested then their effectiveness can be assessed using this strategy. Our current strategy described above was developed and refined using this framework. (However, performance of `mash` is relatively robust to the addition of poorly-estimated U_k because they are typically estimated to have small weight.)

When applying this strategy to the GTEx data we created the test and training data

by randomly selecting half the *genes*, rather than half the rows (gene-SNP pairs). Specifically we used genes on even-numbered chromosomes as the training set, and genes on odd-numbered chromosomes as the test set. This ensures that rows in the test set are independent of rows in the training set.

Visualizing U_k

In our application to the GTEx data $R = 44$, so each U_k is a 44 by 44 covariance matrix, and each component of the mixture (2.1) is a distribution in 44 dimensions. Visualizing such a distribution is challenging, but we can get some insight from the first eigenvector of U_k , v_k say, which captures the principal direction of the effects in component k . If U_k is dominated by this principal direction then we can think of effects from that component as being of the form λv_k for some scalar λ . For example, if the elements of the vector v_k are approximately equal then component k captures effects that are approximately equal in all conditions. Or, if v_k has one large element, with other elements close to 0, then component k corresponds to an effect that is strong in only one condition. See Figure 3.2 for illustration.

Relationship with existing methods

The `mash` method essentially includes many existing methods for joint analysis of multiple effects as special cases. Specifically, many existing methods correspond to making particular choices for the “canonical” covariance matrices \mathbf{U} (and excluding

the data-driven covariance matrices). For example, a simple “fixed effects” meta-analysis – which assumes equal effects in all conditions – corresponds to $K = 1$ with $U_1 = 11'$ (the matrix with all entries 1). (This covariance matrix is singular, but this is allowed within `mash`). A more flexible assumption is that effects in different conditions are normally distributed about some mean, and this also corresponds to a multivariate normal assumption if the mean is assumed to be normally distributed [24]. More flexible still are models that allow that effects may be exactly zero in some subset of conditions, as in [5,6]. These models correspond to using (singular) covariances U_k with 0s in the rows and columns corresponding to the subset of conditions with zero effect.

However, `mash` also goes beyond these previous methods in two ways. First, `mash` includes a large number of scaling coefficients ω_l , which allows it to flexibly capture a range of effect distributions (see [15]). Second, and perhaps more important, `mash` includes data-driven covariance matrices (Step i-a)), making it more flexible and adaptive to patterns in the specific data being analyzed. This innovation is particularly helpful in settings with moderately large R (e.g., in our application here $R = 44$) where it becomes impractical to pre-specify canonical matrices for all patterns of sharing that might occur. For example, [5,6] consider all 2^R different combinations of sparsity in the effects, which works for $R = 9$ [25], but is impractical for $R = 44$. While it is possible to restrict the number of combinations considered (e.g. BMALite in [5]), this comes at an obvious cost in flexibility. The addition of data-driven covariance matrices helps rectify this problem, making `mash` both flexible and computationally tractable for moderately large R .

Definitions of various quantities

RRMSE (accuracy of estimates in simulation studies)

The RRMSEs for estimates \hat{b}_{jr} of b_{jr} reported in Figure 3.1a are computed as

$$\text{RRMSE} = \frac{\sqrt{E((b_{jr} - \hat{b}_{jr}))^2)}}{\sqrt{E((b_{jr} - \hat{b}_{jr})^2)}}. \quad (2.22)$$

ROC curves

For the ROC curves in Figure 3.1b the True Positive Rate and False Positive Rate are computed at any given threshold t as

$$\text{True Positive Rate} := \frac{|CS \cap S|}{|T|} \quad (2.23)$$

$$\text{False Positive Rate} := \frac{|N \cap S|}{|N|} \quad (2.24)$$

where S is the set of significant results at threshold t , CS the set of correctly-signed results, T the set of true (non-zero) effects and N the set of null effects:

$$S := \{j, r : lfsr_{jr} \leq t\}, \quad (2.25)$$

$$CS := \{j, r : E(b_{jr}|D) \times b_{jr} > 0\}, \quad (2.26)$$

$$N := \{j, r : b_{jr} = 0\} \quad (2.27)$$

$$T := \{j, r : b_{jr} \neq 0\}. \quad (2.28)$$

(Thus, to be considered a true positive, we require that the effect be correctly signed and not only significant.)

For the ROC curves in Figure 3.1b the True Positive Rate and False Positive Rate are computed based on treating whole rows j as discoveries. For example, suppose a method produces a p value p_j for testing row j . Then at any threshold t the TPR and FPR are:

$$\text{True Positive Rate} := \frac{|\tilde{S}|}{|\tilde{T}|} \quad (2.29)$$

$$\text{False Positive Rate} := \frac{|\tilde{S}|}{|\tilde{N}|} \quad (2.30)$$

where \tilde{S} is the set of significant results at threshold t , \tilde{T} the set of true (non-zero)

effects and \tilde{N} the set of null effects:

$$\tilde{S} := \{j : p_j \leq t\}, \quad (2.31)$$

$$\tilde{N} := \{j : b_{jr} = 0 \quad \forall r\} \quad (2.32)$$

$$\tilde{T} := \{j : \exists r \text{ s.t. } b_{jr} \neq 0\}. \quad (2.33)$$

Effective sample size

We define the effect sample size for tissue r as

$$n_r^{\text{eff}} := n_r^{\text{orig}} \text{median}_j \frac{\hat{s}_{jr}^2}{\tilde{s}_{jr}^2} \quad (2.34)$$

where \hat{s}_{jr} is the standard error and \tilde{s}_{jr} is the posterior standard deviation for effect j in tissue r .

Normalized effects

We define the normalized effect \tilde{b} in each condition as the ratio of its effect in that condition to the largest effect across all conditions:

$$\tilde{b}_{jr} = \frac{b_{jr}}{b_{jr_0}} \quad (2.35)$$

where

$$r_0 = \arg \max_r |b_{jr}| \quad (2.36)$$

For example, in our eQTL context, a normalized effect $\tilde{b}_{jr} = 0.5$ means that the effect of eQTL j in tissue r is half that of its effect in the strongest tissue.

Pairwise Sharing

To assess pairwise sharing in sign between tissues r and s (Supplementary Figure 3.4) we compute, for QTL that are significant ($lfsr < 0.05$) in at least one of r and s , the fraction that have effect estimates that are of the same sign.

To assess pairwise sharing in magnitude between tissues r and s (Figure 3.5) we compute, for QTL that are significant ($lfsr < 0.05$) in at least one of r and s , the fraction that have effect estimates that are within a factor of 2 of one another.

That is, let

$$\text{QTL}_r := \{j : lfsr_{jr} < 0.05\} \quad (2.37)$$

$$\text{SS}_{rs} := \{j : \text{sign}(\hat{b}_{jr}) = \text{sign}(\hat{b}_{js})\} \quad (2.38)$$

$$\text{SM}_{rs} := \{j : 0.5 \leq \hat{b}_{jr}/\hat{b}_{js} \leq 2\}. \quad (2.39)$$

Then the sharing by sign between r and s is given by:

$$\frac{|SS_{rs} \cap (\text{QTL}_r \cup \text{QTL}_s)|}{|\text{QTL}_r \cup \text{QTL}_s|} \quad (2.40)$$

and sharing by magnitude between r and s is given by:

$$\frac{|SM_{rs} \cap (\text{QTL}_r \cup \text{QTL}_s)|}{|\text{QTL}_r \cup \text{QTL}_s|}. \quad (2.41)$$

ash analyses

For comparison with **mash** we also analyzed the GTEx data using the univariate shrinkage procedure **ash** [15]. We applied **ash** separately on each tissue using the same 20,000 randomly-selected gene-snp pairs as in the **mash** analysis. We then computed the posterior means and *lfsr* for the top SNPs.

mash-bmalite analyses

For comparison with **mash** we implemented a version on **mash-bmalite** ([5]) that outputs effect size estimates and *lfsr* values. This version of **mash-bmalite** can be thought of as a variation of **mash** but without the data driven covariance matrices, and with particular choices for the canonical covariance matrices, and with a smaller grid on ω than **mash** (consistent with the coarse grid used in [5]).

Specifically, the list of U_k for **mash - bmalite** include the 44 singleton configurations

($U_k = e_k e_k'$), and matrices corresponding to the models in [5] with heterogeneity parameters $H = \{0.0, 0.25, 0.5, 1\}$ [5]. (When heterogeneity=0, effects are equal in all conditions; when heterogeneity = 1, effects are independent among conditions.) We use a grid of $\omega \in \{0.1, 0.40, 1.6, 6.4, 25.6\}$ consistent with the coarse grid in [5] and designed to capture the range of the GTEEx Z -statistics.

Simulation Details

“Shared, Structured Effects”

We simulated \mathbf{b}_j from model (2.3) with equal weights on 8 different covariance matrices learned from the GTEEx data, but with the scaling factors ω simulated from a continuous distribution rather than using a fixed grid.

In detail:

1. Take the list of 8 “data-driven” covariance matrices learned from the GTEEx data (see Section 5.2), standardized to have maximum diagonal element 1 (Section 5.2).
2. Simulate 400 ‘true effects’: for each such effect j , a) choose U_j by selecting one of the eight U_k at random, all equally likely; b) simulate ω_j as the absolute value of an $N(0, 1)$ random variable; c) simulate $\mathbf{b}_j \sim N_{44}(\mathbf{0}, \omega_j U_j)$.
3. For 19,600 ‘null effects’ set $\mathbf{b}_j = \mathbf{0}$.

4. For all 20,000 effects, simulate $\hat{\mathbf{b}}_j \sim N(\mathbf{b}_j, V_j)$ where V_j is the diagonal matrix with diagonal elements 0.1^2 . Here, all standard errors are approximately 0.10, consistent with the GTEx dataset.

“Shared, Unstructured Effects”

In these simulations the 400 true effects were all independent and identically distributed: $\mathbf{b}_j \sim N_{44}(\mathbf{0}, 0.1^2 \mathbf{I}_R)$. Other details are as for Shared, Structured Effects.

“Independent Effects”

We also simulated data where effects were entirely independent across conditions; These were simulated as follows:

1. Independently for each $r = 1, \dots, 44$, choose a random set of $400 j \in \{1, \dots, 20,000\}$ to be the ‘true’ effects.
2. For the ‘true effects’ simulate $b_{jr} \sim N(0, \Sigma^2)$ where Σ^2 is chosen with equal probability from the set $\{0.1, 0.5, 0.75, 1\}$ to represent small and large effects within each condition. (All other effects are set to be 0).
3. Simulate $\hat{\mathbf{b}}_j \sim N(\mathbf{b}_j, V_j)$ as in other simulations.

Analysis of simulated data

Each simulated dataset $(\hat{\mathbf{b}}_j, V_j)$ was analyzed using `mash` as detailed above. In particular we re-estimated the $U_k, \boldsymbol{\pi}$ from the data, without making use of the true values for U . We estimated effects by their posterior mean (2.16) and assessed significance by the *lfsr* (2.18). Analyses using `ash` and `mash-bmalite` were performed similarly to the applications on the GTEx data (see above).

CHAPTER 3

APPLICATION OF MASH TO JOINTLY MAP CIS-EQTL IN MULTIPLE TISSUES

Improved effect size estimates

An important novelty of our method, `mash`, is its focus on *estimation of effect sizes*, in contrast with most existing multivariate analysis methods which focus only on *testing* for non-zero effects. Furthermore, `mash` is more than just an extension of existing methods to estimate effect sizes, because the underlying model (2.1) is more flexible than models underlying existing methods – and, indeed, includes existing models as special cases.

To illustrate the potential for multivariate analysis to improve accuracy of effect size estimates we performed simple simulations and compared three approaches to effect size estimation:

1. `mash`, the method we describe here.
2. A simpler version of our method, `mash-bmalite`, which represents an extension of existing methods to estimation of effect sizes. Specifically `mash-bmalite` performs effect size estimation based on the BMALite models from [5], which include both the random effects models (RE and RE2) and fixed effects model (FE) used in the software `metasoft` [14]. These models allow for shared effects

of equal size across all conditions (FE), shared effects of varying size across conditions (RE, RE2), and condition-specific effects (i.e. effects that occur in only one condition). Even though this, in itself, would represent a useful contribution, `mash` is more flexible than this. Specifically, `mash` can learn from the data that some subsets of conditions are more correlated than others, due to its use of data-driven covariance matrices in (2.1).

3. `ash` [15], which is a univariate analogue of `mash` designed to estimate effect sizes using results from a single condition. Results from `ash` are obtained by applying it separately to each condition, and so represent what can be achieved by a simple “condition-by-condition” analysis. This is included as a baseline against which to quantify the benefits of multivariate analysis.

We applied these three methods to estimate effect sizes under two scenarios:

1. “Shared, structured effects”: data were simulated using the model (2.1), based on the fit of this model to the GTEx eQTL data below (see Methods for details). In this scenario effects tend to be shared among many conditions, and furthermore these shared effects are highly “structured”, in that they are often similar in size (or at least sign), with the similarity being greater among some subsets of conditions than others. For example, in the GTEx analysis later we see that effect sizes are often particularly similar among the subset of brain-derived tissues. This scenario will arise frequently in practice, and an important goal of our work is to provide methods that perform well here.

2. “Shared, unstructured effects”: in this scenario effects are shared among all conditions (i.e. either every condition shows an effect, or no condition shows an effect), but the effect sizes and directions of the *non-zero* effects are independent across conditions. We aim to show that even in this unstructured setting `mash` provides improved effect estimates compared with an analogous univariate (condition-by-condition) approach, and in this case acts essentially as an extension of existing methods to estimate effect sizes.

In each case we simulate a 20,000 by 44 matrix of data \hat{B} containing 20,000 estimated effects in each of 44 conditions (and their associated standard errors). We assume that non-null effects are rare: of the 20,000 effects, only 400 are non-null. Thus the matrix of effects is sparse, with non-zero values concentrated in a small number of rows.

Figure 3.1a (See also Supplementary Table 1) compares the accuracy of effect size estimates, as measured by the relative root mean squared error (RRMSE) (2.22), which is the RMSE of the estimates, divided by the RMSE achieved by simply using the original observed estimates \hat{B} for the effects. Thus an $\text{RRMSE} < 1$ indicates that the method produces estimates that are more accurate than the original observations \hat{B} . As expected, the joint (multivariate) methods outperform the univariate method in both scenarios, due to their combining information across conditions. Furthermore, `mash` substantially outperforms the other methods in the “structured effects” scenario, and performs like `mash-bmalite` in the unstructured case. That is, the flexibility of `mash`, which is responsible for its improved performance in the structured

setting, does not decrease performance in this simpler setting.

In all settings, all three methods have $\text{RRMSE} < 1$, indicating a substantial improvement in accuracy compared with the original observed effects \hat{B} . This improvement can come from two sources: i) the methods shrink estimated effects towards zero, which improves average accuracy because most effects are indeed null; ii) in the presence of “structured effects”, the multivariate methods can share information across conditions to improve accuracy. For example, if a particular effect is shared, and similar in size, across a subset of conditions then averaging the observed effects in those conditions will improve estimation accuracy. Both these factors help explain the strong performance of `mash` in the structured effects setting (Supplementary Table 1).

As a check on implementation we also applied the three methods to data simulated under an “Independent effects” scenario, in which *all* effects are entirely independent across conditions, with no greater sharing than expected by chance. (Note that this is very different from the “shared, unstructured” scenario, where only the non-zero effects are independent.) We used this to confirm the intuition that in such settings the univariate method that analyzes each condition independently should perform best, as indeed it does (Supplementary Table 1).

Improved detection of significant effects

In addition to effect estimates, **mash** also provides a measure of significance for each effect in each condition. Specifically **mash** estimates the “local false sign rate” (*lfsr*) [15], which is the probability that the estimated effect has the incorrect sign. The *lfsr* is analogous to the local false discovery rate [23], but more stringent in that it insists that effects be correctly signed to be considered “true discoveries”. Similarly **mash-bmalite** can estimate the *lfsr*, but under its less flexible model; and **ash** can estimate the *lfsr* separately in each condition.

We used the same simulations as above to illustrate the gains in power to detect significant effects that come from the flexible multivariate model in **mash**. Figure 3.1b shows the trade-off between false positive and true positive discoveries for each method as the significance threshold is varied. The relative performance of the methods precisely mirrors the RRMSE results: multivariate methods perform best, and **mash** outperforms other methods for detecting shared structured effects. Further, in the “shared, structured” scenario **mash** is finding essentially all (> 99%) of the signals that the other methods find, plus additional signals (Supplementary Table 2). (And in the “shared, unstructured” scenario **mash** and **mash-bmalite** not only have similar average performance, but are finding almost identical signals; Supplementary Table 2.)

Comparison with `metasoft`

Among existing software packages for this problem, `metasoft` [14] is in some ways the most comparable with `mash`. In particular, it is both generic – requiring only effect estimates and their standard errors – and computationally tractible for $R = 44$. The `metasoft` software implements several different multivariate tests for association analyses, each corresponding to different multivariate models for the effects. For example, the FE model assumes that the effects in all conditions are equal; the RE2 model assumes that the effects are normally distributed about some common mean, with deviations from that mean being independent among conditions [26]; and the BE model is an extension of the RE2 model that allows that some effects are exactly zero [14]. These models are similar to the BMALite models from [5], and none of them capture the kinds of structured effects that can be learned from the data by `mash`. Our comparisons above illustrate the benefits of the more flexible model in `mash`. However, because differences in software implementation sometimes lead to unanticipated differences in performance we also performed some simple direct benchmarks comparing `mash` and `mash-bmalite` with `metasoft`.

Specifically we compared these methods in the simplest type of multivariate test: separating the null from the non-null signals, where here null means zero effect in *all* conditions. Here, for each model (FE, RE2, and BE), `metasoft` produces a p value for each multivariate test, whereas `mash` and `mash-bmalite` produce a Bayes Factor (see Methods); in each case these can be used to rank the significance of the tests. Figure 3.1c) shows the trade-off between false positive and true positive discoveries for

each method as the significance threshold is varied in the same simulation scenarios as above. In both cases `mash` is the most powerful method, again illustrating the benefits of its more flexible model.

