

Subject Section

A convex optimization framework for gene-level tissue network estimation with missing data and its application in disease architecture

Kushal K. Dey ^{1,*}, Rahul Mazumder ^{2,*}

¹Department of Epidemiology, Harvard T. H. Chan School of Public Health, Boston, MA

²Sloan School of Management, Operations Research Center and Center for Statistics, MIT, Cambridge, MA.

* denotes authors to whom correspondence should be addressed.

Associate Editor: XXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: Genes with correlated expression across individuals in multiple tissues are potentially informative for systemic genetic activity spanning these tissues. In this context, the tissue-level gene expression data across multiple subjects from the Genotype Tissue Expression (GTEx) Project is a valuable analytical resource. Unfortunately, the GTEx data is fraught with missing entries owing to subjects often contributing only a subset of tissues. In such a scenario, standard techniques of correlation matrix estimation with or without data imputation do not perform well. To solve this problem, we propose `Robocov`, a novel convex optimization-based framework for robustly learning sparse covariance or inverse covariance matrices for missing data problems.

Results: `Robocov` produces more interpretable visual representation of correlation and causal structure in simulation settings and GTEx data analysis. We also show that `Robocov` estimators have a lower false positive rate than competing approaches for missing data problems. Genes prioritized by the average value of `Robocov` correlations or partial correlations across tissues are enriched for pathways related to systemic activities such as signaling pathways, circadian clock and immune function. SNPs linked to these prioritized genes showed high enrichment and unique information for blood-related traits; in comparison, no disease signal is observed for SNPs characterized analogously using standard correlation estimator.

Availability: `Robocov` is available as an R package <https://github.com/kkdey/Robocov>.

Contact: kdey@hsp.harvard.edu

Supplementary information: Supplementary data are available at *Bioinformatics* online.

1 Introduction

The gene expression data from nearly 50 tissues across more than 500 post-mortem donor individuals from Genotype Tissue Expression (GTEx) project has proved to be a valuable resource for understanding tissue-specific and tissue-shared genetic architecture [? ? ? ?]. Here we are interested in one specific aspect of tissue-shared gene regulation: the correlation and partial correlation in gene expression for different tissue pairs based on individual donor level data. A major challenge in this context is the extensive amount of missing entries in gene expression data—each donor contributes only a subset of tissues for sequencing. Common imputation based methods do not work well here as reported in [?], owing

to stringent assumptions about missing entries being close to some central tendency (median) or adhering to some low-dimensional representation of the observed entries [? ?]. Popular shrinkage and/or sparse correlation or partial correlation estimators such as `corpcor` [? ?], GLASSO [?] or CLIME [?] are not designed for data with missing values.

A recently proposed approach, `CorShrink` [?], co-authored by one of the authors of this paper (Dey), accounts for missing data through adaptive shrinkage [?] of correlations. `CorShrink` does not guarantee a positive semidefinite (PSD) matrix as part of its EM-based framework, and necessitates a post-hoc modification to ensure a PSD correlation matrix. Furthermore, `CorShrink` does not extend to conditional graph or partial correlation estimation. Here, we propose a new approach based on convex optimization: `Robocov` that applies to both covariance and

inverse covariance matrix estimation in the presence of missing data under the following regularization principles: (a) the covariance matrix is sparse (i.e., has a few nonzero entries) or (b) the inverse covariance matrix is sparse.

`Robocov` does not *impute* missing values per-se¹—it directly estimates the covariance or inverse covariance matrices in the presence of missing values. To handle missing values, we consider a loss function that depends upon the pairwise covariance terms (computed based on the observed samples) but incorporates an adjustment to guard against our lack of knowledge regarding the missing observations. For inverse covariance estimation, `Robocov` uses a robust optimization based approach [? ?] that accounts for the uncertainty in estimating the pairwise sample covariance terms (due to the presence of missing values). Interestingly, both lead to convex optimization formulations that are amenable to modern optimization techniques [?]—they are scalable to moderately-large scale instances; and unlike conventional EM methods (that lead to nonconvex optimization tasks), our estimators attain the global solution of the optimization formulations defining the `Robocov` estimators.

Our experiments suggest that `Robocov` estimators for correlation and partial correlation matrices have lower false positive rate compared to competing approaches for missing data problems. When applied to the GTEx gene expression data with $\sim 70\%$ missing data, `Robocov` produced less cluttered and highly interpretable visualization of correlation and conditional graph architecture. From a biological perspective, a gene with high correlation in expression across many tissue pairs is potentially reflective of more systemic biological processes spanning multiple tissues. To this end, we prioritized genes based on the average `Robocov` estimated correlation (partial correlation) across all tissue-pairs; we call them `Robospan` (`pRobospan`) genes. A pathway enrichment analysis of `Robospan` (`pRobospan`) genes showed enrichment in systemic functional pathways and the immune system. SNPs linked to `Robospan` (`pRobospan`) genes were tested for autoimmune disease informativeness by applying Stratified LD-score regression (S-LDSC) to 11 common blood-related traits (5 autoimmune diseases and 6 blood cell traits; average $N=306K$), conditional on a broad set of annotations. `Robospan` and `pRobospan` genes showed high disease informativeness for blood-related traits. In comparison, `Corspan` genes defined similarly using the standard correlation estimator were non-informative. This highlights the biological and disease-level significance of our work.

2 Methods

Let $X_{N \times P}$ be a data matrix with N samples and P features, where some of the entries X_{np} may be missing, denoted here by NA. Let X^f denote the fully-observed version of the partially-observed data matrix² X . We assume that samples are independent and follow a Multivariate Normal distribution: i.e., $X_{n,*}^f \sim \text{MVN}(0, \Sigma)$ where $\Sigma_{P \times P}$ and $\Omega := \Sigma^{-1}$ (also of size $P \times P$) denote the model covariance and inverse covariance matrices respectively. Based on the observed entries, we obtain a matrix $\hat{\Sigma}$ of pairwise covariances such that for all $i, j \in \{1, \dots, P\}$:

$$\hat{\Sigma}_{ij} := \frac{1}{n_{ij} - 1} \sum_{n: X_{ni} \neq \text{NA}, X_{nj} \neq \text{NA}} (X_{ni} - \bar{X}_i)(X_{nj} - \bar{X}_j) \quad (1)$$

where, \bar{X}_k denotes the sample mean of feature k based on the observed entries; and n_{ij} is the number of samples n with non-missing entries

¹Expectation Maximization (EM) [?] methods typically used for estimation with missing values depend upon probabilistic modeling assumptions and lead to highly nonconvex problems posing computational challenges.