Assessing heterogeneity and sharing in effects

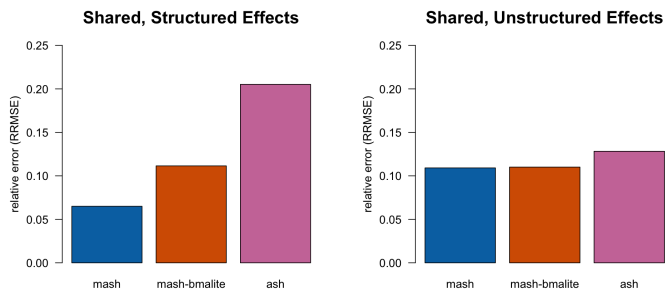
In analyses of effects in multiple conditions, it is often desired to identify effects that are shared across many conditions, or, conversely, those that are specific to one or a few conditions. This turns out to be a particularly delicate task. For example, [5] emphasize that the simplest approach – first identifying significant signals separately in each condition, and then examining the overlap of the significant effects – can very substantially under-estimate sharing. This is due to incomplete power: by chance, a shared effect can easily be significant in one condition and not in another. To address this [5,6] estimate sharing among conditions as a parameter in a joint hierarchical model, which takes account of incomplete power. However, these approaches are infeasible for $R = 44$. Furthermore, even for smaller values of R they have some drawbacks. In particular they are based on a “binary” notion of sharing, i.e. whether or not an effect is non-zero in each condition, and so do not capture differences in magnitude, or even signs, of effects among conditions. If effects that are shared among conditions actually differ greatly in magnitude – for example, being very strong in one condition and weak in all others – then this would seem important to know.

Here we address this problem with a new approach based on assessing *quantitative similarity of effects*. Specifically, we assess sharing of effects in two ways: i) “sharing by sign” (estimates have the same sign); and ii) “sharing by magnitude” (effects are similar in magnitude). Here we define similar in magnitude to mean both the same sign and within a factor of 2 of one another (although other thresholds could be used, and in some settings – for example, where the “conditions” are different phenotypes – the requirement that effects have the same sign may best be dropped.) These measures of sharing can be computed for any pair of conditions, and an overall summary of sharing across conditions can be obtained by assessing how many conditions share with some reference condition (here, we use the condition with the largest estimated effect as the reference).

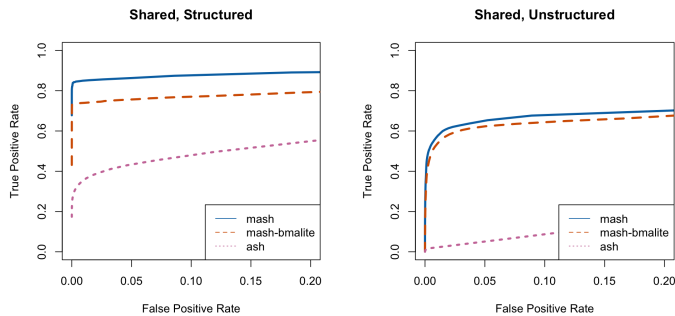
These measures of sharing could be naively estimated from the raw observed effect estimates from each condition; however, errors in these effect estimates will naturally lead to errors in assessed sharing. Because `mash` combines information across conditions to improve effect estimates (see above) it can provide more accurate estimates of sharing. To illustrate this we used the “shared structured” simulations to compare accuracy of overall estimates of sharing from `mash` with those from the raw effect estimates, as well as with `ash` and `mash-bmalite`. Table 3.1 summarizes these results, which confirm the improved accuracy of `mash`. For example, `mash` reduces the error in the estimated number of conditions sharing by sign from 4.7 for the raw estimates to 2.4 for `mash`.

Method	sign	magnitude
raw	4.7	4.3
ash	4.7	5.7
mash-bmalite	3.7	4.4
mash	2.4	2.3

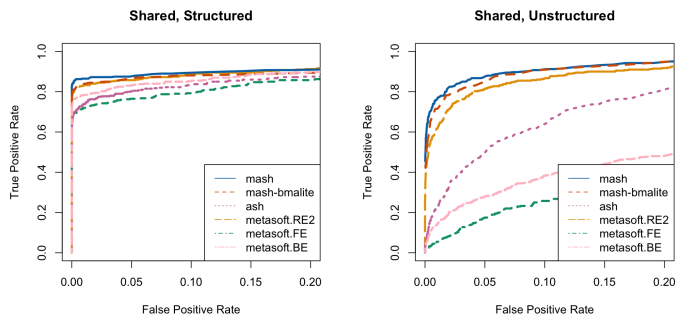
Table 3.1: **Errors in estimates of sharing for simulated data** For each method we computed the mean absolute error of the estimated number of conditions that share effects (by either sign or magnitude) with the condition with the largest estimated effects. Here “raw” indicates the performance of the raw estimates being input into each method. Each number is a mean error across the 400 non-null effects from the “shared structured effects” scenario.



(a) Accuracy of effect estimates (RRMSE).



(b) Detection of non-null effects in *each* condition (ROC curves).



(c) Detection of non-null effects in *any* condition (ROC curves).

Figure 3.1: **Comparison of methods on simulated data.** Results are shown for two simulation scenarios: “shared structured” effects, where the non-zero effects are shared among the 44 conditions in complex structured ways similar to patterns in the GTEx data; and “shared unstructured” effects, where the non-zero effects are shared among the 44 conditions with effect sizes that are independent among conditions. Panel (a) shows accuracy of effect estimates. Panels (b) and (c) show ROC curves for detecting significant effects. The primary difference between (b) and (c) is that in (b) each effect is treated as a separate discovery in each condition (which requires condition-specific measures of significance), whereas in (c) each effect is treated as a single discovery across all conditions (which requires only a single measure of significance, as in traditional meta-analyses). In (b) we require the estimated sign (+/-) of each significant effect to be correct to be considered a “true positive”. Our new method (**mash**) outperforms other methods, particularly for “shared structured” effects, a scenario expected to be common in genomics applications.

GTE_x cis-eQTL analysis

To illustrate the benefits and flexibility of `mash` in a substantive application we applied it to analyse expression Quantitative Trait Loci (eQTLs) across 44 human tissues/cell-types, using data from the Genotype Tissue Expression (GTEx) project [25]. The GTEx project aims to provide insights into the mechanisms of gene regulation by studying human gene expression and regulation in multiple tissues from health individuals. One fundamental question is which SNPs are eQTLs (i.e. associated with expression) in which tissues. Answering this could help distinguish regulatory regions and mechanisms that are specific to a few tissues vs shared among many tissues. It could also help with analyses that aim to integrate eQTL results with GWAS results to help identify the tissues that are most relevant to any specific complex disease (e.g. [25, 27]).

As input to `mash` we use a matrix of eQTL effect estimates \hat{b}_{jr} , and corresponding standard errors \hat{s}_{ij} , where the rows j index different SNP-gene pairs and the columns r index tissues (or cell types). We used the effect estimates and standard errors for candidate local (“cis”) eQTLs for each gene, distributed by the GTEx project (v6 release). These were obtained by (univariate) single-SNP analyses in each tissue by applying `MatrixEQTL` [28] on expression levels that have been corrected for population structure (using genotype principal components [29]) and for other confounding factors affecting expression data (both measured factors such as age and sex, and unmeasured factors using factor analysis [30]), and then rank-transformed to the corresponding quantiles of a standard normal distribution. Thus the effect size es-

estimates are in units of standard deviations on this transformed scale. Because, like most eQTL analyses, these estimates were obtained by single-SNP analysis, the estimated effects for each SNP actually reflect the effects of both the SNP itself and other SNPs in LD with it. Thus our analyses here do not distinguish causal eQTLs from SNPs that are in LD with the causal eQTLs; see Discussion.

We analysed the 16,069 genes for which univariate effect estimates were available for all 44 tissues we considered; the filtering criteria used ensure that these genes show at least some indication of expression in all 44 tissues.

Increased flexibility of mash improves model fit

Since the true effects are unknown we cannot compare models based on accuracy of effect estimates. Therefore, we instead illustrate the gains of the more flexible `mash` model using cross-validation: we fit each model to a random subset of the data (“training set”) and assessed model fit by its log-likelihood computed on the remaining data (“test set”). Comparing `mash` and `mash-bmalite` in this way we found that `mash` with a correlated residual framework improved the test set likelihood by 23,725 log-likelihood units, indicating a very substantial improvement in fit. Further, `mash` placed 79% of the mixture component weights on the data-driven covariance matrices, indicating that our methods for estimating these matrices are sufficiently effective that they capture most effects better than do the canonical matrices used by existing methods.

Identification of data-driven patterns of sharing

The increased flexibility of `mash` comes from its use of “data-driven” components to capture the main patterns of sharing (actually, covariance) of effects. This is illustrated in Figure 3.2, which shows the majority component that `mash` identifies in these data (relatively frequency 34%). The main patterns captured by this component are: i) effects are positively correlated among all tissues; ii) the brain tissues (and, to some extent, testis and pituitary) are particularly strongly correlated with one another, and less correlated with other tissues; iii) effects in whole blood tend to be somewhat less correlated with other tissues. Other components identified by `mash` are shown in Supplementary Figure 3.2. Some of these components also have positive correlations among all tissues and/or highlight heterogeneity between brain tissues and other tissues, confirming these as very common features in these data. However, other components also capture rarer patterns, such as effects that are appreciably stronger in one tissue than others (Supplementary Figure 3.5).

Patterns of sharing inform effect size estimates

Having estimated patterns of sharing from the data, `mash` exploits these patterns to improve effect estimates at each putative eQTL. Although we cannot directly demonstrate improved average accuracy of effect estimates in the real data (for this, see simulations above), individual examples can provide helpful intuition into the way that `mash` achieves improved accuracy. In this vein, Figure 3.3 shows three

Majority component in GTEx data
(relative frequency = 0.34)

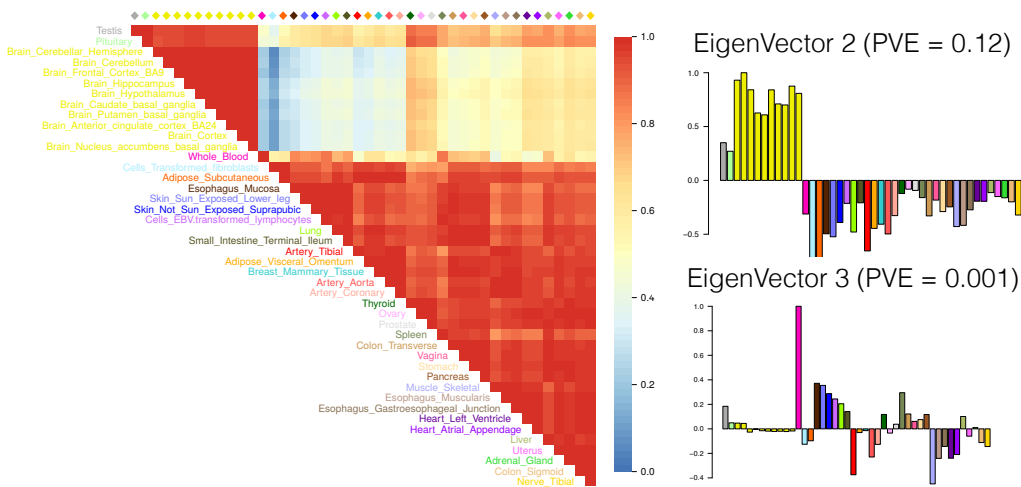


Figure 3.2: **Summary of primary patterns identified by mash in GTEx data.** Shown are the heatmap of the correlation matrix, and barplots of the first 3 eigenvectors, of the covariance matrix U_k corresponding to the dominant mixture component identified by mash. This component accounts for 34% of all weight in the GTEx data. In all cases, tissues are color-coded as indicated in the heatmap legend. The first eigenvector reflects broad sharing among tissues, with all effects in the same direction; the second eigenvector captures differences between brain (and, to a less extent, testis and pituitary) vs other tissues; the third eigenvector primarily captures effects that are stronger in whole blood than elsewhere.

illustrative examples, which we discuss in turn.

In the first example, the vast majority of effect estimates are positive in each tissue, with the strongest signals in a subset of brain tissues. Based on the patterns of sharing learned in the first step, **mash** estimates the effects in all tissues to be positive – even those with negative observed effects. This is because the few modest negative effects at this eQTL are outweighed by the strong background information that effects are highly correlated among tissues. Humans are notoriously bad at weighting background information against specific instances [31] – they tend to underweight background information when presented with specific data – so this behavior may or may not be intuitive to the reader. But **mash** performs this weighting using Bayes rule, which is ideally suited to this job. The **mash** effect estimates are also appreciably larger in brain tissues than in other tissues. Again, this is the result of using Bayes rule to combine the effect estimates for *this* eQTL with the background information on heterogeneity among brain and non-brain effects learned from *all* eQTLs.

In the second example, the effect estimates in non-brain tissues are mostly (30/34) positive, but modest in size, and only one effect is, individually, nominally significant ($p < 0.05$). However, combining information among tissues, **mash** effect estimates in non-brain tissues are all positive, and mostly “significant” ($lfsr < 0.05$). In contrast the data in brain tissues are inconsistent, with a mix of both positive and negative effect estimates. **mash** concludes that we cannot be confident of the eQTL effect sign in brain tissues. This example illustrates how **mash** can learn from the data how to group conditions, rather than treating them equally. In this case **mash** has learned

that effects in brain tissues are sometimes different from the other tissues, and hence avoids jumping to strong conclusions in the brain based on signal in other tissues.

In the final example, effect estimates vary in sign, and are modest except for a very strong signal in whole blood. While whole-blood-specific effects are estimated to be rare, **mash** (again, through Bayes theorem) recognizes that the strong data at this eQTL outweigh this background information, and estimates a strong effect in blood with insignificant effects in other tissues. This illustrates how **mash**, although focussed on combining information among tissues, can still recognize – and clarify – tissue-specific patterns when they exist.

Increased identification of significant effects

Our simulations demonstrated that the more flexible model behind **mash** can increase power to detect significant effects. To illustrate the effects of this here we compare the number of significant eQTLs detected by **mash** with those detected by our modified **mash-bmalite** and **ash**. To avoid double-counting of eQTLs in the same gene that are in LD with one another we assess the significance of only the “top SNP” in each gene, which we define to be the SNP with the largest (univariate) $|Z|$ -statistic across all tissues. Thus we focus on 16,069 putative eQTLs, each with effect estimates in 44 tissues, for a total of 707,036 effects.

The vast majority of top SNPs show a very strong signal in at least one tissue (97% have a maximum $|Z|$ score exceeding 4), consistent with most of these genes

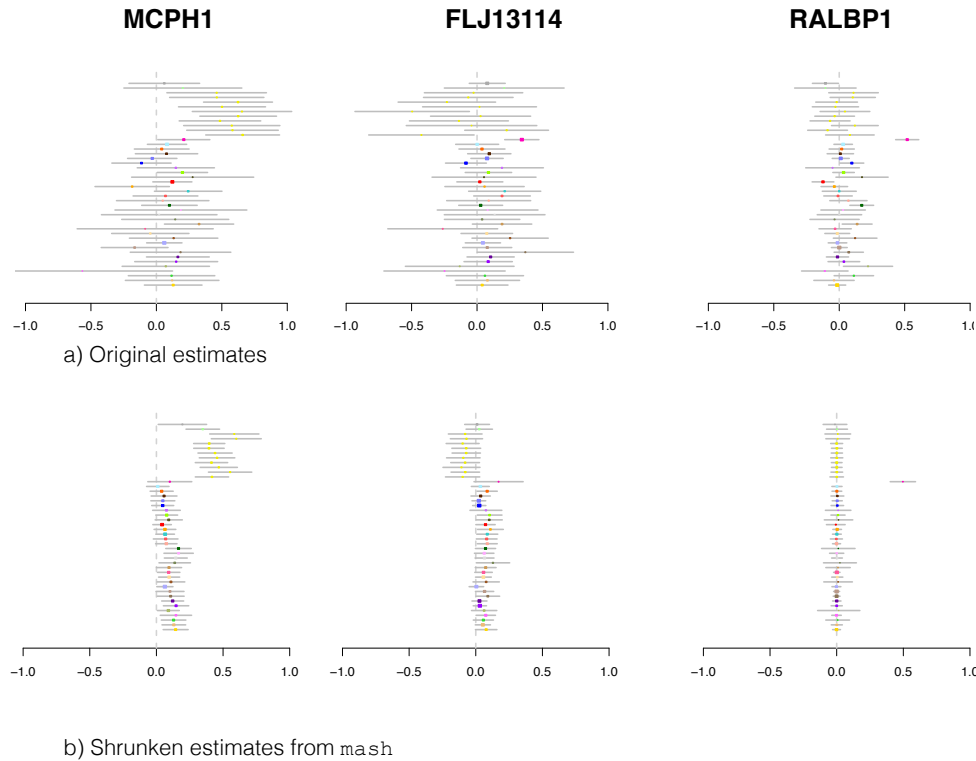


Figure 3.3: **Examples illustrating of how `mash` uses patterns of sharing to inform effect estimates in the GTEx data.** In panel a) each colored dot shows the raw effect estimate for a single tissue (color-coded as in Figure 3.2), with grey bar indicating ± 2 standard errors. These are the data input into `mash`. Panel b) shows the corresponding estimates output by `mash` (posterior mean, ± 2 posterior standard deviations). In each case `mash` combines information across all tissues, using the background information – patterns of sharing – it has learned from data on all eQTLs, to produce more precise estimates. Together, these three examples illustrate the flexibility of `mash` in combining information across different subsets of tissues for different eQTLs, depending on how their data match different patterns of sharing identified in the overall data. See main text for detailed discussion.

containing at least one eQTL in at least one tissue. However, the univariate tissue-by-tissue analysis (**ash**) identifies only 13% of these effects as “significant” at $lfsr < 0.05$; that is, the univariate analysis is highly confident in the sign of the effect in only 13% of cases. In comparison **mash-bmalite** identifies 39% as significant at the same threshold, and **mash** identifies 47%. As in the simulations, the significant associations identified by **mash** include the vast majority (96%) of those found significant by either of the other methods (Supplementary Table 3). Thus, the multivariate methods identify the most significant effects, with **mash** identifying the most.

Overall, **mash** found 76% (12,189/16,069) of the top SNPs to be significant in at least one tissue. We refer to these as the “top eQTLs” in subsequent sections.

Sharing of effects among tissues

To investigate sharing and heterogeneity of the top eQTLs among tissues we used the quantitative measures of sharing introduced above: sharing of effects by sign and by magnitude. The results are summarized in Table 3.2 and Figure 3.4. Because a major feature of these data is that brain tissues generally show more similar effects than non-brain tissues we also show results separately for these subsets of tissues. The results confirm extensive eQTL sharing among tissues, particularly among the brain tissues. Sharing in sign exceeds 85% in all cases, and is as high as 98% among the brain tissues. (Furthermore, these numbers may underestimate the sharing in sign of actual causal effects, because of the potential effects of multiple eQTLs per

gene in LD; see Supplementary Text.) Sharing in magnitude is inevitably lower, because sharing in magnitude implies sharing in sign. Overall, on average 37% of tissues show an effect within a factor of 2 of the strongest effect at each top eQTL. However, within brain tissues this number increases to 78%. That is, not only do eQTLs tend to be shared among the brain tissues, but the effect sizes tend to be quite homogeneous. Because these results are based on only the top eQTLs at each gene they reflect patterns of sharing among strong cis eQTLs; it is possible that weaker eQTLs may show different patterns of heterogeneity among tissues.

Of course, some tissues share eQTLs more than others. Figure 3.5 summarizes eQTL sharing by magnitude between all pairs of tissues (see Supplementary Figure 3.4 for sharing by sign). In addition to strong sharing among brain tissues, **mash** also identifies increased sharing among other biologically-related groups, including: arteries (tibial, coronary and aortal), two groups of gut tissues (one group containing esophagus and sigmoid colon; the other containing stomach, terminal ilium of the small intestine and transverse colon), skin (sun-exposed and non-exposed), adipose (Subcutaneous and Visceral-Omentum) and heart (left ventricle and atrial appendage). This figure also reveals that the main source of heterogeneity in effect sizes among brain tissues is in cerebellum vs non-cerebellum tissues, and also emphasizes sharing between the pituitary and brain tissues.

Different levels of effect sharing among tissues means that effect estimates in some tissues gain more precision than others from the joint analysis. To quantify this we computed an “effective sample size” (ESS) for each tissue that reflects the typical pre-

cision of its effect estimates (Supplementary Figure 3.1). The ESS values are smallest for tissue that show more “tissue-specific” behaviour (e.g. testis, whole blood; see below), and are largest for coronary artery, reflecting its stronger correlation with other tissues.

Tissue-specific eQTLs

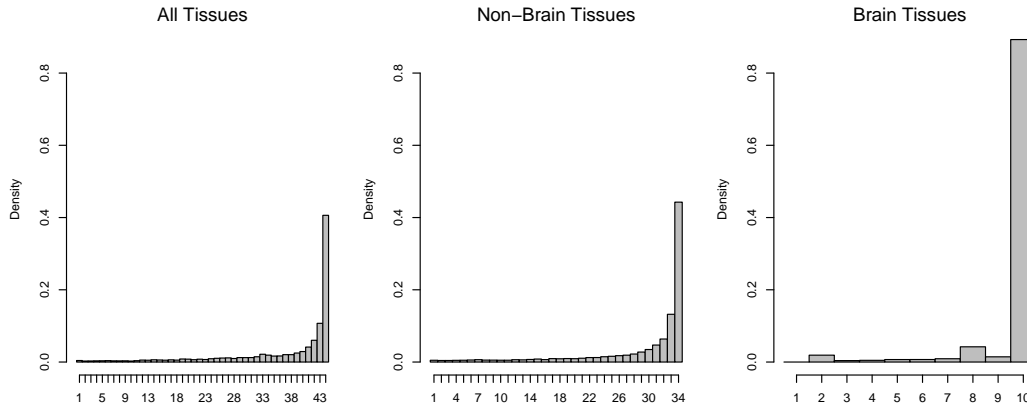
Despite high average levels of sharing of eQTLs among tissues, **mash** also identifies eQTLs that are relatively “tissue-specific”. Indeed, the distribution of the number of tissues in which an eQTL is shared by magnitude has a mode at 1 (Figure 3.4), representing a subset of eQTLs that have much stronger effect in one tissue than in any other (henceforth “tissue-specific” for brevity). Breaking down this group by tissue (Supplementary figure 3.5) identifies Testis as the tissue with the most tissue-specific effects. Testis also stands out, with whole blood, as having lower pairwise sharing of eQTLs with other tissues (Figure 3.5). Other tissues showing stronger-than-average tissue specificity (in either Supplementary Figure 3.5 or Figure 3.5) include skeletal muscle, thyroid, and transformed cell lines (fibroblasts and LCLs).

One possible explanation for tissue-specific eQTLs is tissue-specific expression. That is, if a gene is strongly expressed only in one tissue this could explain why an eQTL for that gene might show a strong effect only in that tissue. Whether or not a tissue-specific eQTL is due to tissue-specific expression could considerably impact biological interpretation. Thus we assessed whether tissue-specific eQTLs identified

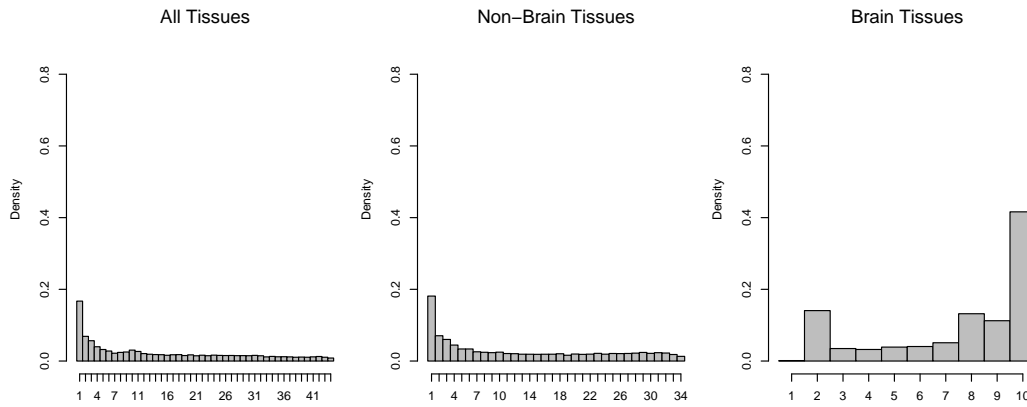
here could be explained by tissue-specific expression. Specifically, we took genes with tissue-specific eQTLs, and examined the distribution of expression in the eQTL-affected tissue relative to expression in other tissues. We found this distribution to be similar to genes without tissue-specific eQTLs (Supplementary Figure 3.6). Thus most tissue-specific eQTLs identified here are not simply reflecting tissue-specific expression.

Data	All Tissues	Non-Brain	Brain
Shared by Sign ($\tilde{b} > 0$)	0.85	0.85 (0.88)	0.96 (0.98)
Shared by Magnitude: ($\tilde{b} > 0.5$)	0.36	0.40 (0.44)	0.76 (0.85)

Table 3.2: **Summary of sharing among top eQTLs.** Numbers show the proportion of effects meeting a given sharing criterion. “Shared by sign” requires that the effect has the same sign as the strongest effect among tissues. “Shared by Magnitude” requires that the effect is also within a factor of 2 of the strongest effect. Numbers in parentheses are obtained by a secondary `mash` analysis of subsets of tissues.



(a) Number of Tissues Shared By Sign.



(b) Number of Tissues shared by Magnitude

Figure 3.4: Histogram showing estimated number of tissues in which top eQTLs are “shared” by two different definitions, a) sign and b) magnitude. Sharing by sign means that eQTLs have the same sign of effect; Sharing by magnitude means that they also have similar effect size (within a factor of 2). Left: All tissues; Center: non-brain tissues; Right: brain tissues.

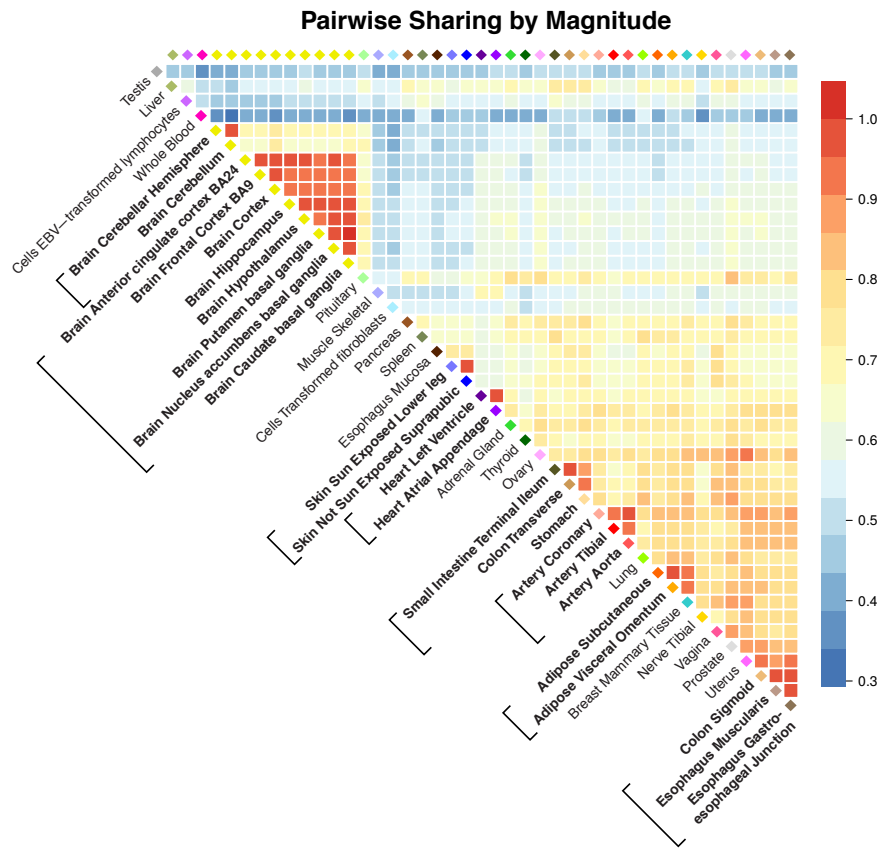


Figure 3.5: **Pairwise sharing by magnitude of eQTL among tissues.** For each pair of tissues we consider the top eQTLs that are significant in at least one of the two tissues, and plot the proportion of these that are “shared in magnitude” – that is, have effect estimates that are the same sign and within a factor of 2 of one another. Pink triangles highlight groups of biologically-related tissues mentioned in the text as showing particularly high levels of sharing.

Supporting Information

Here, I include information supplemental to the main result, but nonetheless informative to and understanding of eQTL sharing.

Effects of Linkage Disequilibrium

Linkage Disequilibrium (LD) between SNPs has two distinct effects. First, LD causes correlations in the observations of effects for near-by SNPs in the same gene. This issue is likely minor here. Although, when estimating g , `mash` ignores correlations between rows of \hat{B} , this can be justified as a “composite likelihood” approach [19], and composite likelihood methods tend to perform well at point estimation.

Second, effect estimates we obtain for each SNP from single-SNP analysis are not actually the individual causal effects of that SNP; rather they are the *combined effects of all SNPs that are in LD with that SNP*, weighted by their LD [32], [33]. This issue is more important, because of the likely presence of multiple eQTLs in some or many genes. It also applies to all single-SNP eQTL analyses, which is the vast majority of all published eQTL analyses, and not just `mash`. Ideally one would develop multi-SNP multi-tissue methods for association analysis at each gene to avoid this issue. And indeed, we see `mash` as a first step towards this more ambitious goal. However, for now we limit ourselves to highlighting one specific feature of our results that we believe may be a consequence of the use of single-SNP effect estimates, and that may

change in multi-SNP analyses that better account for LD.

Specifically, LD among multiple causal SNPs can cause single-SNP analyses to identify eQTL that appear to have strong effects of opposite sign in different tissues. One example is shown in Supplementary Figure 3.3: this eQTL has strong positive Z scores in brain tissues, and negative Z scores in most other tissues, initially suggesting that this eQTL might have causal effects in opposite directions in brain vs non-brain tissues. However, the Z scores could also have a different explanation: there could be two eQTLs in LD with one another, one of which (A say) has a strong effect in brain tissues, and the other of which (B say) has a strong effect in other tissues. If the expression-increasing-allele at A is in negative LD with the expression-increasing-allele at B then the single SNP Z scores for either SNP will show opposite signs in brain vs non-brain. Indeed, closer examination of the data at this gene suggests that this explanation is likely correct in this case (Supplementary Figure 3.3). A similar example is discussed in [25] (their Supplementary Figure S14).

For this reason we believe that estimates of sharing in sign given above are likely to be underestimates of the sharing in sign of actual causal effects, and we caution against over-interpreting eQTLs that show significant effects of different signs in different tissues.

Increase in effective sample size due to multivariate analysis

A particular emphasis of our work here is improved quantitative estimates of effect sizes in each condition. When estimating effects in a condition, `mash` uses the data not only from that condition but also from other “similar” conditions. In this way `mash` effectively increases the sample size available, and this improves both accuracy and precision of estimates. The improvement will be strongest for conditions that are similar to many other conditions, and weaker for conditions with more “condition-specific” effects.

To illustrate this effect in the GTEx data we compute an “effective sample size” (ESS) for each tissue based on the standard deviations of the `mash` estimates. The ESSs (Supplementary Figure 3.1) vary from 241 for testis to 1926 for coronary artery. Other tissues with relatively smaller ESS include liver, pancreas, spleen and brain cerebellum. Identifying tissues with smaller ESS could help guide prioritization of (effectively) under-represented tissues in future experimental efforts.

For testis the ESS of 241 represents only a small (1.4-fold) increase compared with actual sample size, reflecting that its effects are more “tissue specific”, or, more precisely, that they are less correlated with other tissues. Other tissues showing a similarly small gain in ESS include transformed fibroblasts and whole blood, which are also highlighted as showing more “tissue specific” signals above. In contrast, the ESS for coronary artery represents a 14-fold increase compared with the actual sample size for this tissue, reflecting its stronger correlation with other tissues. On

average, across all tissues, `mash` provides a 6-fold increase in ESS for estimating these (strongest) eQTL effects, reflecting the overall moderate to large correlation among effect sizes across tissues.

One caveat here is that ESS reflects *average* gains in precision for a tissue: in practice effects that are shared across many tissues will benefit more than effects that are tissue-specific. For example, if one were particularly interested in effects that are specific to uterus (which has the smallest actual sample size here), then the substantial ESS for uterus may not be as useful as it would first seem. More generally, detecting tissue-specific effects will inevitably benefit most from collecting more samples in that particular tissue.

Supplementary Tables

Method	Simulation Framework	RRMSE ^{All}	RRMSE ^{Non-null}	RRMSE ^{Null}
mash	Shared, structured	0.06	0.44	0.015
mash-bmalite	Shared, structured	0.11	0.78	0.018
ash	Shared, structured	0.21	1.34	0.076
mash	Shared, unstructured	0.14	1.00	0.014
mash-bmalite	Shared, unstructured	0.15	1.03	0.014
ash	Shared, unstructured	0.21	1.37	0.078
mash	Independent	0.28	1.82	0.112
mash-bmalite	Independent	0.28	1.82	0.118
ash	Independent	0.21	1.37	0.076

Table 3.3: **Supplemental Table 1: Comparison of accuracy of effect size estimates for each method.** Results show the RRMSE for all effects (RRMSE^{all}), and for the subsets of effects that are truly non-null ($\beta \neq 0$; RRMSE^{Non-null}) and truly null ($\beta = 0$, RRMSE^{Null}). Values of RRMSE^{Null} < 1 indicate how shrinkage towards zero is helping improve the estimates of null effects. Values of RRMSE^{Non-null} < 1 indicate how pooling information across conditions can improve accuracy of estimates of non-null effects. (In the Independent simulations the shrinkage of all methods improves overall performance, despite hurting performance for the non-null effects, because most effects are null.)

Method	Simulation Framework		
	Shared, structured	Shared, unstructured	Independent
mash only	3889	622	32
ash only	0	0	740
mash-bmalite only	37	9	79
mash & ash	7	0	44
mash & mash-bmalite	5777	336	70
ash & mash-bmalite	0	0	10
ALL	3477	2	5962

Table 3.4: **Supplemental Table 2: Comparison of Identified Associations in Simulations.** Results show the overlap of associations identified by each combination of methods. In both “Shared” scenarios **mash** captures the vast majority of the associations identified by the other methods.

Method	Associations Identified
mash only	63956
ash only	2383
mash-bmalite only	11789
mash & ash	665
mash & mash-bmalite	176572
ash & mash-bmalite	248
ALL	88459

Table 3.5: **Supplemental Table 3: Comparison of Identified Associations in GTEx Data.** Results show the overlap of identified associations with each method. **mash** captures the vast majority of the associations identified by the other methods, in addition to other associations.

Table 3.6: Qualitative comparison of some available software packages for analyzing effects in multiple conditions.

	flexibility ^a	tractable? ^b	effect estimates? ^c	adaptive? ^d	generic? ^e
mash	highest	yes	yes	yes	yes
eQTL-BMA ¹	high	no	no	yes	somewhat
eQTL-BMA (lite)	moderate	yes	no	yes	somewhat
metasoft ² (BE)	moderate	somewhat ⁵	no	no	yes
metasoft (RE2)	moderate	yes	no	no	yes
metasoft (RE)	low	yes	no	no	yes
metasoft (FE)	low	yes	no	no	yes
corMotif ³	moderate-high ⁴	yes	no	yes	no

^a A qualitative assessment, based on what types of effect sharing is captured by the underlying model.