²Note that X is a restriction of X^f to the observed entries.

in both features i and j . Let n_i denote the number of observed samples (i.e., not missing) for feature i . For our analysis, we will assume³ that $n_{ij} > 2$ for all i, j . We note that the matrix of all pairwise covariance terms: $\hat{\Sigma} = ((\hat{\Sigma}_{ij}))$, as defined in (1), need not be positive semidefinite due to the presence of missing values in the data matrix.

2.1 Robocov covariance estimator

We first present the `Robocov` covariance matrix estimator—this leads to an estimate of Σ via the following regularized criterion:

$$\min \sum_{i < j} |\Sigma_{ij}| \quad \text{s.t.} \quad \Sigma \succeq 0, \quad |\hat{\Sigma}_{ij} - \Sigma_{ij}| \leq C_{ij}, \quad \forall i, j \quad (2)$$

where Σ is the optimization variable and C_{ij} s are data-driven constants that control the amount by which Σ_{ij} can differ from the sample version $\hat{\Sigma}_{ij}$. Problem (2) minimizes a convex penalty function (this encourages sparsity [?] in Σ_{ij} s) subject to convex constraints (note that Σ is positive-semidefinite i.e., $\Sigma \succeq 0$). Problem (2) is a convex semidefinite optimization problem [?]; and can be solved efficiently by modern semidefinite optimization algorithms for moderately large instances (e.g, $P \sim 1000$) using (for example) the SCS solver in CVX software [? ? ? ?].

We compute C_{ij} based on the Fisher’s Z-Score [? ?] (for a complete derivation see Supplementary Note):

$$C_{ij} = \hat{\sigma}_i \hat{\sigma}_j \min \left(2, \eta(n_{ij}) \left\{ 3(1 - \hat{R}_{ij}^2) + 2\sqrt{3}\eta(n_{ij}) \right\} \right) \quad (3)$$

where $\eta(n_{ij}) = \sqrt{1/(n_{ij} - 1) + 2/(n_{ij} - 1)^2}$; and \hat{R} is the pairwise sample correlation matrix derived from $\hat{\Sigma}$.

In summary, we note that our proposed `Robocov` estimator does not impute missing values per-se — it directly leads to an estimate for the covariance matrix Σ while taking into account the presence of missing-values in the data matrix.

While (2) leads to a covariance matrix estimator, this can be modified to deliver a correlation matrix instead of a covariance matrix:

$$\begin{aligned} \min \quad & \sum_{i < j} |\mathcal{R}_{ij}| \\ \text{s.t.} \quad & \mathcal{R} \succeq 0, \mathcal{R}_{ii} = 1, \forall i, \quad |\hat{R}_{ij} - \mathcal{R}_{ij}| \leq C_{ij}^{(R)}, \quad \forall i, j \end{aligned} \quad (4)$$

where \mathcal{R} is the optimization variable and $C_{ij}^{(R)} = \frac{\hat{C}_{ij}}{\hat{\sigma}_i \hat{\sigma}_j}$. See the Supplementary Note for additional details.

In the Supplementary Note, we present a framework given by the minimization of a regularized loss function that generalizes the estimator presented in (2).

2.2 Robocov inverse covariance estimator

We present a regularized likelihood framework to estimate the inverse covariance matrix (Ω) under a sparsity constraint. Our optimization criterion is convex in Ω (and not Σ which was the case in Section 2.1).

Recall that GLASSO minimizes an ℓ_1 -norm regularized negative log-likelihood criterion (fully observed case); and is given by:

$$\min_{\Omega \succ 0} -\log \det(\Omega) + \langle \tilde{\Sigma}, \Omega \rangle + \lambda \sum_{ij} |\Omega_{ij}| \quad (5)$$

where, $L(\Omega; \tilde{\Sigma}) := -\log \det(\Omega) + \langle \tilde{\Sigma}, \Omega \rangle$ is the negative log-likelihood (ignoring irrelevant constants), $\tilde{\Sigma}$ is the fully observed sample covariance

³If necessary, as a pre-processing step, we remove features so that the condition $n_{ij} > 2$ is satisfied for all i, j .

matrix and $\lambda \geq 0$ is the regularization parameter. Replacing $\tilde{\Sigma}$ by the observed matrix $\hat{\Sigma}$ in (5) is problematic due to the error in estimating the pairwise covariances arising from the missing values (different cell entries of the sample covariance matrix involve different effective sample sizes n_{ij} s leading to varying accuracies in estimating $\tilde{\Sigma}_{ij}$ s). To account for this uncertainty, we use ideas from robust optimization [?]—to the best of our knowledge, this approach has not been used earlier in the context of sparse inverse covariance estimation (in the presence of missing values). Our robust optimization approach minimizes the worst-case loss arising from the errors in estimating the cell entries $\tilde{\Sigma}_{ij}$ s. This leads to a min-max optimization problem of the form:

$$\min_{\Omega \succeq 0} \max_{|\Delta_{ij}| \leq D_{ij}, \forall i,j} \left\{ -\log \det(\Omega) + \langle \Omega, \hat{\Sigma} + \Delta \rangle \right\} + \lambda \sum_{ij} |\Omega_{ij}|. \quad (6)$$

which is convex [?] in Ω (See Supplementary Note). Convexity ensures that a global minimum to the problem can be obtained reliably—making our approach different from traditional missing data techniques based on the EM algorithm [?] that often lead to complex nonconvex optimization tasks with multiple local solutions.

In words, the inner maximization over Δ in Problem (6) gives the largest (or worst-case) value of the negative log-likelihood— $\max_{\Delta} L(\Omega; \hat{\Sigma} + \Delta)$ where, Δ captures the uncertainty involved in estimating the entries of the sample covariance matrix $\tilde{\Sigma}$ due to the presence of missing values. The outer minimization problem (wrt Ω) considers the minimum of the *adjusted* loss function in addition to an ℓ_1 -penalization on Ω that encourages a sparse estimate of Ω .

The so-called uncertainty set [?] in Δ is given by: $|\Delta_{ij}| \leq D_{ij}$ (for all i, j) where, the upper bound D_{ij} arises from a probability computation using the Fisher’s Z-score criterion (see Supplementary Note):

$$\begin{aligned} D_{ij} &= C_{ij} + \tilde{C}_{ij} \\ \tilde{C}_{ij} &= \hat{\sigma}_i \hat{\sigma}_j \min \left\{ 2, \eta(N) \left\{ 3(1 - \hat{R}_{ij}^2) + 2\sqrt{3}\eta(N) \right\} \right\}. \end{aligned} \quad (7)$$

Above, the value of the error D_{ij} will be large if n_{ij} is small, and equal to zero when $n_{ij} = n$ (with no missing entries).

The seemingly complicated min-max optimization problem in (6) reduces to a cousin of the GLASSO criterion (See Supplementary Note for details)—we use a weighted version of the ℓ_1 -norm penalty on Ω :

$$\min_{\Omega \succeq 0} \left\{ -\log \det(\Omega) + \langle \Omega, \hat{\Sigma} \rangle + \sum_{ij} (\lambda + D_{ij}) |\Omega_{ij}| \right\}. \quad (8)$$

Problem (8) is a nonlinear semidefinite optimization problem in Ω —and the constraint $\Omega \succeq 0$ leads to a positive semidefinite inverse covariance matrix⁴. Problem (8) uses a weighted ℓ_1 -norm on Ω where the penalty weights are adjusted to account for the uncertainty due to the presence of missing values. Note that the penalty parameter λ accounts for the sparsity in Ω arising from our prior sparsity assumption on Ω —the overall penalty weight for the (i, j) -th entry: $(\lambda + D_{ij})$ adds further regularization due to the presence of missing values.

Note that, as in Section 2.1, the Robocov inverse covariance estimator, bypasses the task of imputing the missing values. Our main goal is to directly estimate Ω from a partially observed data-matrix X . In this way, we can potentially mitigate the limitations of a sub-optimal imputation procedure. See Section 3 for an empirical validation.

The inverse covariance estimate Ω from (8), can be used to obtain the partial correlation estimator \mathcal{W} as follows

⁴We get a positive semidefinite (PSD) estimate for Ω even if $\hat{\Sigma}$ is not PSD. The log det-term in the objective encourages an optimal solution to (8) to be positive definite (i.e., of full rank).

$$\mathcal{W}_{ij} := -\frac{\Omega_{ij}}{\sqrt{\Omega_{ii}\Omega_{jj}}}. \quad (9)$$

Problem (8) was solved using R implementation of the CVX software [?]. This was sufficient for the problem-scales we are dealing with.

In all our subsequent analysis and numerical results, we use the Robocov correlation estimator \mathcal{R} (see Problem (4)) and partial correlation estimator \mathcal{W} (9).

3 Results

Simulation Experiments: Synthetic and Real Data

We applied Robocov on simulated multivariate normal data from three population correlation structure models (hub, Toeplitz and 1-band precision matrix) with N samples, P features and π proportion of missing entries randomly distributed throughout the data matrix (Supplementary Note). Figure 1 shows results for all three model-settings with $N = 500$, $P = 50$, $\pi = 0.5$. In all cases, Robocov generated a sparse estimate of the population correlation \mathcal{R} (Section 2.1) or partial correlation \mathcal{W} (Section 2.2). The Robocov correlation estimator captured population structure more effectively for all three models compared to the standard pairwise sample correlation estimator (Figure S1). The Robocov partial correlation estimator also accurately captured the causal structure in the hub and 1-band precision matrix models; for the Toeplitz matrix, it recovered the high partial correlation band immediately flanking the diagonal but not the other alternating positive and negative low correlation bands (Figure 1).

Recent work [?] has shown hub-like patterns in expression correlation across tissue pairs for most genes. To this end, we applied Robocov on simulated data for hub population correlation matrix structure for different settings of N , P and π (Supplementary Note). Two metrics of particular interest were the false positive rate (FPR) and the false negative rate (FNR) (Supplementary Note). We used these metrics to compare Robocov correlation estimator with both the pairwise sample correlation estimator and CorShrink[?]. Across different (N, P, π) -settings, the Robocov correlation estimator had lower FPR than CorShrink. In comparison, for data with a large number of missing entries (i.e., high π), FNR for Robocov was worse compared to CorShrink (Table 1). We did not compare against other shrinkage-based correlation estimators such as PDSCE[?] and corpcor[? ?] as (i) they do not account for missing entries in the data and have been shown to be sub-optimal to CorShrink for fully observed data (see Figure 4 from ref.[?]).

Next, we assess the performance of the Robocov partial correlation estimator for the same simulation settings (Table 1). We are not aware of a sparse conditional graph or partial correlation estimation method that directly takes into account missing entries. Nevertheless, we compare the Robocov partial correlation estimator with (i) GLASSO on the pairwise sample correlation estimator $\hat{\Sigma}$ and (ii) CLIME on an imputed data matrix where, the imputation is performed using SoftImpute [?]. In the presence of missing data, Robocov partial correlation estimator showed better FPR and FNR compared to both GLASSO and CLIME-based estimators (Table 1). The underperformance of CLIME may be attributed to the error arising from the imputation step (Table 1).

Next, we evaluate the predictive performance of Robocov correlation estimator with pairwise sample correlation estimator and CorShrink. We considered the GTEx gene expression data for an example gene (ARHGAP30) across 544 donors and 53 tissues with close to 70% missing data owing to subjects contributing only a small fraction of tissues. We split the individual by tissue data for the gene into two equal groups and compared the estimated correlation matrix (we used different estimators: Robocov, CorShrink and pairwise sample correlation

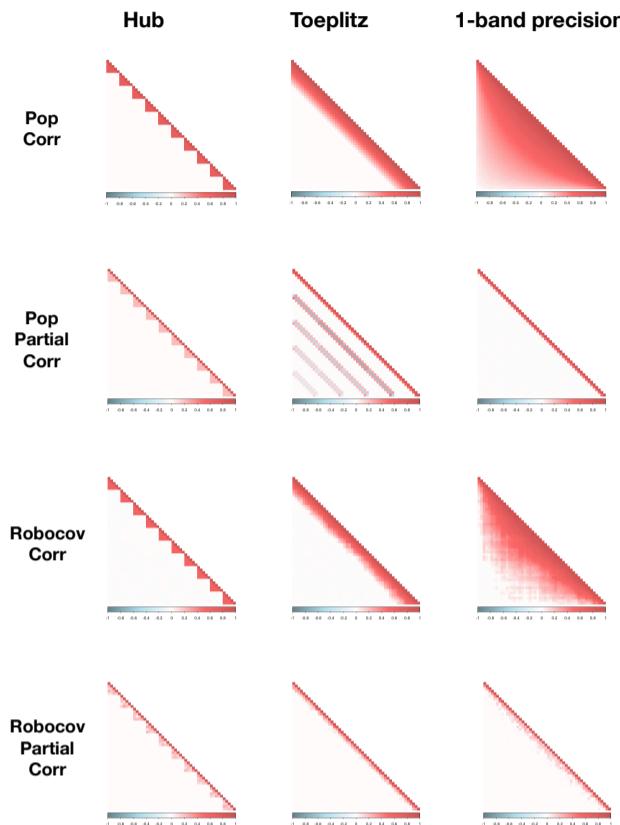


Fig. 1. We applied `Robocov` correlation and partial correlation estimators on data generated from Hub, Toeplitz or 1-band precision matrix based population models (Supplementary Note) with $N = 500$ samples, $P = 50$ features and $\pi = 0.5$ proportion of missing data. We present the population correlation matrix, population partial correlation matrix, `Robocov` correlation matrix and `Robocov` partial correlation matrix sequentially from first to last row.

matrix) computed on one half of the individuals with the pairwise sample correlation matrix computed from the other half. Both `Robocov` and `CorShrink` estimators considerably outperformed the pairwise sample correlation estimator, with `CorShrink` having slightly better predictive accuracy (Figure S2 and Table S1). As `Robocov` and `CorShrink` predictive performances are similar, the former may be preferable as it results in sparse estimates, leading to better interpretability.