^b Does computation cost scale non-exponentially with R ? Provides indication of whether software should be tractable for moderate R , $R > 15$ say.

^c Does the software output effect estimates (other than univariate estimates)?

^d Does the software learn hyperparameters from data (yes) or are they fixed (no)?

^e Can the software be directly applied to summary data?

¹ eQTL-BMA implements methods in [5]

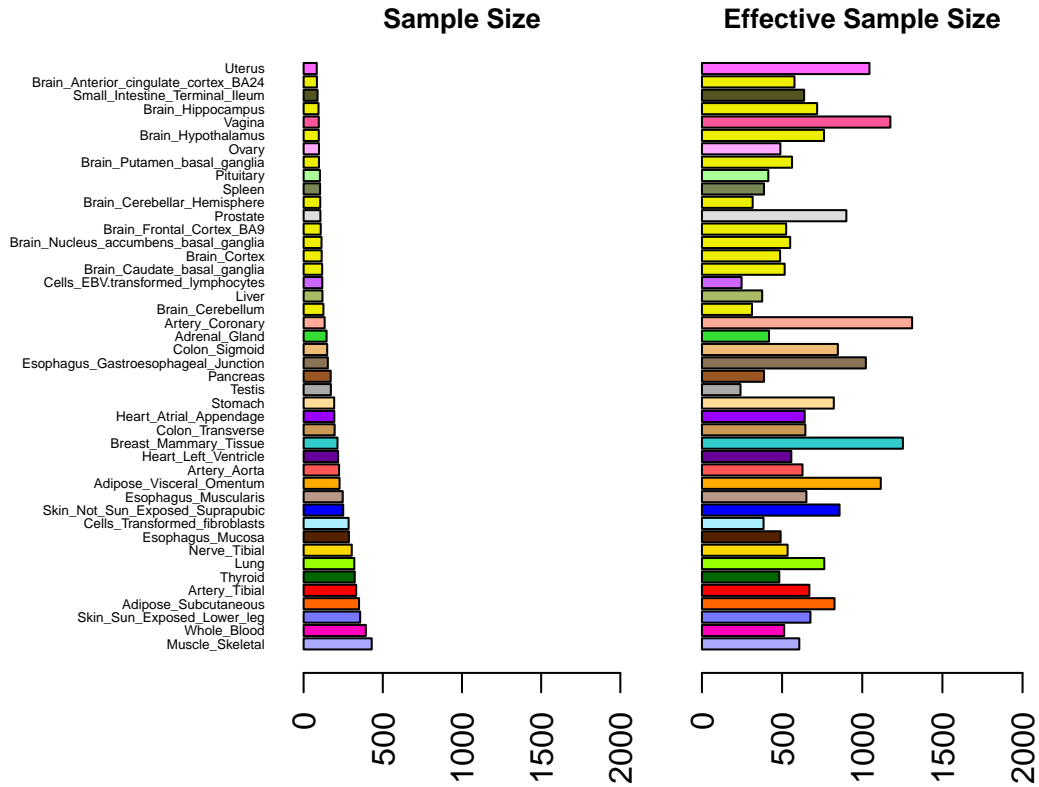
² metasoft implements methods in [9, 14]

³ corMotif implements methods in [12]

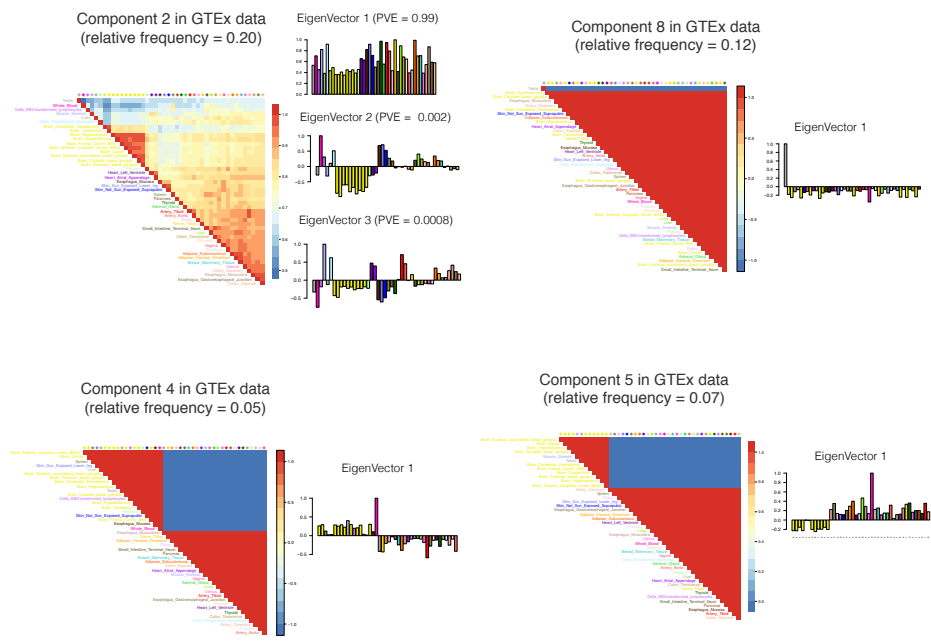
⁴ corMotif is qualitatively different than other methods considered here. It allows for sharing of effects in arbitrary subsets of conditions, but not for positive correlation in the effects sizes.

⁵ The model scales exponentially, but computation is done using a sampling-based method whose tractability is difficult to predict.

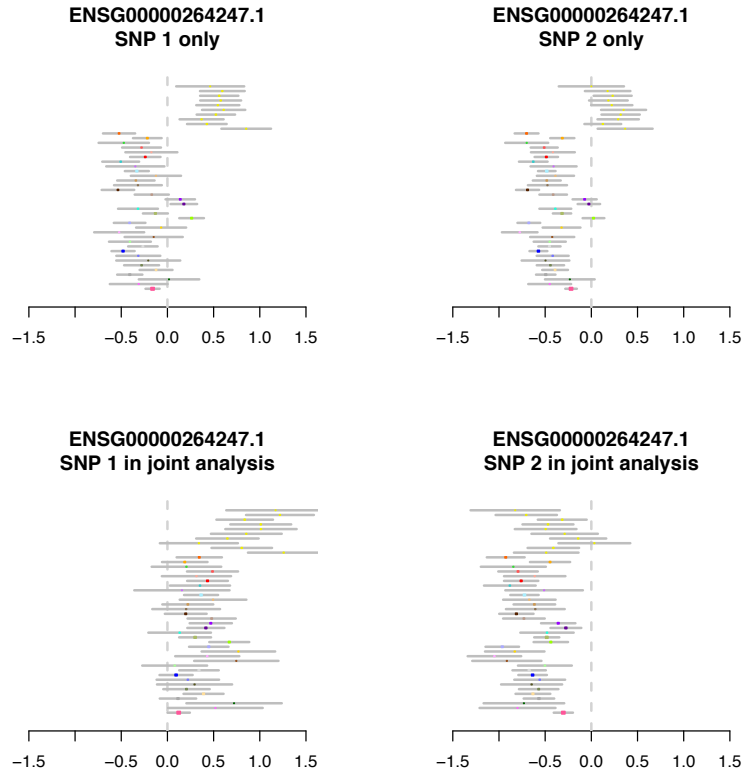
Supplementary Figures



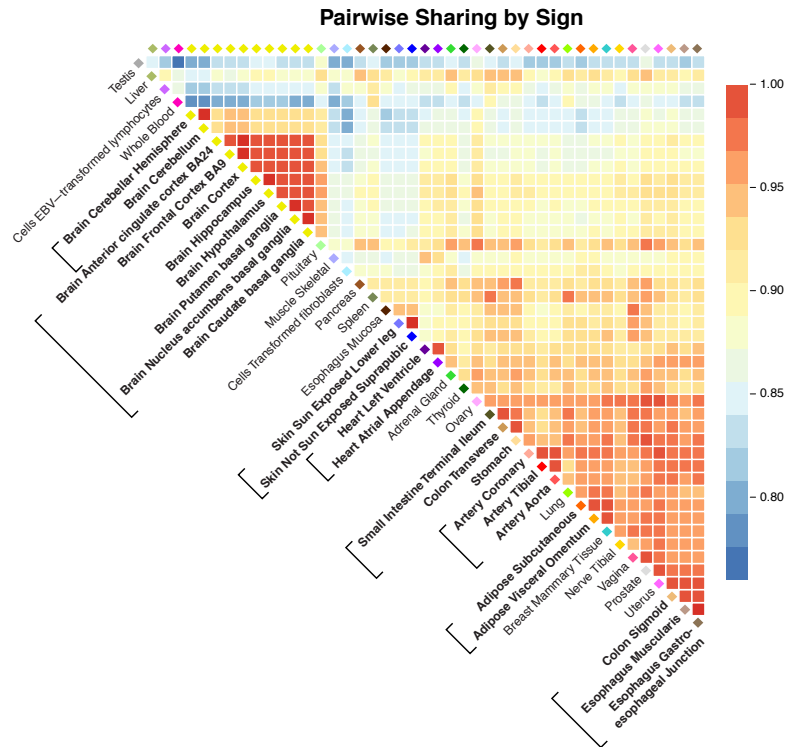
Supplementary Figure 3.1: **Sample sizes and effective sample sizes from mash analysis across tissues.** Left: sample size for each tissue; Right: median effective sample size for each tissue. Tissues are ordered by their original sample size. Effective sample sizes are consistently higher than actual sample sizes, primarily due to sharing of information among tissues.



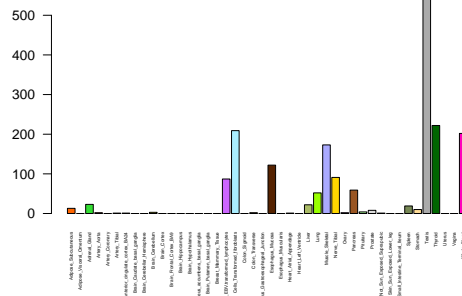
Supplementary Figure 3.2: **Summary of covariance matrices U_k with largest estimated weight (> 1%) in GTEx data.** Component 2 largely captures qualitatively similar effects to the component highlighted in Figure 3.2, although with quantitative differences. Component 8 captures testis-specific effects. Components 4 and 5 primarily capture effects that are stronger in Whole Blood than other tissues.



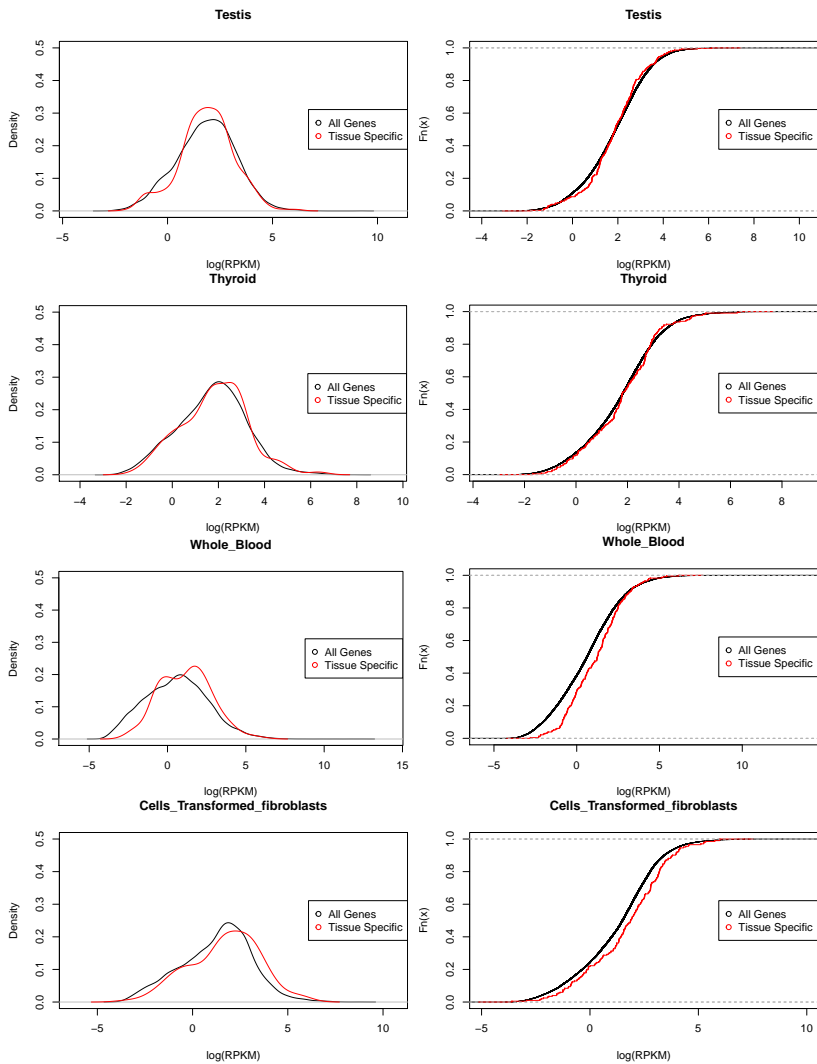
Supplementary Figure 3.3: **Illustration of how Linkage Disequilibrium can impact effect estimates.** This gene was chosen as an example where the effect estimates in the “top eQTL” were opposite in sign in brain vs non-brain tissues, and where further investigation suggested that this is likely due to multiple eQTLs in LD. Specifically, SNP1 and SNP2 are the SNPs that show the strongest eQTL association in brain and non-brain tissues respectively. The top panels show effect estimates for these SNPs from a simple (1-SNP) regression model in each tissue, $Y = \mu + \hat{B}_i g_i$ where $i \in \{1, 2\}$ indexes the two SNPs. The bottom panels show effects from a multiple (2-SNP) regression model in each tissue, $Y = \mu + \hat{B}_1 g_1 + \hat{B}_2 g_2$. The simple regression estimates show apparent opposite-sign effects in brain vs non-brain tissues (with testis and pituitary clustering with brain in one case). However, the multiple regression results suggest that in fact there are (at least) two eQTLs in this gene, because both SNPs show a significant effect that excludes 0 in most tissues. Furthermore, for both SNP1 and SNP2 the multiple regression effect estimates are consistent in sign across all tissues.



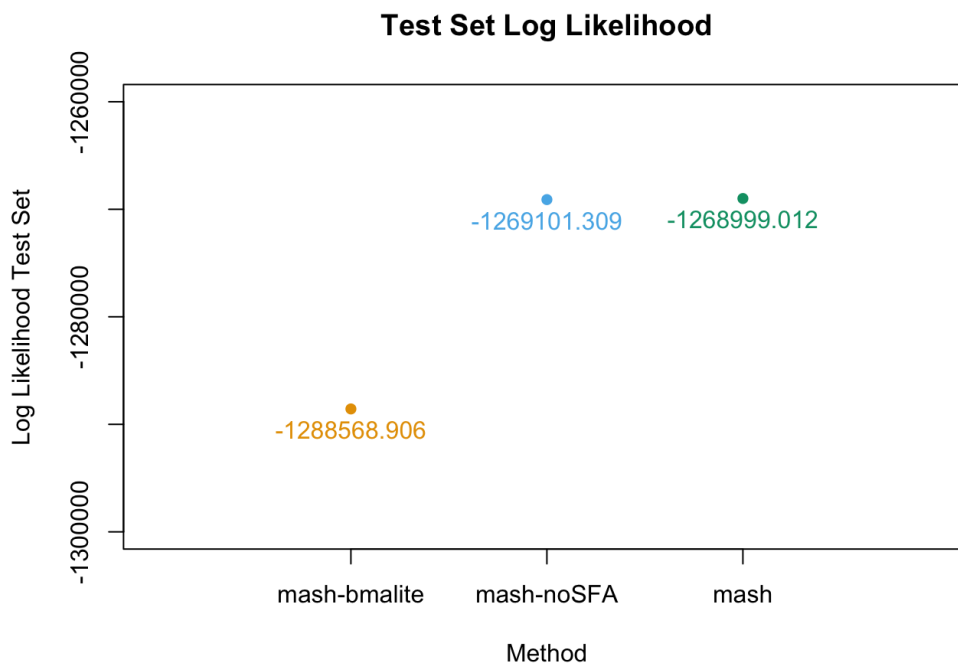
Supplementary Figure 3.4: **Pairwise sharing by sign.** For each pair of tissues we consider the top eQTLs that are significant in at least one of the tissues, and estimate the proportion that have effect sizes that are the same sign. These proportions are displayed in this heatmap.



Supplementary Figure 3.5: **Number of “tissue-specific eQTLs” in each tissue.** Here “tissue-specific” is defined to mean that the effect is at least 2-fold larger in one tissue than in any other (i.e. $\tilde{b}_{jr} > 0.5$ in only one tissue).



Supplementary Figure 3.6: **Expression levels in genes with “tissue-specific eQTLs” are similar to those in other genes.** The plots compare the densities (left) and cumulative distribution functions (right) of the expression level for all genes (black) and for genes identified as having a “tissue-specific” eQTL (red) in each of Testis, Thyroid, Whole Blood and Transformed Fibroblasts. In each case the distribution functions are reasonably similar, demonstrating that tissue-specific eQTLs are not simply reflecting tissue-specific expression. Expression is here defined as median across individuals of the log Reads per Kilobase Mapped (RPKM).



Supplementary Figure 3.7: **Increase in log-likelihood on Test Set as new U_k are added.** The figure shows the log-likelihood on the test set for different “models” (choices of U_k). From left to right the models are: **mash-bmalite** (no data-driven U_k); **mash-no-SFA** (the combination of canonical and data-driven covariances, excluding the rank-1 matrices derived from SFA); **mash** (the full combination of canonical and data-driven covariances described here). The result illustrates how, as more data-driven covariances are added, the log-likelihood on the test set typically increases. (Note that the point for **mash** is approximately 100 log-likelihood units higher than **mash-no-SFA**, although it is difficult to see the difference on this scale.)