An alternative to `Robocov`, we may consider an estimator obtained by first imputing the missing entries in the data matrix and then estimating the correlation or partial correlation matrix for the complete data. For the same ARHGAP30 gene, we performed imputation by either a low rank factorization (SoftImpute[?], with or without scaling) or a median based approach (replacing the missing entries of a feature by the median value of the observed entries). The correlation matrix obtained by SoftImpute (both with and without scaling) showed artificial high negative and positive correlation sweeps between brain and non-brain tissues that were not observed in the pairwise correlation matrix (Figure S3). One possible explanation of this is that the data matrices in our case do not seem to have a low rank representation based on eigenvalue analysis (Figure S4). The median based imputation method on the other hand, is prone to showing false positives—for example, we see a high correlation between Fallopian tube and Cervix-Ectocervix, which is a consequence of only 3 individuals contributing both the tissues (Figure S3). `Robocov` can effectively get rid of these edge cases and generate sparser and more robust results compared to these imputation based approaches.

Table 1. We compare three metrics: FP2 (False Positive 2-norm), FPR (False Positive Rate) and FNR (False Negative Rate) (Supplementary Note) to compare (i) the `Robocov` correlation estimator (Cor) against `CorShrink` and the standard pairwise sample correlation estimator; and (ii) the `Robocov` partial correlation estimator (PCor) against estimators available from GLASSO and CLIME. Data was generated for different (N , P , π) settings and results were averaged over 50 replications from same model. Optimal λ was chosen by cross-validation.

Hub: N = 50, P=50										
Type	Method	$\pi=0$			$\pi=0.25$			$\pi=0.5$		
		FP2	FPR	FNR	FP2	FPR	FNR	FP2	FPR	FNR
Cor	<code>Robocov</code>	0.05	0	0	0.14	0	0.14	0.26	0	0.19
	<code>CorShrink</code>	1.4	0.01	0	2.2	0.04	0.03	4	0.07	0.09
	Standard	6.7	0.24	0	8.8	0.30	0	15	0.28	0
PCor	<code>Robocov</code>	0.08	0	0.07	0.27	0.01	0.13	0.47	0	0.09
	<code>GLASSO</code>	0.12	0	0.15	0.29	0.01	0.15	0.59	0.02	0.12
	<code>CLIME</code>	1.5	0.09	0.07	1.4	0.07	0.08	1.3	0.08	0.07
Hub: N = 100, P=50										
Type	Method	$\pi=0$			$\pi=0.25$			$\pi=0.5$		
		FP2	FPR	FNR	FP2	FPR	FNR	FP2	FPR	FNR
Cor	<code>Robocov</code>	0.05	0	0	0.06	0	0	0.18	0	0.15
	<code>CorShrink</code>	0.9	0	0	1.3	0.02	0	2.9	0.03	0.01
	Standard	4.8	0.17	0	6.2	0.20	0	10	0.31	0
PCor	<code>Robocov</code>	0.23	0	0.06	0.21	0	0.09	0.18	0.03	0.11
	<code>GLASSO</code>	0.11	0	0.16	0.23	0	0.22	0.29	0.01	0.24
	<code>CLIME</code>	1.8	0.12	0.08	1.8	0.14	0.09	1.8	0.16	0.11
Hub: N = 500, P=50										
Type	Method	$\pi=0$			$\pi=0.25$			$\pi=0.5$		
		FP2	FPR	FNR	FP2	FPR	FNR	FP2	FPR	FNR
Cor	<code>Robocov</code>	0.03	0	0	0.01	0	0	0.08	0	0
	<code>CorShrink</code>	0.21	0	0	0.32	0	0	0.83	0	0
	Standard	2.1	0.01	0	2.8	0.05	0	4.4	0.14	0
PCor	<code>Robocov</code>	0.12	0	0.11	0.16	0	0.12	0.11	0	0.14
	<code>GLASSO</code>	0.16	0	0.19	0.29	0	0.20	0.19	0.02	0.20
	<code>CLIME</code>	2.1	0.11	0.16	2.0	0.14	0.18	2.0	0.15	0.17

Gene Expression correlation analysis across tissue pairs

We applied `Robocov` to each of 16,069 cis-genes (genes with at least one significant cis-eQTL) from the GTEx v6 project [?] (see URLs). For each gene, the data matrix had 544 rows (post-mortem donors), 53 columns (tissues) and comprised of ~ 70% missing entries. Figure 2 presents a visual comparison of `Robocov` correlation and partial correlation estimators with standard pairwise sample correlation matrix for two example genes (ARHGAP30 and GSTM1)—the `Robocov` estimators are sparse and visually less cluttered than the standard approach. The `Robocov` correlation structure across tissue pairs varied from one gene to another: some genes showed high correlation across all tissues (e.g. HBB, RPL9), some showed little to no correlation across tissues (e.g. NCCRP1), some showed high intra-Brain correlation but relatively low inter-Brain correlation (e.g. ARHGAP30) (Figures 3, S5 and S2). Additionally, two genes with similar correlation profiles may have very distinct expression profiles. For example, HBB and RPL9 both showed high correlation across all tissue pairs, but they had very distinct tissue-specific expression profiles. HBB showed high expression in Whole Blood relative to other tissues, while RPL9 had a more uniform expression profile across tissues (Figure 3). A similar pattern was observed also for two genes with negligible correlation across tissues, NCCRP1 and RPL21P11 (Figure S5).