CHAPTER 4

MASHCOMMONBASELINE: COMPARING CONDITIONS TO A COMMON CONTROL

We chose to apply this method for flexibly estimating effects across multiple tissues to a setting in which there existed no obvious ‘effect’ in each condition, but we might want to estimate the change in some quantity computed in multiple conditions over a common control for example. Such analyses are common in practice in case and control literature ([34]) and in gene expression studies, ([35, 36]). Scientists are often interested in identifying patterns of differential gene regulation under many conditions, compared to a condition in which there exists no external stimuli. In these cases, differential expression in any condition is defined as a difference in expression over a common control. For instance, ([1]) was interested in identifying gene expression changes in human innate immune cells that are specific to infection with mycobacteria, where the change in gene expression was defined relative to an uninfected control. In these cases, such a pattern might implicate genes specific to defense against some infections but not others. Thus infection similarity may be defined by similar elicited response, and genes may share defense roles in some conditions but not others. As in `mash`, typical approaches to solving this problem ([37]; [38]) focusing on analyzing differential expression over control in each condition separately, and thus vastly underestimate sharing. While we wish to take advantage of correlation and thus boost power as in `mash`, we now must consider the additional burden of comparing all subsequent conditions to the same reference. Specifically,

comparing all condition estimates to the same base condition without accounting for the correlation in errors artificially induces many false positives, and thus any successful joint method needs to account for this correlation in the error structure.

We now consider the problem of estimating changes that occur in multiple conditions relative to some common baseline. Such applications are extremely common in genomics. Consider, for instance, the case of gene expression, in which we observe an \mathbf{R} vector of ‘noisy’ estimates of gene expression in R subgroups. For every subgroup, there is some underlying mean gene expression which may be different from a common control group. For example, Blischak et al ([1]) estimated changes in gene expression in cells infected with 8 different strains of Mycotuberculin bacterium, and compared these patterns in gene expression with a common uninfected sample. Simply estimating the expression by subtracting the mean expression in the control condition from every subsequent condition might create false positives, because of the inherent correlation in errors. However, analyzing each condition separately fails to exploit the power of a joint analysis. Furthermore, exploring effect sizes consider the quantitative heterogeneity of effect sizes among groups.

Our goal is to fit a model for the distribution of a $R - 1$ dimensional quantity of deviations using a set of J observational data points $\hat{\boldsymbol{\delta}}_j$ that measure the observed deviations over the measured control expression..

Here, the use of bold-face notation indicates a vector, while matrix quantities are typeset in capital but not boldface letters.

4.1 Defining the Model

Now, we observe for each gene \mathbf{j} a vector of uncentered noisy mean feature expression $\hat{\mathbf{c}}_j$ across R conditions:

$$\hat{\mathbf{c}}_j | \mathbf{c}_j \sim N_R(\mathbf{c}_j, \hat{V}). \quad (4.1)$$

where \mathbf{c}_j represents the ‘true’ means across R subgroups.

Thus \tilde{L} denote the $(R-1) \times R$ matrix of contrasts which removes the true mean ‘control’ expression from each subsequent condition. Here, we let $R = 8$:

$$\tilde{L} = \begin{bmatrix} -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ -1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ -1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ -1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Then,

$$\boldsymbol{\delta}_j = \tilde{L}\mathbf{c}_j. \quad (4.2)$$

Here, we can express $\boldsymbol{\delta}_j$ as a zero-centered mixture of multivariate normals where each covariance matrix U_k now represents the underlying covariance matrix from which the ‘true’ deviations $\boldsymbol{\delta}_j$ are thought to arise.

$$\boldsymbol{\delta}_j | \boldsymbol{\pi}, \mathbf{U} \sim \sum_{\mathbf{k}, \mathbf{l}} \pi_{\mathbf{k}, \mathbf{l}} N_{\mathbf{R}}(\mathbf{0}, \omega_{\mathbf{l}} \mathbf{U}_{\mathbf{k}}). \quad (4.3)$$

Then we can write the observed noisy deviations $\hat{\boldsymbol{\delta}}_j$ from baseline (control) expression as

$$\begin{aligned} \tilde{L} \hat{\mathbf{c}}_j &= \tilde{L} \mathbf{c}_j + \tilde{L} \mathbf{E} \\ \hat{\boldsymbol{\delta}}_j &= \boldsymbol{\delta}_j + \mathbf{E}^*. \end{aligned} \quad (4.4)$$

where $\mathbf{E}^* \sim \mathcal{N}(0, \tilde{L} \hat{V} \tilde{L}')$. Critically, even if \hat{V} is diagonal and thus the observed noisy mean expression measurements in each condition are independent, $\tilde{L} \hat{V} \tilde{L}'$ and thus accounts for the induced correlation in errors.

Critically, our quantity of interest now, $\boldsymbol{\delta}_j$ represents the true ‘deviations’ from mean gene expression which are the effects of interest across each condition and can be seen as the effects in `mash`.

For intuition, we present a simulated matrix \hat{C} of R length independently simulated mean expression measurements $\hat{\mathbf{c}}_j$ across 8 subgroups (see Figure 4.1 for under-

standing). We display the correlation heatmap of this matrix to illustrate that such measurements are independent among subgroups. After contrasting by the 7×8 matrix \tilde{L} , we then demonstrate (4.1) how the matrix of observed deviations derived from the original independent measurements $\hat{\Delta}$ is inherently correlated because the same noisy estimate of the common control is subtracted from every subsequent column.

4.2 Likelihood with `mashcommonbaseline`

Combining 4.3 and 4.4 for each gene J at each component k , integrating over δ_j ,

$$\hat{\delta}_j \sim \mathcal{N}(0, U_k + \tilde{L}\hat{V}\tilde{L}'). \quad (4.5)$$

Using our original `mash` framework, we assume that the observed 'control' contrasted estimates directly approximate δ_j , that is $\hat{\delta}_j = \delta_j + \mathbf{E}$ where $\mathbf{E} \sim N(0, V)$, assuming that the variance of the residuals V simply represents the sum of the variances of the uncontrasted measurements.

Hence `mashcommonbaseline` allows that the residuals are correlated, and we hope controls for false positives - i.e., erroneously identifying associations simply because the corrected errors are correlated.

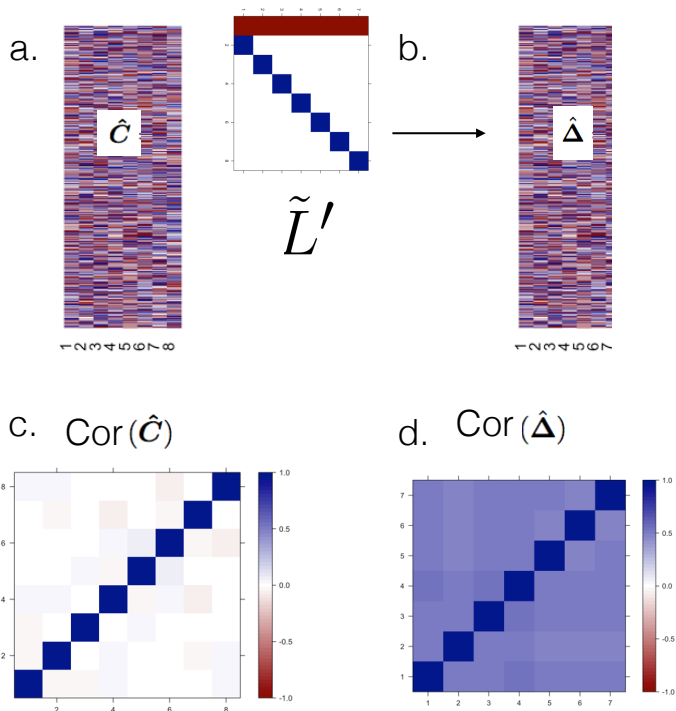


Figure 4.1: **Intuition behind mashcommonbaseline.** We illustrate the necessity of our approach by first simulating a matrix in (a) \hat{C} in which each row is simulated independently, i.e., $\hat{c}_j \sim N(0, I_8)$. It's corresponding correlation matrix (c) is diagonal. After removing the control condition gene expression measurement from every subsequent condition by premultiplying by the 7x8 contrast matrix \tilde{L} , we then observe the matrix of estimate deviations $\hat{\Delta}$ shown in (b). The correlation matrix of this matrix (d) has non-zero diagonal elements, indicating correlation between the elements of $\hat{\delta}_{j..}$.

4.3 mashcommonbaseline Simulation

Suppose that rather than comparing to an overall baseline for the gene, we are interested in producing a method that accurately captures true ‘contrasts’, that is, deviation from a control. Such an application is often desired in biology - i.e., measuring gene expression or some dependent variable on treated versus untreated controls, and detecting whether such differences are significant. We seek first to demonstrate that failing to account for the inherent correlations induced by subtracting the same noisy observed control measurement from each subsequent condition inflates our identification of true deviations even when no underlying deviations exist. We want to show that if you ignore these correlations you might get false positives.

In the simplest case, suppose no true deviation exists, but we approximate contrasts as simply the subtraction from observed controls. Thus for a given gene j , we need to simulate only one scalar μ which represents identical mean expression in all R conditions.

$$\mathbf{c}_j = \mu \mathbf{1}. \tag{4.6}$$

Thus $\mathbf{c}_{j2\dots R}$ is identically \mathbf{c}_{j1} in all R conditions (i.e. $\delta_{jr} = 0 \forall \mathbf{r}$).

Then the observed mean gene expression measurements are:

$$\hat{\mathbf{c}}_j \sim N(\mathbf{c}_j, 1/2I_R). \quad (4.7)$$

Then the observed deviations from the control $\tilde{L}\hat{\mathbf{c}}_j = \hat{\boldsymbol{\delta}}_j$:

$$\hat{\boldsymbol{\delta}}_j = N(0, \frac{1}{2}\tilde{L}\tilde{L}'). \quad (4.8)$$

However, one might naively assume the subtracted observations are simulated $N(0, I_{R-1})$ and thus merely the sum of the variances of the original observations $\hat{\mathbf{c}}_j$:

$$\begin{aligned} \mathbb{V}(\hat{\mathbf{c}}_{j2} - \hat{\mathbf{c}}_{j1}) &= \mathbb{V}(\hat{\mathbf{c}}_{j2}) + \mathbb{V}(\hat{\mathbf{c}}_{j1}) \\ \mathbb{V}(\hat{\boldsymbol{\delta}}_j) &= I_{R-1} \end{aligned} \quad (4.9)$$

We expect that such an assumption will produce many false positives because such a model fails to adequately capture the correlation in residuals and attributes all correlation to ‘true’ effects.

4.4 Simulation Design

For every gene J , the j^{th} row of the matrix \hat{C} let $\hat{\mathbf{c}}_{j1}, \hat{\mathbf{c}}_{j2} \dots \hat{\mathbf{c}}_{jR}$ be simulated $N(\mathbf{c}_j, 1/2I_R)$. A naive analysis, and indeed the assumption behind available methods ([37]; [39]) might assume that the observed differences over the control,

$$\hat{\boldsymbol{\delta}}_j = \hat{\mathbf{c}}_{j2-1}\hat{\mathbf{c}}_{j3-1}\dots\hat{\mathbf{c}}_{jR-1}, \quad (4.10)$$

are simulated $N(0, I_{R-1})$.

Our method, `mashcommonbaseline` allows that $\hat{\boldsymbol{\delta}}_j \sim N(0, \frac{1}{2}\tilde{L}\tilde{L}')$. First, we compare to the number of false positives computed when `mash` is applied to a matrix containing vectors of truly independent random variables, $\mathbf{z} \sim N(0, I_7)$ and show that the assumption $\hat{\boldsymbol{\delta}}_j \sim N(0, I_7)$ indeed produces many more false positives than analysing data in which the observed deviations are truly independent.

Then, we show that incorporating the correlated residuals, i.e., correctly modeling $\tilde{L}\hat{\mathbf{c}}_j$ as $\sim N(0, \frac{1}{2}\tilde{L}\tilde{L}')$ reduces the number of false positives by over 80%.

4.5 Simulation with Signal

In the simulation above, we assume that the true deviation from baseline $\boldsymbol{\delta}_j$ is 0. Now, we add signal to a number of ‘non-null’ simulations such that for any subgroup $r = 2 \dots R$ an effect different than the control exists:

$$\mathbf{c}_{j2 \dots R} = \mathbf{c}_{j1} + \boldsymbol{\delta}_j. \quad (4.11)$$

We simulate such that for 10% of the associations, the non-null effects $\boldsymbol{\delta}_j$ arise from

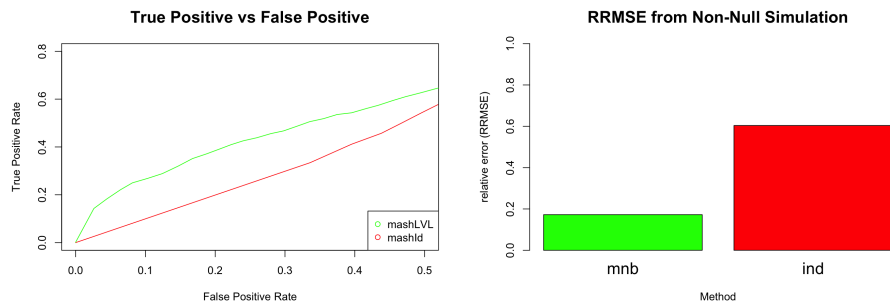


Figure 4.2: **Improvement in both Power and Accuracy.** The number of false positives versus true positive associations detected using `mashcommonbaseline` (accounting for residual correlation among errors) improves. We also demonstrate improved accuracy of the estimated deviations in the (RRMSE) as in equation 2.22.

one of three components, i.e., $K=3$, and $\omega = 1$ according to our model 4.3.

We show that the relative root mean square error (RRMSE) and the ROC curve both improve using `mashcommonbaseline` (4.2).

4.6 Application of Deconvolution

As in `mash`, we choose the ‘top’ deviations by which to initiate our covariance matrices for denoising.

Then we can apply Extreme Deconvolution ([17]) to the maximum estimates of $\hat{\delta}_j$; that is, those $\hat{\delta}_j$ that have a maximum observed deviation $\hat{\delta}_{maxj} \geq 2.5$, where $\hat{\delta}_{max}$ is defined as :

$$\hat{\delta}_{maxj} := \max_r \frac{\hat{\delta}_{maxjr}}{\hat{s}_{jr}}. \quad (4.12)$$

We define the set of genes T which have a $\hat{\delta}_{maxj} \geq 2.5$:

$$T := \{j : \hat{\delta}_{maxj} \geq 2.5\}. \quad (4.13)$$

As before we use the EM algorithm from [17] to obtain the U_k that maximize the likelihood over all j in T :

$$\begin{aligned} L(\theta) &:= p(\hat{\delta}|\pi, V, U_k) \\ &= \prod_{j \in T}^J p(\hat{\delta}_j|\pi, V, U_k) \\ &= \prod_{j \in T}^J \sum_k^K \pi_k N_R(\hat{\delta}_j; \mathbf{0}, U_k + \tilde{L}V_j\tilde{L}'). \end{aligned} \quad (4.14)$$

We will then rescale each of the U_k by choosing a set of ω that are appropriate to $\hat{\Delta}$ (see section (:grid)) to comprise a set of $P = KxL$ prior covariance matrices Σ that maximizes the likelihood over all J , and not just those with a maximum deviation.

$$\begin{aligned}
L(\pi) &:= p(\hat{\boldsymbol{\delta}}_j | \pi, V, \Sigma) \\
&= \prod_{j=1}^J p(\hat{\boldsymbol{\delta}}_{jj} | \pi, V, \Sigma) \\
&= \prod_{j=1}^J \sum_p \pi_p N_R(\hat{\boldsymbol{\delta}}_{jj}; \mathbf{0}, \Sigma_p + \tilde{L}V_j\tilde{L}').
\end{aligned}$$

(4.15)

We assemble a matrix of likelihoods that will compute the following likelihood at each of the P components:

$$\begin{aligned}
\hat{\boldsymbol{\delta}}_j &\sim \mathcal{N}(0, \Sigma_p + \tilde{L}\hat{V}_j\tilde{L}') \\
T_{jp} &= \Sigma_p + \tilde{L}\hat{V}_j\tilde{L}' \\
\hat{\boldsymbol{\delta}}_j &\sim \mathcal{N}(0, T_{jp})
\end{aligned} \tag{4.16}$$

the posterior means (equation 2.16) and posterior covariance (equation 2.17) are then computed as described.

4.7 Real Data Analysis: Method Application

To illustrate the power of our approach, we applied it to data in which gene expression across multiple subgroups had indeed been compared to a control. Blischak *et al* [1] analyzed gene expression in cells infected with 8 strains of Tuberculin bacteria and sought to understand changes in gene expression in comparison to uninfected controls.

We consider the gene expression patterns identified 18 hours post-infection in our analysis. The gene expression readings (see [1] for details) represent batch-corrected \log_2 counts per million for the 12,728 Ensemble genes analyzed in this study for each of the 156 samples gene expression data across the 8 conditions and controls. To obtain summary statistics for each gene-condition pair, we used the Empirical Bayes linear model method Limma ([40]), to estimate $\hat{\mathbf{c}}_j$ of mean gene expression and corresponding standard errors. In the `mashcommonbaseline` analysis, the $[j, r]$ entry of the matrix $\hat{\mathbf{C}}$ represents the mean gene expression of gene j for individuals infected with bacteria r .

Let the matrix of observed deviations, $\hat{\Delta} = \hat{\mathbf{C}}\tilde{\mathbf{L}}'$. We use the rows of this matrix corresponding to the 1000 genes containing the top absolute $\hat{\delta}_{max}$ to obtain T as in equation (4.6).

As in equation (4.14), we use these maximum statistics to estimate the 9 data-driven covariance matrices and use the larger 12,728 gene data set and to expand by a grid which ranged from $\omega = 1 \times 10^{-7}$ to $\omega = 487$, with $m=2$ as in section 2.2.

We use all 12,728 genes to fit the EM algorithm and thus estimate π and compute posteriors for all 12,728 genes.

4.8 Results: Patterns of Sharing

We first examine the primary patterns of sharing identified by our analysis. We see that the majority of the hierarchical weight falls on two patterns which reflect broad sharing of both sign and magnitude across conditions. However, we can see that in both patterns, *Yersinia* and *Salmonella* seem to share effects more closely.

Similar to the eQTL analysis above, we now address this problem of describing sharing with a new approach based on assessing *quantitative similarity of effects*. Specifically, we assess sharing of effects in two ways: i) “sharing by sign” (estimates have the same sign); and ii) “sharing by magnitude” (effects are similar in magnitude). Here we define similar in magnitude to mean both the same sign and within a factor of 2 of one another (although other thresholds could be used, and in some settings – for example, where the “conditions” are different phenotypes – the requirement that effects have the same sign may best be dropped.) These measures of sharing can be computed for any pair of conditions, and an overall summary of sharing across conditions can be obtained by assessing how many conditions share with some reference condition (here, we use the condition with the largest estimated effect as the reference).

We can see (4.4) that similar to the eQTL analysis, deviations tend to be shared

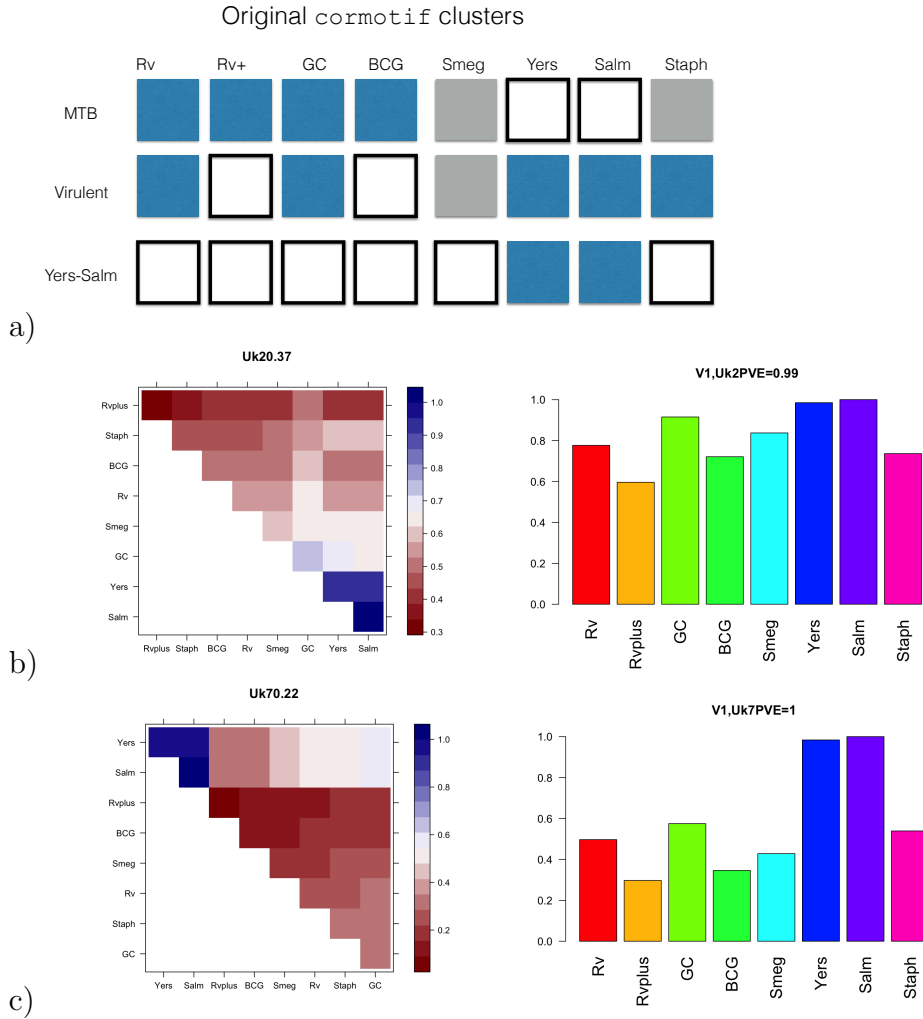
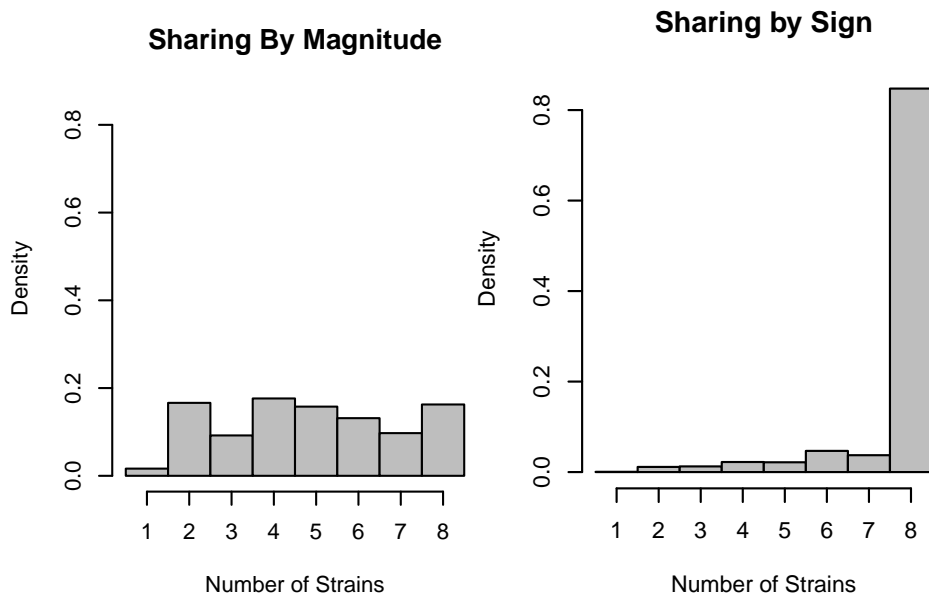


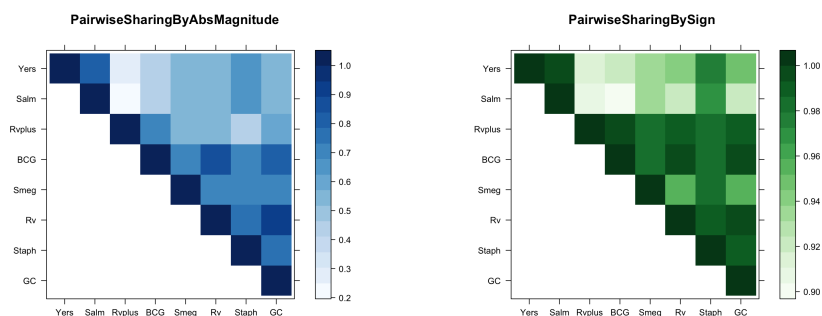
Figure 4.3: **Inferred Patterns of Sharing in data from [1]:** In the top panel (a), *cormotif* identifies three separate clusters that are collapsed into the primary patterns of sharing in *mashcommonbaseline*. This figure is adapted from ([1]). We display the two most common patterns of sharing (b) and (c), U_k as defined by the proportion of hierarchical weight received. At left is a heatmap of each scaled covariance matrix (i.e., $U_k/\max(\text{diag}(U_k))$) and at right it's first eigenvector. The two patterns of sharing which receive the majority of the weight both represent deviations which are broadly shared in both sign among subgroups, though the deviations appear strongest in *Yersinia* and *Salmonella*, with weaker correlations among *BCG* and *RV+*.

by sign much more commonly than they are shared by magnitude, reflecting the fact that to be shared by magnitude must be shared by sign. However, we can also observe that contrary to the eQTL analysis, there is a range of intermediate sharing by magnitude. That is, genes can be of the same sign (and thus deemed differentially expressed over controls in the same direction) and yet demonstrate more subtle patterns of quantitative heterogeneity.

Noting the sizable proportion of genes which share either 2 or 4 conditions by magnitude and the majority of genes sharing by sign (Figure 4.4) we can then examine which subgroups tend to be shared more closely than others. Again, we observe that *Yersinia* and *Salmonella* tend to share effects more closely than other subgroups, and we notice a subpopulation of sharing among BCG and RV+, which is biologically reassuring given that both are attenuated vaccine strains. Lastly, we note that the more virulent Staph and GC strains tend to share closely by sign. All these patterns were identified by the `corMotif` analysis conducted by Blischak et al, [1]. Rather than restrict genes to a pattern that simply related the binary outcome of differential expression in each subgroup, we describe a pattern of continuous effects; e.g., differential expression in the same direction but of varying magnitude among conditions. This is consistent with the primary eigenvectors of the predominant patterns of sharing (Figure 4.3) which suggest that a gene can be modestly differentially expressed in GC and Staph and strongly differentially expressed in *Yersinia* and *Salmonella* for example, with corresponding effect sizes reflecting this information.



(a) Distribution of Sharing by Magnitude and Sign



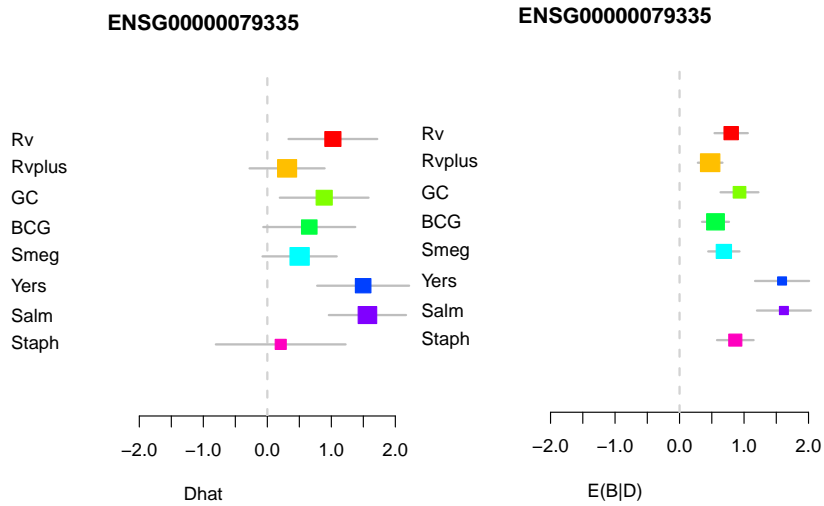
(b) Patterns of Sharing Heatmaps

Figure 4.4: **Visualizing Quantitative Sharing using mashcommonbaseline:** In panel a and b, we display the distribution of sharing by magnitude (left) and sign (right) in the tuberculin data. In (b) we display pairwise sharing by magnitude (left) and sign (right) using `mashcommonbaseline` as computed in 3.5 and 3.4 For each pair of conditions, we consider the effects that are significant in at least one of the tissues, and estimate the proportion that have effect sizes that are within 2-fold of one another and the same sign (a) or of the same sign only (b)

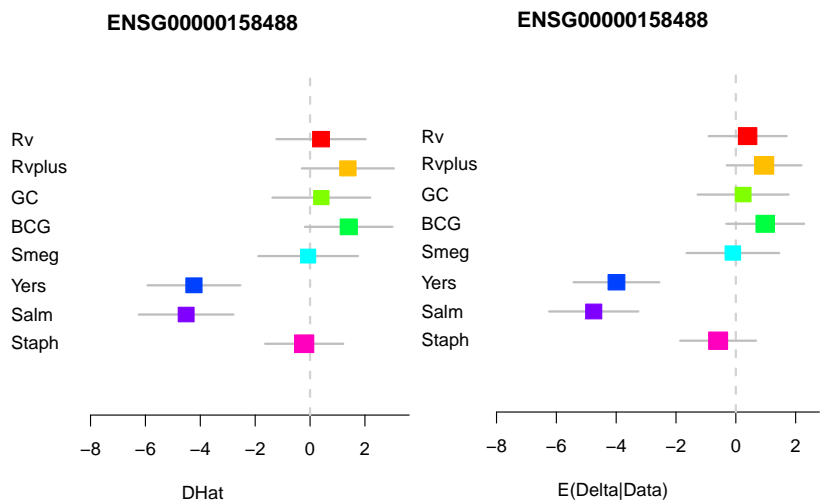
4.9 Improvements over Existing Methods

One of the benefits of estimating effects rather than simply assigning genes to a particular group or pattern as in existing methods (see [39], [5]) is that intermediate levels of expression between subgroups can be captured. Both [39], [5] assign genes to a particular configuration, with configuration being represented as an arrangement of binary ‘on-off’ designation among subgroups. For example, in the data above, a common pattern identified in the analysis by Blischak et al [1] (representing roughly 9% of genes) was “Yersinia-Salmonella” in which genes were classified as having strong differential expression over controls in only these two bacterially-infected conditions(see 4.5 . However, our analysis identified more subtle patterns (Fig. 4.3) as evidenced by the eigenvectors of the patterns which receive the most weight. In fact, in figure 4.5 we examine examples in which `corMotif` assigned the data to a Yersinia-Salmonella pattern while the raw data (and accompanying `mashcommonbaseline` posterior effects) showed deviations with strong evidence in Yersinia and Salmonella, and weaker but positively correlated evidence of differential expression in the additional subgroups.

In Figure 4.5, we examine two different examples of `mashcommonbaseline`’s ability to adaptively capture patterns of deviation. In the first situation, `corMotif` assigned the gene to be differentially expressed only in Yersinia and Salmonella, while the real data show non-zero estimated deviations in most tissues. Indeed, our posterior analysis borrows information across conditions to recognize that in fact, the gene is differentially expressed in all conditions, though the posterior estimated effects



(a) All Effects Shared



(b) Subgroup Specific

Figure 4.5: Examples illustrating how `mashcommonbaseline` can capture more subtle patterns of sharing than restricting effects to a binary outcome in conditions. Similar to Figure 3.3, we plot the original noisy estimates (± 2 SD) at left and the posterior means (± 2 Posterior SD) at right in each panel. In the top panel, `mashcommonbaseline` recognizes differential expression in all groups, while in the bottom panel, `mashcommonbaseline` can still preserve subgroup specificity when it exists.

are larger in *Yersinia* and *Salmonella* than in other conditions. This is a benefit of reporting quantitative continuous effects, rather than binary outcomes of differential expression.

In the second example, we illustrate that though `mashcommonbaseline`, like `mash`, aims to emphasize sharing by incorporating information across subgroups, it can still recognize ‘subgroup’ specificity when it exists. Here, both `mashcommonbaseline` and `corMotif` show significant effects in only *Yersinia* and *Salmonella*, and `mashcommonbaseline` successfully shrinks the insignificant effects to zero.

CHAPTER 5

APPLICATION OF MASH TO GWAS OF MULTIPLE PHENOTYPES

One of the challenges of modern genetics has been the paucity of genetic variants identified to explain any sizable proportion of the variation in disease risk for complex diseases ([41]; [42]). Given the availability of summary statistics for a variety of univariate analyses, it is only natural to consider a joint analysis to augment the power of such an approach. The assumption here is that diseases may share underlying genetic underpinnings, and thus any information about the association of a genetic variant with one disease may be augmented by incorporating information about the allele's effect on an additional disease. While these diseases may sometimes seem clinically related, it is possible that diseases with seemingly unrelated clinical features might still borrow genetic information. To explore such an approach, we considered analyzing an assortment of 16 different conditions (5.1) initially classified as Psychiatric (anx,bip,mdd,scz,an,as), morphologic (bmi, height), Immune (cd,uc,ibd) and Endocrinological (fg,fn,fa,ls,aam).

5.1 Consortium

We aggregated data from over 800,000 SNPs across 16 traits including schizophrenia, bipolar and major depressive disorder, anxiety, height, anorexia-nervosa, autism

spectrum disorder, fore arm bone-density, femoral neck bone-density, lumbar spine bone-density, age-at-menarche, body mass index, Crohn’s disease, irritable bowel disease, Ulcerative-Colitis and Fasting-Glucose available on LD Hub ([43]).

Consortium	Abbreviation	Disease
pgc	scz	schizophrenia
pgc	bip	bipolar
pgc	mdd	major depressive disorder
angst	anx	anxiety
giant	height	height
pgc	an	anorexia nervosa
pgc	as	autism spectrum disorder
gefoc	fa	fore arm bone density
gefoc	fn	femoral neck bone density
gefoc	ls	lumbar spine bone density
reprogen	aam	age at menarche
giant	bmi	body mass index
ibdgenetics	cd	crohns disease
ibdgenetics	ibd	irritable bowel disease
ibdgenetics	uc	ulcerative colitis
magic	fg	fasting glucose

Table 5.1: **Studies Considered** We describe the consortium and the disease referenced by each abbreviation.

We analyzed the data as described below.

5.2 Data Analysis Procedure

Generate data-driven covariance matrices U_k for Consortium

As in the `mash` application to the GTEx data above, we first identify rows j of the matrix \hat{B} that likely have an effect in at least one condition. Here, j indexes the SNP and r indexes the study. There are $R = 16$ studies.

In the GWAS data, we chose the rows corresponding to the “top” 1000 SNPs, which we define to be the SNPs with the highest value of Z_j^{\max} where

$$Z_j^{\max} := \max_r \hat{\mathbf{b}}_{jr} / \hat{s}_{jr}. \quad (5.1)$$

(We used max here, rather than, say, the sum, to try to include effects that are very strong in a single condition and not only effects that are shared among conditions.)

Next we fit a mixture of MVN distributions to these strongest effects, using methods from [17]. Specifically results in [17] provide an EM algorithm for fitting a model very similar to (2.3) and (2.2) with the crucial difference that there is no scaling parameters on the covariances. That is,

$$p(\mathbf{b}_j | \boldsymbol{\pi}, \mathbf{U}) = \sum_k \pi_k N_R(\mathbf{b}_j; \mathbf{0}, U_k). \quad (5.2)$$

As in our `mash` application to the eQTL data above, we again initialized the EM

with $K = 3$ and

- $\tilde{U}_1 = \frac{1}{J} \tilde{Z}' \tilde{Z}$, the empirical covariance matrix of \tilde{Z} .
- $\tilde{U}_2 = \frac{1}{J} \sum_{p=1}^P \lambda_p^2 v_p v_p'$, which is a rank P approximation of the covariance matrix of \tilde{Z} . Here we used $P = 3$.
- $\tilde{U}_3 = \frac{1}{J} (LF)'(LF)$ which is a rank Q approximation of the covariance matrix of \tilde{Z} .