Next, we assign to each gene, a prioritizing score defined by the average value of `Robocov` correlation (*Robospan-score*) or partial correlation (*pRobospan-score*) across all tissue pairs. Similarly, we also computed the

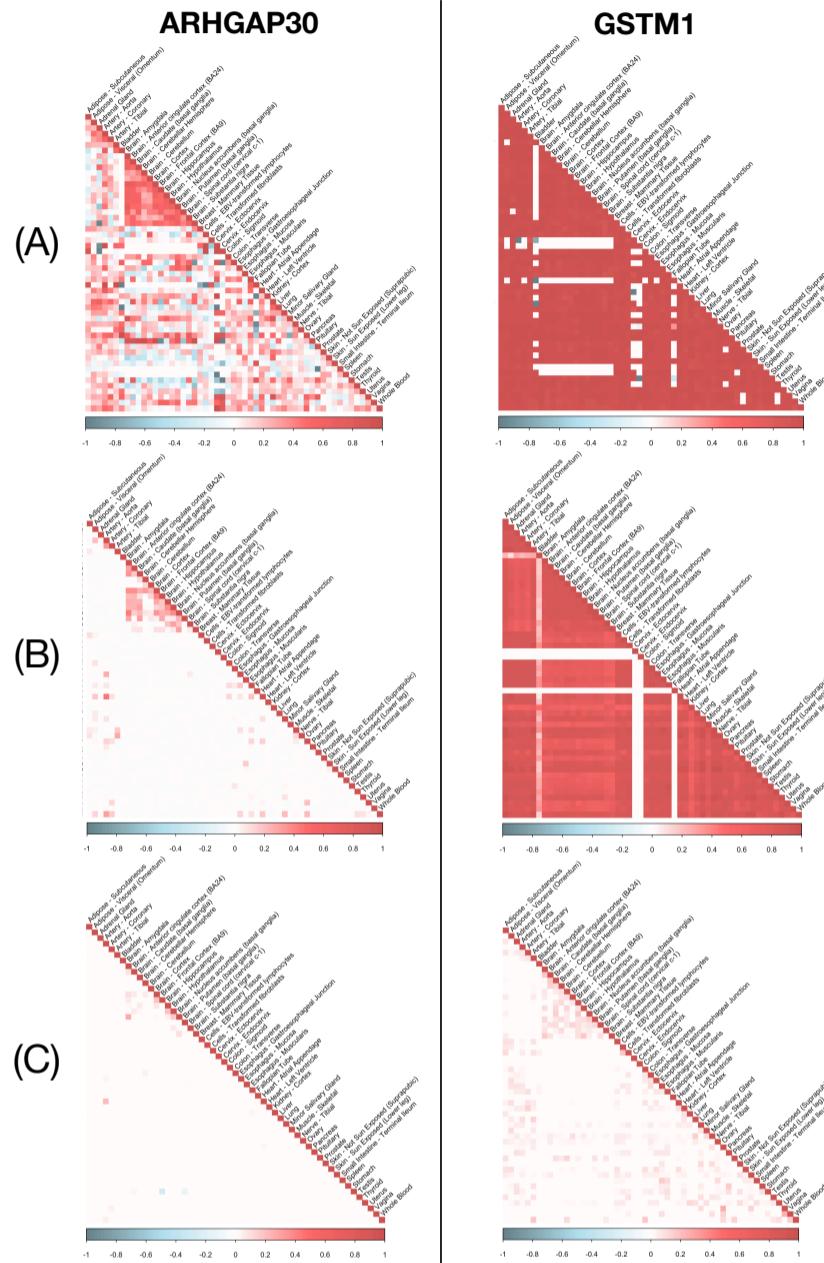


Fig. 2. Illustrative examples of pairwise sample correlation estimator, Robocov correlation and partial correlation estimators for 2 genes: (Left column) ARHGAP30 gene and (Right column) GSTM1 gene. Each column shows the (A) pairwise sample correlation estimator, (B) Robocov correlation estimator and (C) partial correlation estimator stacked from top to bottom.

average value of the pairwise sample correlation (*Corspan-score*) across tissues. Then we tested these gene scores for functional relevance. Contrary to expectation, none of the three scores showed significant enrichment in 3,804 housekeeping genes[?] (0.84x, 0.48x and 0.72x for Robospan-score, pRobospan-score and Corspan-score respectively). We compared these 3 gene scores with constraint-based metric of gene essentiality such as the absence of loss-of-function(LoF) variants (pLI[?] and s_het[?]). For each of the 50 quantile bins of pLI and s_het, we computed the median of each of these scores; and compared with the mid-value of the quantile bin. We observed a slight negative trend in all 3 scores with increasing quantile bins of both pLI and s_het (Figure S6). One possible explanation may be that genes with highly correlated expression across all tissues

may be driven by tissue-shared regulation machinery which imposes lower selective constraints on these genes. The top 10% genes from each of the three gene prioritizing scores were used to define gene sets; we call them Robospan, pRobospan and Corspan genes. In a pathway enrichment analysis[?] of these gene sets, the top enriched pathways comprised of immune system, interferon signaling, heat stress factor (Table S2). Though not among the top 5 pathways, other interesting significant pathways included different signaling pathways (interleukin mediated signaling, NFKB signaling) and circadian clock related pathways (see URLs). The significance of pathway enrichment was stronger for Robospan and pRobospan genes compared to Corspan (Table S2). The enrichment of immune related pathways was further backed by high enrichment of

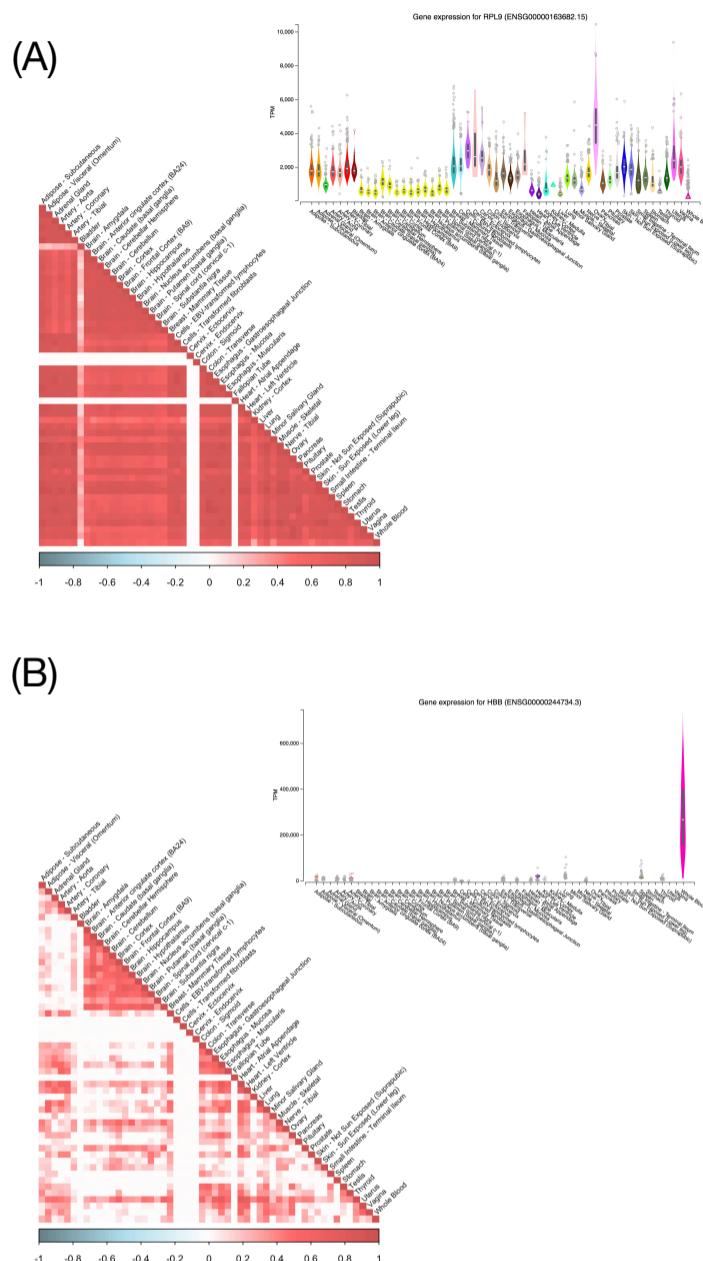


Fig. 3. Examples of genes with high average Robcov correlation across all tissue pairs but with distinct expression profiles. (A) RPL9 gene has uniformly high TPM (transcripts per million) values across most tissues (inset picture). (B) HBB shows high expression specifically in Whole Blood (inset picture). The expression profile plots for the genes have been fetched from the GTEx Portal (<https://gtexportal.org/home/>).

these genes in top 10% specifically expressed genes in Whole Blood (SEG-Blood[?]) (Robospan: 1.48x, pRobospan: 2.50x, Corspan: 1.45x). One may conjecture that this enrichment is an artifact caused by contamination of blood with GTEx tissue samples. This, however, is countered by examples of genes that have high correlation across all tissues but expression-wise, are specific to tissues that are not Whole Blood (Figure S7). We also see examples of specifically expressed genes in Whole Blood that have low Robospan-score (Figure S8).

Heritability analysis of blood-related traits

The enrichment of Robospan, pRobospan and Corspan genes with SEG-Blood genes and immune related pathways prompted us to test

whether these genes are uniquely informative for blood-related complex diseases and traits.

For each gene set, we define SNP-level annotations to test for disease heritability. We define an *annotation* as an assignment of a numeric value to each SNP with minor allele count ≥ 5 in a 1000 Genomes Project European reference panel[? ?]. For each gene set X, we generate two binary SNP-level annotations – we assign a value of 1 to a SNP if it lies within 5kb or 100kb window upstream and downstream of a gene in the gene set and 0 otherwise; this strategy has been used in several previous works[? ? ?].

We assessed the informativeness of SNP annotations for disease heritability by applying stratified LD score regression (S-LDSC) [7] conditional on 86 baseline annotations comprising of coding, conserved,

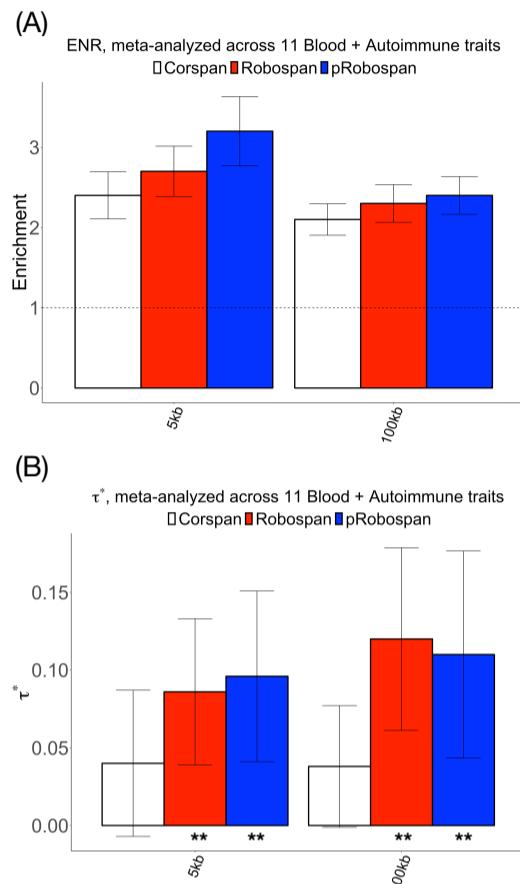


Fig. 4. Disease informativeness of 5kb and 100kb SNP annotations for **Corspan**, **Robospan** and **pRobospan** gene sets: (A) Heritability enrichment, conditional on baseline-LD model (v2.1). The base enrichment level is 1. (B) Standardized effect size (τ^*) conditional on baseline-LD model for **Corspan** (left column, white), **Robospan** (middle column, red) and **pRobospan** (right column, blue) gene sets. Results are meta-analyzed across 11 blood and autoimmune traits. ** denotes annotations that are significant after Bonferroni correction ($P < 0.05/8$) where 8 is the total number of SNP annotations tested. Error bars denote 95% confidence intervals. Numerical results are reported in Table S4.

epigenomic and LD related annotations (this is called the baseline-LD model; here we use version 2.1[?]). S-LDSC results were meta-analyzed across 11 relatively independent blood-related traits (5 autoimmune diseases and 6 blood traits (Table S3). We considered two S-LDSC metrics for comparison: enrichment and standardized effect size (τ^*) (Supplementary Note). Enrichment is defined as the proportion of heritability explained by SNPs in an annotation divided by the proportion of SNPs in the annotation[?]. Standardized effect size (τ^*) is defined as the proportionate change in per-SNP heritability associated with a 1 standard deviation increase in the value of the annotation, conditional on other annotations included in the model[? ?]; unlike enrichment, τ^* quantifies effects that are unique to the focal annotation and is a better metric for disease informativeness[? ? ?].

All 6 annotations (5kb and 100kb for the 3 gene scores) were significantly enriched when meta-analyzed across 11 blood and autoimmune traits. However, SNP annotations corresponding to **Robospan** and **pRobospan** gene sets showed higher enrichment than **Corspan** genes (Figure 4 and Table S4). More importantly, 2 **Robospan**, 2 **pRobospan** and 0 **Corspan** annotations showed significant τ^* conditional on the baseline-LD annotations after Bonferroni correction

(Figure 4 and Table S4). When restricted to the 5 autoimmune traits, 2 **Robospan**, 0 **pRobospan** and 0 **Corspan** SNP annotations showed unique signal (Table S5). Even when these annotations were modeled jointly with SEG-Blood[?] genes and subjected to forward stepwise elimination similar to ref.[? ?], 1 **Robospan** annotation (100kb) still remains significantly informative, suggesting unique disease information over SEG-Blood genes (Table S6).