In addition to the covariance matrices obtained from this EM algorithm, we added some more matrices based on the SFA results, specifically

- The 5 matrices $F'_q L'_q L_q F'_q$, which are each rank 1 matrices that reflect the effects captured by the q th factor in the SFA analysis ($q = 1, \dots, 5$).

In total this procedure produces 8 data-driven covariance matrices for our GWAS analyses which we append to a set of canonical matrices (5.2).

Generate canonical covariance matrices U_k

To these “data-driven” covariance matrices we add the following “canonical” matrices:

1. The matrix \mathbf{I}_R . This represents the situation where the effects in different conditions are independent, which may be unlikely in some applications (like

the GTEEx application here), but seems useful to include if only to exclude it.

2. The R rank-1 matrices $\mathbf{e}_r \mathbf{e}_r'$ where \mathbf{e}_r denotes the unit vector with 0s everywhere except for element r which is a 1. These represents effects that occur only in a single condition.
3. The rank-1 matrix $\mathbf{1}\mathbf{1}'$ where $\mathbf{1}$ denotes the R -vector of 1s. That is, the matrix of all 1s. This represents effects that are identical among all conditions.

The user can, if desired, add additional canonical matrices. For example, if R is moderate then one could consider adding the 2^R canonical matrices that correspond to shared (equal) effects in each of the 2^R subsets of conditions.

In total this procedure produces 18 canonical covariance matrices for our GWAS analyses.

Standardize covariance matrices

Since (2.3) uses the same grid of scaling factors ω we standardize the matrices U_k obtained above so that they are similar in scale. Specifically, for each k , we divide every element of U_k by the maximum diagonal element of U_k (so that the maximum diagonal element of the rescaled matrix is one). These rescaled matrices provide the \hat{U} , completing step i)-a of `mash`.

Define grid of ω_l values

We choose a dense grid of ω_l ranging from “very small” to “very large”. [15] provides a specific way to select suitable limits $(\omega_{\min}, \omega_{\max})$ for this grid in the univariate case; we simply apply this method to each condition r in turn and take the smallest ω_{\min} and the largest of the ω_{\max} as the grid limits. The internal points of the grid are then obtained as in the univariate case [15], by setting $\omega_l = \omega_{\max}/m^{l-1}$, for $l = 1, \dots, L$, where $m > 1$ is a user-tunable parameter that affects the grid density and L is chosen to be just large enough so that $\omega_L < \omega_{\min}$. In the GWAS data we use $m = 2$.

Estimate $\boldsymbol{\pi}$ by maximum likelihood

Given $\hat{\boldsymbol{U}}, \boldsymbol{\omega}$, we estimate the mixture proportions $\boldsymbol{\pi}$ by maximum likelihood.

To simplify notation, let $\Sigma_{k,l} := \omega_l \hat{U}_k$, and replace the double index k, l with a single index p which ranges from 1 to $P := KL$. Thus the prior (2.3) becomes:

$$p(\mathbf{b}_j | \boldsymbol{\pi}, \boldsymbol{\Sigma}) = \sum_p \pi_p N_R(\mathbf{b}_j; \mathbf{0}, \Sigma_p). \quad (5.3)$$

Combining the prior (5.3) with the likelihood (2.2), we have that each row of \hat{B}

comes from a mixture of MVNs:

$$p(\hat{\mathbf{b}}_j | \boldsymbol{\pi}, V, \boldsymbol{\Sigma}) = \sum_p^P \pi_p N_R(\hat{\mathbf{b}}_j; \mathbf{0}, \Sigma_p + V_j). \quad (5.4)$$

This essentially comes from the fact that the sum of two MVNs is MVN.

Assuming independence of rows of \hat{B} , the likelihood is given by

$$\begin{aligned} L(\boldsymbol{\pi}) &:= p(\hat{B} | \boldsymbol{\pi}, V, \boldsymbol{\Sigma}) \\ &= \prod_{j=1}^J p(\hat{\mathbf{b}}_j | \boldsymbol{\pi}, V, \boldsymbol{\Sigma}) \\ &= \prod_{j=1}^J \sum_p^P \pi_p N_R(\hat{\mathbf{b}}_j; \mathbf{0}, \Sigma_p + V_j). \end{aligned} \quad (5.5)$$

In our GWAS application we used 50,000 randomly chosen rows. (It is important that this is a random subset, and not the \tilde{J} rows of strong effects used to generate the data-driven \hat{U}_k ; use of the strong effects in this step would be a mistake as it would bias estimates of $\boldsymbol{\pi}$ towards large effect sizes.)

5.3 Results

We identify 4 times as many single associations using MASH as using a univariate method like Ash, 5.3, and about 50% more SNPs that are significant in at least one condition.

Table 5.2: **Associations Identified**

Method	mash	ash
Total Effects ($lfsr_{jr} \leq 0.05$)	.093	0.006
In at Least one ($lfsr_{j.} \leq 0.05$)	0.167	0.025

Table 5.3: **Summary of Significance By Method** Numbers show the proportion of effects meeting a given sharing criterion. “Total Effects” reports the total proportion of significant effects, while ‘In at Least One’ requires that the effect be significant in at least one subgroup across conditions.

Similar to our eQTL analysis, we might wish to understand the primary pattern of sharing evident in the data. In fact, the matrix which received the majority of the loading - U_2 reveals several expected trends in the data: The immune phenotypes, the bone density traits (forearm density, femoral neck, and lumbosacral spine), and the psychological traits (mdd, anx, bip, and scz) also tend to cluster together. Furthermore, we see an inverse relationship between age at menses and BMI. Some of the more unexpected trends are the inverse relationship between autism spectrum disorder (as) and height.

Based on the heavy level of correlation in sign we might want to further investigate which pairs of diseases tend to share effects more closely. Similar to the GTEx eQTL analysis, we can see that most effects are shared by sign but that effects are much less likely to be shared by the absolute value of the magnitude, alluding to the quantitative heterogeneity in magnitude we have described previously.

MASH GWAS: Primary Patterns of Sharing

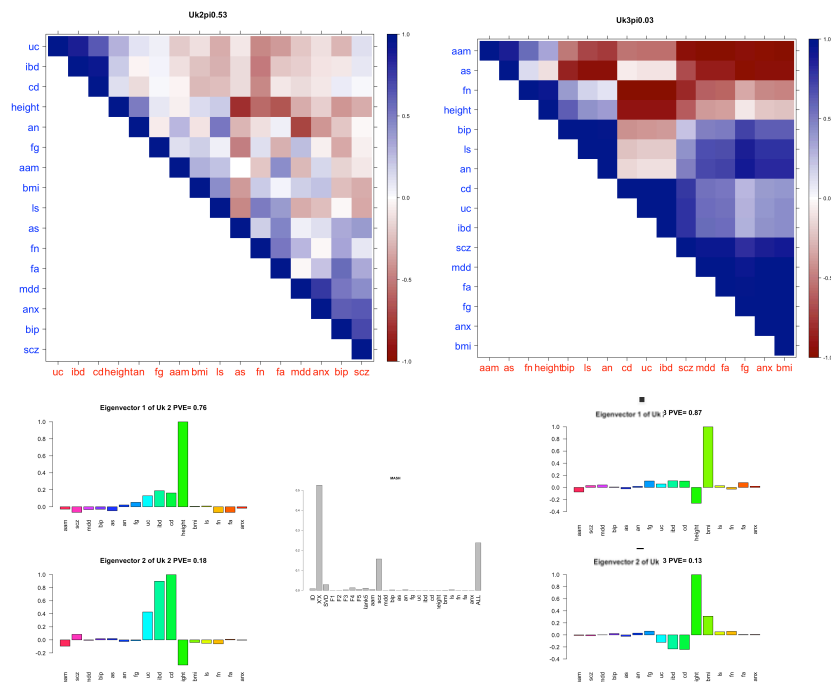


Figure 5.1: **Pairwise Sharing of Effects Across Traits** At left we demonstrate the learned patterns which receive the majority of the weight.

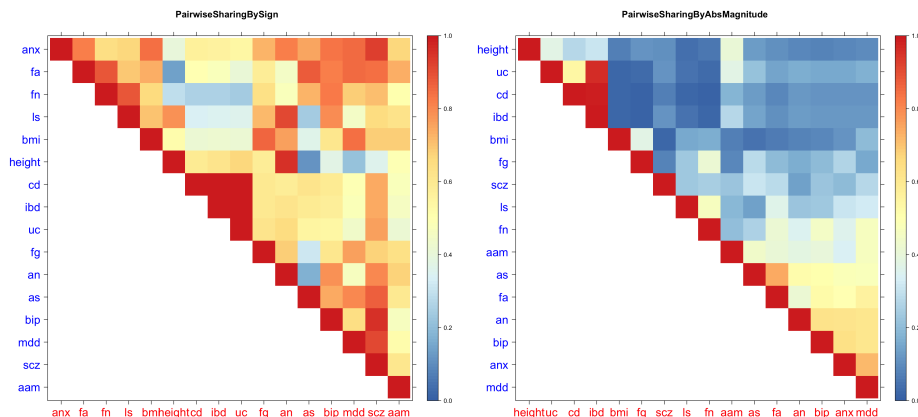


Figure 5.2: **Pairwise Sharing of Effects Across Traits** At left we demonstrate pairwise sharing by sign (left) and absolute value of magnitude (right) across 16 GWAS traits.

5.4 Too much significance?

Our initial concern when analyzing these data was that perhaps we were identifying too many associations. Indeed, if we consider the number of SNPs naively called significant in each trait, we see that the number exceeds 5000 in many cases (Figure 5.3).

We sought to ‘prune’ our list so to remove any associations that were merely within LD of one another. Thus for 12 of the 16 studies, we ascertained published GWAS hits from the GWAS catalog [44] and intersected them with our list of ‘significant hits’ in the corresponding studies in the following manner:

1. Remove anything in our list within 1 Mb of published hits

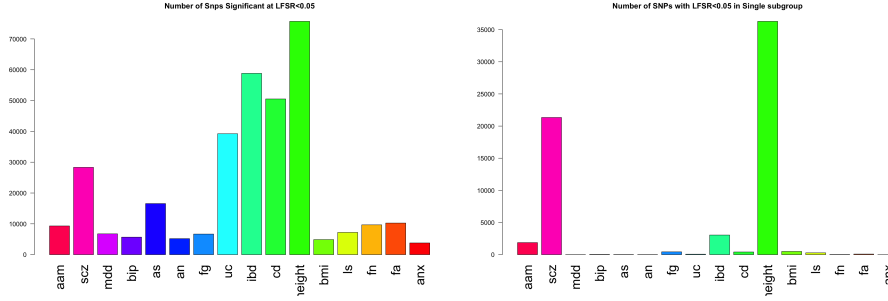


Figure 5.3: **Number of SNPs Found Significant** We display the number of significant (i.e., $lfsr \leq 0.05$) in each condition. At right, we display the number of 'disease-specific' SNPs in any condition - i.e., those SNPs that are significant in that condition alone

2. With remaining hits, select the SNPs with the minimum $lfsr$; remove anything within 1Mb.
3. Continue the process until no SNPs remain.

Here we display the list of significant 'pruned' findings (5.4), and as expected, some of the traits that show the most findings are those to benefit most from the correlated analysis, i.e., Crohn's disease and Ulcerative colitis. Importantly, we note nearly 1500 significant hits in Schizophrenia, a Psychiatric disease that has been difficult to assess in the past due to small sample sizes.

One can then examine the SNPs that might be identified as significant in a joint approach but not in a univariate approach as those that benefit most from sharing. For example, a SNP that showed a small effect in Crohn's disease benefited from a joint analysis (Fig 5.5) in which the larger effect in highly correlated traits (i.e., ulcerative colitis and IBD) boosted our posterior estimates. Similarly, a modest ef-

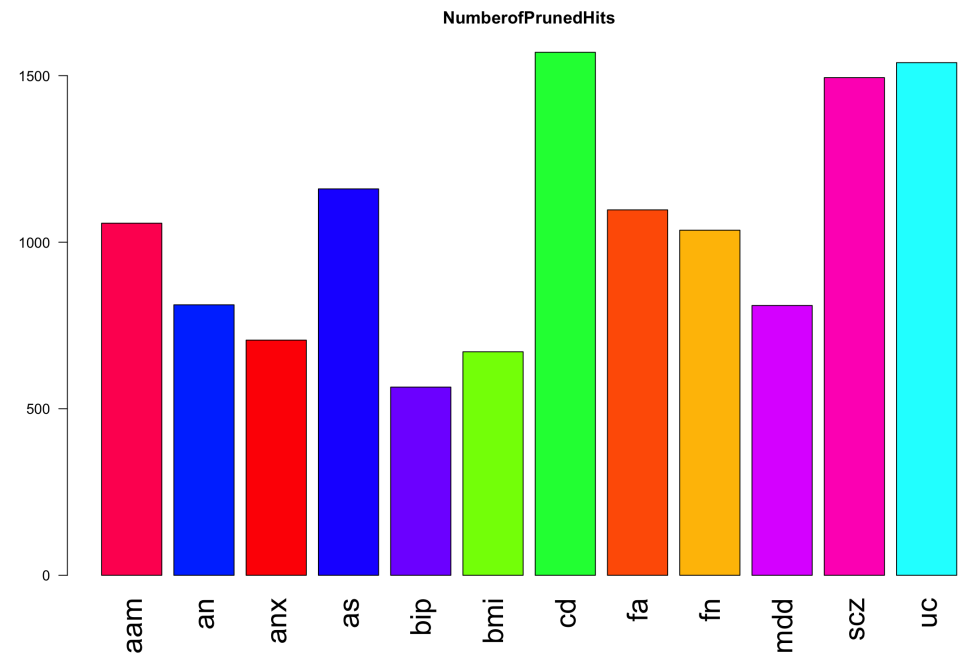


Figure 5.4: **Pruned Significance.** We demonstrate the number of remaining significant results after the pruning exercise described in section 5.4.

fect in femoral neck bone density benefited from sharing with other bone density traits (forearm and lumbosacral spine). Finally, a negligible effect in bipolar disease benefited from boosting by larger effects in highly correlated anxiety and major depression. Taken together, these results emphasize the profound power improvement of a joint analysis.

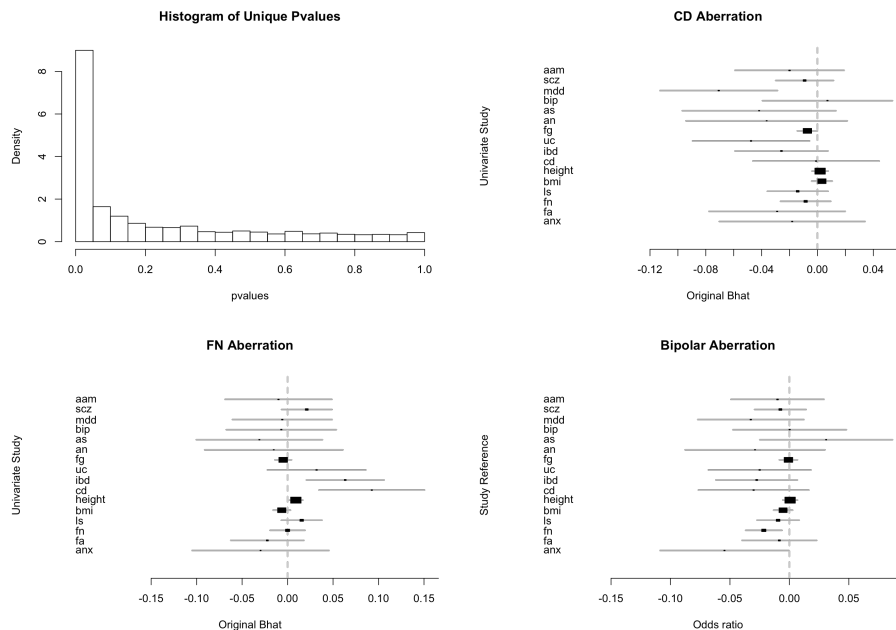


Figure 5.5: **Aberrations** Above, we demonstrate effects for whom univariate analyses concluded lack of significance while our joint analysis found significance. Top left, we examine the distribution of initial univariate p-values for the identified hits. We then demonstrate a metaplot of the initial univariate $\hat{\beta}$ and $\hat{\gamma}$ across subgroups for SNPs initially identified as insignificant in CD, BP and FN (clockwise, top right).

One additional problem that we did not consider is the overlap of individuals between studies. For instance, if there were shared individuals in the CD and IBD cohort, that would create additional correlation in the residual structure that we did not account for. Furthermore, this represents a very early attempt to ‘prune’ for Linkage

Disequilibrium. Further analyses might want to correct for this by identifying LD blocks and removing all but one sentinel marker for identification of an associated block.

CHAPTER 6

CONCLUSION

In the preceding chapters, we have illustrated how a mixture of multivariate normals can be effectively used to powerfully and accurately assess the space of multiple effects. We emphasize novel patterns of sharing that can be identified, as defined by sharing of both magnitude and sharing of sign. While we use these methods to understand genetic effects on gene-expression across multiple types and genetic determinants of disease, such methods can also be used when no obvious baseline is under assessment. Here, we consider the case of gene expression across multiple tissues in which we wish to assess true deviation in expression rather than simply the mean-centered estimates for the gene.