4 Discussion

Here we present **Robocov**—a novel convex optimization-based framework for sparse estimation of covariance (correlation) and inverse covariance (partial correlation) matrix, given a data matrix with missing entries. Our approach does not rely on missing data imputation and hence mitigates the possible shortcomings of a sub-optimal imputation procedure (e.g., based on a low-rank assumption). Instead, **Robocov** directly estimates the correlation or partial correlation matrix of interest via a regularized loss minimization framework. Although here we focus our analysis on gene expression analysis, **Robocov** is a stand-alone generic tool that can be applied to any data with missing entries.

We have assessed the significance of our proposed **Robocov** framework over standard methods from a methodological, biological and disease analysis perspective. **Robocov** leads to sparse estimates and has a lower false positive rate compared to other competing methods. **Robocov** estimator is visually less cluttered and captures more robust biological signal. In terms of disease informativeness, **Robospan** and **pRobospan** gene sets, generated from the **Robocov** estimated correlation and partial correlation matrices, perform considerably better than the analogous **Corspan** gene set defined from standard correlation estimator.

There are several directions for future research. One such direction would be to incorporate covariate information underlying structured missing-ness to inform **Robocov** estimators. For GTEx data, donor metadata such as cause of death, age, gender etc can serve as important covariates. Second, we are interested in modifying **Robocov** to learn shared correlation structure between gene expression and other genetic and epigenomic data such as transcript level expression, ATAC-seq data etc. Third, from application standpoint, **Robocov** can also be used as an ingredient in item response models for large scale participant data that may contain extensive amount of missing entries, as in UK Biobank [? ?].

- Robocov software

<https://github.com/kkdey/Robocov>

- GTEx v6 data analysis, gene list, pathway enrichment results, gene sets, annotations

<https://github.com/kkdey/Robocov-pages>

- Baseline-LD annotations:

<https://data.broadinstitute.org/alkesgroup/LDScore/>

- Summary statistics:

https://data.broadinstitute.org/alkesgroup/sumstats_formatted/

Acknowledgements

We thank Alkes L. Price, Bryce van de Geijn and Rajarshi Mukherjee for helpful comments. Rahul Mazumder was partially supported by the Office of Naval Research ONR-N000141512342, ONR-N000141812298 (Young Investigator Award), the National Science Foundation (NSF-IIS-1718258) and IBM.

Supplementary Note

Fisher Z-score

The population Fisher Z-score[?] is defined as

$$Z_{ij} = \frac{1}{2} \log \left[\frac{1 + R_{ij}}{1 - R_{ij}} \right] \quad (10)$$

where R is the population correlation matrix. The corresponding empirical Fisher Z-score is defined as follows

$$\hat{Z}_{ij} = \frac{1}{2} \log \left[\frac{1 + \hat{R}_{ij}}{1 - \hat{R}_{ij}} \right] \quad (11)$$

For bivariate normally distributed random variables X_i and X_j , the empirical Fisher Z-score \hat{Z}_{ij} (based on n_{ij} -many samples) is normally distributed given the population counterpart Z_{ij} [?]:

$$\hat{Z}_{ij} | Z_{ij} \sim N \left(Z_{ij}, \frac{1}{n_{ij} - 1} + \frac{2}{(n_{ij} - 1)^2} \right); \quad (12)$$

and the Z-scores are conditionally independent. Dey and Stephens [?] assume an adaptive shrinkage prior on the population Fisher Z-scores for each pair of variables. Here we use property (12) in the context of directly estimating Σ or Ω with an ℓ_1 -norm penalty.

Derivation of C

Here we show how we derive the analytical form of the upper bound C in (3) appearing in Problem (2).

Lemma 1. Let $X_{N \times P}^f$ be the fully observed version of the data matrix X ; and let every sample $X_{n_{i,j}}^f$ follow a Multivariate Gaussian distribution with covariance matrix Σ and correlation matrix R . The samples are independent. Then, for any fixed $\epsilon > 0$ and for sufficiently large n_{ij} , there exists a $C'_{ij}(\epsilon)$ such that

$$\Pr \left(|\hat{R}_{ij} - R_{ij}| \leq C'_{ij}(\epsilon) \mid R_{ij} \right) > (1 - \epsilon) \quad (13)$$

where

$$C'_{ij}(\epsilon) := \min \left(2, \eta(n_{ij}) M(\epsilon) \left\{ (1 - \hat{R}_{ij}^2) + \frac{2M(\epsilon)}{3\sqrt{3}} \eta(n_{ij}) \right\} \right) \quad \forall i \neq j \quad (14)$$

and

$$\eta(n_{ij}) := \sqrt{\frac{1}{n_{ij} - 1} + \frac{2}{(n_{ij} - 1)^2}} \quad (15)$$

and $M(\epsilon)$ is a sufficiently large finite number.

Corollary 1. For $\epsilon = 0.001$, $M(\epsilon)$ can be taken to be 3 in Lemma 1. Then

$$\Pr \left(|\hat{R}_{ij} - R_{ij}| < C'_{ij} \mid R_{ij} \right) \approx 1 \quad (16)$$

where

$$C'_{ij} := \min \left(2, \eta(n_{ij}) \left\{ 3(1 - \hat{R}_{ij}^2) + 2\sqrt{3}\eta(n_{ij}) \right\} \right) \quad \forall i \neq j \quad (17)$$

If n_i and n_j are sufficiently large, in which case $\hat{\sigma}_i \approx \sigma_i$ and $\hat{\sigma}_j \approx \sigma_j$, then Corollary 1 leads to the following probability inequality for the pairwise sample covariance:

$$\Pr \left(|\hat{\Sigma}_{ij} - \Sigma_{ij}| < C_{ij} \mid \Sigma_{ij} \right) \approx 1 \quad (18)$$

where

$$C_{ij} := \hat{\sigma}_i \hat{\sigma}_j C'_{ij}. \quad (19)$$

Proof of Lemma 1 and Corollary 1

If a random variable $W \sim N(0, 1)$, then for any small $\epsilon > 0$, we can get a number $M(\epsilon)$ such that

$$\Pr(|W| < M(\epsilon)) > (1 - \epsilon) \quad (20)$$

Using (12) and (20), we have

$$\Pr \left(|\hat{Z}_{ij} - Z_{ij}| < M(\epsilon) \eta(n_{ij}) \mid Z_{ij} \right) > (1 - \epsilon). \quad (21)$$

The estimated and population correlations \hat{R}_{ij} and R_{ij} (respectively) can be written in terms of the Z-scores using (10) as follows:

$$\hat{R}_{ij} = \frac{\exp(2\hat{Z}_{ij}) - 1}{\exp(2\hat{Z}_{ij}) + 1}, \quad R_{ij} = \frac{\exp(2Z_{ij}) - 1}{\exp(2Z_{ij}) + 1}. \quad (22)$$