The statistical benefits of joint multivariate analyses compared with univariate analyses are well documented, and increasingly widely appreciated ([45], [5]). As we have shown, not only do we identify additional genetic associations in comparison with univariate approaches (see GTE_x results), we also identify novel biological patterns of sharing. For example, in the GTE_x approach, we show more similarity among brains than non-brain tissues, and identify several additional clusters of biologically-related tissues (see 3.5). We believe this potential nonetheless remains under-exploited in practice. Our aim here is to provide a set of flexible and general tools to help in such analyses, and we designed **mash** with this aim in mind. In particular, **mash** is *generic* and *adaptive*. It is generic in that it can take as input any matrix of Z scores (or,

better, a matrix of effect estimates and their corresponding standard errors) testing many effects in many conditions. For example, the effect estimates we used in our GTEx analysis came from simple linear regressions, but it would be perfectly possible to use `mash` with estimates from other approaches, such as generalized linear models or linear mixed models for example [14]. And `mash` is adaptive in that it learns patterns of sharing of multivariate effects from the data, allowing it to maximize power and precision for each setting. Consequently `mash` should be very widely applicable. Indeed, although genomics applications form our primary motivation, `mash` could be useful in any setting involving testing and estimation of multivariate effects.

Perhaps most importantly, `mash` allows for the estimation of *effects* rather than simply reporting significance, thus giving the user some understanding of the quantitative heterogeneity evident across conditions. These effects are also learned adaptively, in that the pattern from which every vector of observed effects is thought to arise emerges from a mixture of distributions in which many components are learned from the data and reflect data-based patterns.

For instance, the primary patterns of sharing identified in our eQTL analysis (see GTEx results) illustrated that a common pattern present in the data is one of effects that are shared by all tissues but heterogeneous with regard to variation in size between brain and non-brain tissues. Even a joint analysis which simply reported the effects as significant or insignificant might say effects are ‘shared’ by all tissues when in fact there exists a continuum of sizes - biologically, such effect types might help to illustrate genetic effects that are more important in particular tissue types and

thus implicate a particular biological process, while not eliminating the possibility of alternative tissue involvement. We might imagine a genetic effect that had strong activity in certain psychiatric traits with more modest effects in immune traits, suggesting related pathways among both groups of traits with greater similarity within each group.

At its core, **mash** uses an Empirical Bayes hierarchical model, and so is related to other methods that use this approach, including [5, 6, 12]. Indeed, the **mash** framework essentially includes these previous methods as special cases (as well as simpler methods such as “fixed effects” and “random effects” meta-analyses [9, 24]). However, one key feature that distinguishes **mash** from these previous methods is that **mash** puts greater focus on *quantitative* estimation and assessment of effects. More specifically, whereas previous methods have focussed on “binary” models for effects – that is, effects are either present or absent in each condition – **mash** focusses instead on allowing for and assessing quantitative variation among effects. This move away from binary-based models has at least two advantages. First, allowing for all possible binary configurations can create computational challenges. Second, in practice we have found that data often show widespread sharing of effects among many conditions, and that in such settings binary-based methods tend to conclude that effects are non-zero in most or all conditions, even when the signal is very modest in some conditions. This conclusion may not be technically incorrect – for example, in our GTEx analysis it is not impossible that all eQTLs are somewhat active in all tissues. However, as our analysis here illustrates, a more quantitative focus can reveal variation in effect sizes that may be of considerable biological importance.

One important limitation of eQTL analysis here is that it does not distinguish between SNPs that may be causally affecting expression, and those that are merely associated with expression due to being in LD with a causal SNP. This limitation also applies to most previous multi-tissue eQTL analyses, and indeed also to most (but not all) single-tissue eQTL analyses. This issue is particularly important to appreciate when cross-referencing, say, GWAS associations with eQTL effect estimates: a GWAS-associated SNP may be a “significant” eQTL simply because it is in LD with another causal SNP. For single-tissue eQTL mapping, this problem has been addressed in several ways. These include the development of (single-phenotype) fine-mapping methods that attempt to distinguish causal from non-causal effects [46–51], and also co-localization methods [52–54] that attempt to assess whether the same causal SNP may explain an observed association signal in two different phenotypes (e.g. GWAS and gene expression). For multi-tissue analysis, only more limited attempts exist to address this problem. For example, eQTLBMA [5] implements a Bayesian approach to fine-mapping under the simplifying assumption of at most causal SNP per gene [46, 47]. It would be straightforward to adapt `mash` to also perform fine-mapping under this assumption. However, although this simplifying assumption seems a reasonable starting point, it becomes decreasingly plausible in analyses that involve large numbers of tissues, and we view the development of more flexible fine-mapping multi-tissue eQTL methods as an important and challenging problem for future work.

One potentially powerful extension of `mash` would be to allow for the patterns of each effect to depend on covariates. For example, in an eQTL context, one might wish

to allow functional annotations – such as the distance of the SNP from the transcription start site, or its coding/non-coding status – to affect the prior distributions on patterns of sharing or sizes of effects. Furthermore, one would want to estimate the effects of these covariates from the data [47, 55]. One possible way forward here would be to allow the mixture proportions π in **mash** to depend on covariates through a logistic link. However, this appears a challenging problem, and a fully satisfactory solution may require considerable further ingenuity.

Furthermore, we might envision a **mash** iterative approach, in which subgroups that seem to behave in a correlated fashion are added to the assortment of canonical configurations - for example, one could collapse all brain tissues into a single subgroup or add an additional canonical matrix that included non-zero effects in only brain tissues.

Dealing with multiple tests is often described as a “burden”. This description likely originates from the fact that controlling family-wise error rate (the probability of making even one false discovery) requires more and more stringent thresholds as the number of tests increases. However, most modern analyses prefer to control the false discovery rate (FDR) [56], which (under weak assumptions) does not depend on the number of tests [57]. Consequently the term “burden” is inaccurate and unhelpful. Indeed, we believe that the availability of results of many tests in many conditions should be viewed not as a burden, but an *opportunity*: specifically, an opportunity to learn about the relationships among underlying effects, and consequently to make data-driven decisions that help improve both power to detect effects and precision of

effect estimates. Approaches along these lines will inevitably, it seems, involve modeling assumptions, and the goal should be flexible models that are capable of dealing with a wide range of situations that can occur in practice. The methods presented here represent a substantial step towards this goal.

Software implementing our method is available at <http://github.com/stephenslab/mashr>.

Scripts for generating results from the paper are at

https://github.com/surbut/gtexresults_mash.

REFERENCES

- [1] Blischak, J. D., Tailleux, L., Mitrano, A., Barreiro, L. B. & Gilad, Y. Mycobacterial infection induces a specific human innate immune response. *Scientific Reports* **5** (2015). URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4653619/>.
- [2] Ferguson, J. P., Cho, J. H. & Zhao, H. A New Approach for the Joint Analysis of Multiple Chip-Seq Libraries with Application to Histone Modification. *Statistical applications in genetics and molecular biology* **11**, Article–1 (2012). URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3770480/>.
- [3] Pickrell, J. K. *et al.* Understanding mechanisms underlying human gene expression variation with RNA sequencing. *Nature* **464**, 768–772 (2010). URL <http://www.nature.com/nature/journal/v464/n7289/full/nature08872>.
- [4] Dimas, A. S. *et al.* Common regulatory variation impacts gene expression in a cell type-dependent manner. *Science (New York, N.Y.)* **325**, 1246–1250 (2009).
- [5] Flutre, T., Wen, X., Pritchard, J. & Stephens, M. A Statistical Framework for Joint eQTL Analysis in Multiple Tissues. *PLoS Genet* **9**, e1003486 (2013). URL <http://dx.doi.org/10.1371/journal.pgen.1003486>.
- [6] Li, G., Shabalin, A. A., Rusyn, I., Wright, F. A. & Nobel, A. B. An Empirical Bayes Approach for Multiple Tissue eQTL Analysis. *arXiv:1311.2948 [stat]* (2013). URL <http://arxiv.org/abs/1311.2948>. ArXiv: 1311.2948.

- [7] Petretto, E. *et al.* New insights into the genetic control of gene expression using a bayesian multi-tissue approach. *PLOS Computational Biology* **6**, 1–13 (2010). URL
- [8] Wen, X. & Stephens, M. Using Linear Predictors to Impute Allele Frequencies From Summary Of Pooled Genotype Data. *The annals of applied statistics* **4**, 1158–1182 (2010). URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3072818/>.
- [9] Han, B. & Eskin, E. Random-Effects Model Aimed at Discovering Associations in Meta-Analysis of Genome-wide Association Studies. *The American Journal of Human Genetics* **88**, 586–598 (2011). URL [http://www.cell.com/ajhg/abstract/S0002-9297\(11\)00155-8](http://www.cell.com/ajhg/abstract/S0002-9297(11)00155-8).
- [10] Stephens, M. A Unified Framework for Association Analysis with Multiple Related Phenotypes. *PLoS ONE* **8**, e65245 (2013). URL <http://dx.doi.org/10.1371/journal.pone.0065245>.
- [11] Sul, J. H., Han, B., Ye, C., Choi, T. & Eskin, E. Effectively identifying eqtls from multiple tissues by combining mixed model and meta-analytic approaches. *PLOS Genetics* **9**, 1–13 (2013). URL
- [12] Wei, Y., Tenzen, T. & Ji, H. Joint analysis of differential gene expression in multiple studies using correlation motifs. *Biostatistics* **16**, 31–46 (2015).
- [13] Pickrell, J., Berisa, T., Segurel, L., Tung, J. Y. & Hinds, D. Detection

- and interpretation of shared genetic influences on 40 human traits. *bioRxiv* (2015). URL <http://biorxiv.org/content/early/2015/05/27/019885>.
<http://biorxiv.org/content/early/2015/05/27/019885.full.pdf>.
- [14] Han, B. & Eskin, E. Interpreting Meta-Analyses of Genome-Wide Association Studies. *PLOS Genetics* **8**, e1002555 (2012).
- [15] Stephens, M. False Discovery Rates: A New Deal. *bioRxiv* 038216 (2016). URL <http://biorxiv.org/content/early/2016/01/29/038216>.
- [16] Urbut, S. M., Wang, G. & Stephens, M. Flexible statistical methods for estimating and testing effects in genomic studies with multiple conditions. *bioRxiv* 096552 (2016). URL <http://biorxiv.org/content/early/2016/12/24/096552>.
- [17] Bovy, J., Hogg, D. W. & Roweis, S. T. Extreme deconvolution: Inferring complete distribution functions from noisy, heterogeneous and incomplete observations. *The Annals of Applied Statistics* **5**, 1657–1677 (2011). URL <http://projecteuclid.org/euclid.aoas/1310562737>.
- [18] Engelhardt, B. E. & Stephens, M. Analysis of Population Structure: A Unifying Framework and Novel Methods Based on Sparse Factor Analysis. *PLOS Genet* **6**, e1001117 (2010).
- [19] Larribe, F. & Fearnhead, P. Composite likelihood methods in statistical genetics. *Statistica Sinica* **21**, 43–69 (2011).

- [20] Wakefield, J. Bayes factors for genome-wide association studies: comparison with P-values. *Genetic Epidemiology* **33**, 79–86 (2009).
- [21] Dempster, A. P., Laird, N. M. & Rubin, D. B. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B* **39**, 1–38 (1977).
- [22] Varadhan, R. & Roland, C. Squared Extrapolation Methods (SQUAREM): A New Class of Simple and Efficient Numerical Schemes for Accelerating the Convergence of the EM Algorithm. *Johns Hopkins University, Dept. of Biostatistics Working Papers* (2004). URL <http://biostats.bepress.com/jhubiostat/paper63>.
- [23] Efron, B. Local false discovery rates (2005).
- [24] Wen, X. & Stephens, M. Bayesian methods for genetic association analysis with heterogeneous subgroups: From meta-analyses to gene-environment interactions. *The Annals of Applied Statistics* **8**, 176–203 (2014). URL <http://arxiv.org/abs/1111.1210>. ArXiv:1111.1210 [stat].
- [25] Consortium, T. G. The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Science* **348**, 648–660 (2015). URL <http://science.sciencemag.org/content/348/6235/648>.
- [26] Lebec, J. J., Stijnen, T. & van, H. H. C. Dealing with Heterogeneity between Cohorts in Genomewide SNP Association Studies. *Statistical Applications in*

Genetics and Molecular Biology **9** (2010).

- [27] Nicolae, D. L. *et al.* Trait-associated snps are more likely to be eqtls: Annotation to enhance discovery from gwas. *PLoS Genet* **6**, 1–10 (2010). URL
- [28] Shabalin, A. A. Matrix eQTL: ultra fast eQTL analysis via large matrix operations. *Bioinformatics* **28**, 1353–1358 (2012). URL <http://bioinformatics.oxfordjournals.org/content/28/10/1353>.
- [29] Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nature Genetics* **38**, 904–909 (2006).
- [30] Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of expression residuals (PEER) to obtain increased power and interpretability of gene expression analyses. *Nature protocols* **7**, 500–507 (2012). URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3398141/>.
- [31] Tversky, A. & Kahneman, D. Judgment under uncertainty: Heuristics and biases. *Science* **185**, 1124–1131 (1974).
- [32] Bulik-Sullivan, B. *et al.* LD Score Regression Distinguishes Confounding from Polygenicity in Genome-Wide Association Studies. *bioRxiv* (2014). URL <http://biorxiv.org/content/early/2014/02/21/002931>.
- [33] Zhu, X. & Stephens, M. Bayesian large-scale multiple regression with summary statistics from genome-wide association studies. *bioRxiv* 042457 (2016). URL <http://biorxiv.org/content/early/2016/03/04/042457>.

- [34] Dunnett, C. W. A multiple comparison procedure for comparing several treatments with a control. *Journal of the American Statistical Association* **50**, 1096–1121 (1955). URL <http://www.jstor.org/stable/2281208>.
- [35] Katsel, P., Davis, K. L., Gorman, J. M. & Haroutunian, V. Variations in differential gene expression patterns across multiple brain regions in schizophrenia. *Schizophrenia Research* **77**, 241–252 (2005). URL <http://www.sciencedirect.com/science/article/pii/S0920996405001313>.
- [36] McCarthy, D. J., Chen, Y. & Smyth, G. K. Differential expression analysis of multifactor rna-seq experiments with respect to biological variation. *Nucleic Acids Res* **40**, 4288–4297 (2012). URL <http://dx.doi.org/10.1093/nar/gks042>.
- [37] Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140 (2010). URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2796818/>.
- [38] López-Kleine, L. & González-Prieto, C. Challenges Analyzing RNA-Seq Gene Expression Data. *Open Journal of Statistics* **06**, 628 (2016).
- [39] Wei, Y., Tenzen, T. & Ji, H. Joint analysis of differential gene expression in multiple studies using correlation motifs. *Biostatistics* **16**, 31 (2015).
- [40] Smyth, G. K., Ritchie, M., Thorne, N., Wettenhall, J. & Shi, W. Limma:

linear models for microarray data. In *Bioinformatics and Computational Biology Solutions using R and Bioconductor*, 397–420 (Springer, 2005).

- [41] Manolio, T. A. *et al.* Finding the missing heritability of complex diseases. *Nature* **461**, 747–753 (2009). URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2831613/>.
- [42] Eichler, E. E. *et al.* Missing heritability and strategies for finding the underlying causes of complex disease. *Nature Reviews Genetics* **11**, 446–450 (2010).
- [43] Zheng, J. *et al.* Ld hub: a centralized database and web interface to perform ld score regression that maximizes the potential of summary level gwas data for snp heritability and genetic correlation analysis. *Bioinformatics* **33**, 272 (2017).
- [44] Welter, D. *et al.* The nhgri gwas catalog, a curated resource of snp-trait associations. *Nucleic Acids Research* **42**, D1001 (2014).
- [45] Skol, A. D., Scott, L. J., Abecasis, G. R. & Boehnke, M. Joint analysis is more efficient than replication-based analysis for two-stage genome-wide association studies. *Nature Genetics* **38**, 209–213 (2006).
- [46] Servin, B. & Stephens, M. Imputation-Based Analysis of Association Studies: Candidate Regions and Quantitative Traits. *PLoS Genet* **3**, e114 (2007). URL <http://dx.plos.org/10.1371/journal.pgen.0030114>.
- [47] Veyrieras, J.-B. *et al.* High-Resolution Mapping of Expression-QTLs Yields Insight into Human Gene Regulation. *PLoS Genet* **4**, e1000214 (2008). URL

<http://dx.doi.org/10.1371/journal.pgen.1000214>.

- [48] Hormozdiari, F., Kostem, E., Kang, E. Y., Pasaniuc, B. & Eskin, E. Identifying Causal Variants at Loci with Multiple Signals of Association. *Genetics* genetics.114.167908 (2014). URL <http://www.genetics.org/content/early/2014/08/06/genetics.114.167908>.
- [49] Kichaev, G. *et al.* Integrating functional data to prioritize causal variants in statistical fine-mapping studies. *PLOS Genetics* **10**, 1–16 (2014). URL
- [50] Chen, W. *et al.* Fine mapping causal variants with an approximate bayesian method using marginal test statistics. *Genetics* **200**, 719–736 (2015). URL <http://www.genetics.org/content/200/3/719>. <http://www.genetics.org/content/200/3/719.full.pdf>.
- [51] Moyerbrailean, G. A. *et al.* Which genetics variants in dnase-seq footprints are more likely to alter binding? *PLOS Genetics* **12**, 1–27 (2016). URL
- [52] Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. *PLOS Genetics* **10**, 1–15 (2014). URL
- [53] Fortune, M. D. *et al.* Statistical colocalization of genetic risk variants for related autoimmune diseases in the context of common controls. *Nature Genetics* **47**, 839–846 (2015).

- [54] Nica, A. C. *et al.* Candidate causal regulatory effects by integration of expression qtls with complex trait genetic associations. *PLOS Genetics* **6**, 1–11 (2010). URL
- [55] Conesa, A. *et al.* A survey of best practices for RNA-seq data analysis. *Genome Biology* **17**, 13 (2016). URL <http://dx.doi.org/10.1186/s13059-016-0881-8>.
- [56] Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society. Series B (Methodological)* **57**, 289–300 (1995). URL <http://www.jstor.org/stable/2346101>.
- [57] Storey, J. D. The positive false discovery rate: a Bayesian interpretation and the q-value. *The Annals of Statistics* **31**, 2013–2035 (2003). URL <http://projecteuclid.org/euclid.aos/1074290335>.