Applying a Taylor series expansion to R_{ij} as a function of Z_{ij} around \hat{Z}_{ij} , we get:

$$\begin{aligned} \frac{\exp(2Z_{ij}) - 1}{\exp(2Z_{ij}) + 1} &= \frac{\exp(2\hat{Z}_{ij}) - 1}{\exp(2\hat{Z}_{ij}) + 1} + 4 \frac{\exp(2\hat{Z}_{ij})}{\exp(2\hat{Z}_{ij}) + 1} (\hat{Z}_{ij} - Z_{ij}) \\ &\quad + 4 \frac{\exp(2\xi)(\exp(2\xi) - 1)}{(\exp(2\xi) + 1)^3} (\hat{Z}_{ij} - Z_{ij})^2 \end{aligned} \quad (23)$$

where ξ is a value between Z_{ij} and \hat{Z}_{ij} . We can place an upper bound on the coefficient of the last term in (23):

$$\left| \frac{\exp(2\xi)(\exp(2\xi) - 1)}{(\exp(2\xi) + 1)^3} \right| \leq \frac{1}{6\sqrt{3}}. \quad (24)$$

Using Equations (22), (23) and (24), we can write

$$\begin{aligned} |\hat{R}_{ij} - R_{ij}| &\leq 4 \frac{\exp(2\hat{Z}_{ij})}{(\exp(2\hat{Z}_{ij}) + 1)^2} |\hat{Z}_{ij} - Z_{ij}| \\ &\quad + \frac{2}{3\sqrt{3}} |\hat{Z}_{ij} - Z_{ij}|^2 \end{aligned} \quad (25)$$

Using the definition of \hat{Z}_{ij} in Equation (11), we get

$$\frac{\exp(2\hat{Z}_{ij})}{(\exp(2\hat{Z}_{ij}) + 1)^2} = \frac{(1 - \hat{R}_{ij}^2)}{4}. \quad (26)$$

Using the above expression in (25), we get:

$$|\hat{R}_{ij} - R_{ij}| \leq (1 - \hat{R}_{ij}^2) |\hat{Z}_{ij} - Z_{ij}| + \frac{2}{3\sqrt{3}} |\hat{Z}_{ij} - Z_{ij}|^2 \quad (27)$$

Using (21) and (27), we have:

$$\begin{aligned} \Pr \left(|\hat{R}_{ij} - R_{ij}| < (1 - \hat{R}_{ij}^2) M(\epsilon) \eta(n_{ij}) + \frac{2}{3\sqrt{3}} M^2(\epsilon) \eta^2(n_{ij}) \mid R_{ij} \right) \\ > (1 - \epsilon). \end{aligned}$$

Since, \hat{R}_{ij} and R_{ij} are both correlation terms, they lie between -1 and $+1$ and hence with probability one:

$$|\hat{R}_{ij} - R_{ij}| \leq 2 \quad (28)$$

Combining Equations (27) and (28), we get

$$\Pr \left(|\hat{R}_{ij} - R_{ij}| < \min \{2, B\} \mid R_{ij} \right) > (1 - \epsilon) \quad (29)$$

where,

$$B = (1 - \hat{R}_{ij}^2) M(\epsilon) \eta(n_{ij}) + \frac{2}{3\sqrt{3}} M^2(\epsilon) \eta^2(n_{ij})$$

which completes the proof of Lemma 1.

In (20), if we choose $\epsilon = 0.001$, we have $M(\epsilon) \approx 3$ —hence, (29) leads to:

$$\begin{aligned} \Pr \left(|\hat{R}_{ij} - R_{ij}| < \min \{2, 3(1 - \hat{R}_{ij}^2) \eta(n_{ij}) + 2\sqrt{3}\eta^2(n_{ij})\} \mid R_{ij} \right) \\ > (1 - \epsilon) \end{aligned} \quad (30)$$

which proves Corollary 1. Usually this result holds good [?] for any $n_{ij} > 3$. If however $n_{ij} \rightarrow \infty$ for all (i, j) pairs, then the bound on $|\hat{R}_{ij} - R_{ij}|$ in (29) approaches 0 and \hat{R}_{ij} would be close to R_{ij} .

A General Likelihood Framework for Robocov Covariance Matrix Estimation

We propose a generalization of the Robocov covariance matrix estimation framework presented in Section 2.1 – the loss function presented here is directly motivated by the Fisher’s Z-score framework discussed above, but differs from that appearing in Section 2.1.

Recall that the estimators in Section 2.1 are special cases of the following regularized loss minimization framework:

$$\min_{\Sigma \geq 0} \mathcal{L}(\Sigma) + \lambda \xi(\Sigma) \quad (31)$$

where \mathcal{L} is the data fidelity function and ξ is the penalty function, and λ is a tuning parameter that controls the trade-off between data-fidelity and regularization. We can choose $\mathcal{L}(\Sigma; \hat{\Sigma}) = \sum_{ij} \mathcal{L}_{ij}(\Sigma_{ij}, \hat{\Sigma}_{ij})$ with $\mathcal{L}_{ij}(\Sigma_{ij}, \hat{\Sigma}_{ij}) = \max\{|\hat{\Sigma}_{ij} - \Sigma_{ij}| - C_{ij}, 0\}$ for all i, j . This leads to a regularized convex optimization problem of the form:

$$\min \frac{1}{\lambda} \mathcal{L}(\Sigma; \hat{\Sigma}) + \sum_{i < j} |\Sigma_{ij}|. \quad (32)$$

In the limiting case, $\lambda \rightarrow 0+$ i.e., $1/\lambda \rightarrow \infty$, estimator obtained from Problem (32) will reduce to the estimator available from (2). This is because, for sufficiently large values of $1/\lambda$, an optimal solution to (32) will lead to a zero loss— $\mathcal{L}(\Sigma; \hat{\Sigma}) = 0$ which implies that $\mathcal{L}_{ij}(\Sigma_{ij}; \hat{\Sigma}_{ij}) = 0$ for all i, j —these are the data-fidelity constraints in (2). We note that in our numerical experiments, estimator (2) had a performance which was roughly similar to that of the general estimator (32).

A quadratic loss alternative to covariance estimation problem

We present below (See (33)) a convex quadratic loss function $\mathcal{L}(\Sigma)$. While this differs from the loss function considered in (2), in practice, the performances of these two estimators were found to be similar (at least on the datasets we experimented on).

To derive the loss function, we make use of Lemma 2—which presents the (conditional) mean and variance of \hat{R}_{ij} (given R_{ij}). This leads to a loss function of the form:

$$\sum_{ij} \frac{(\hat{R}_{ij} - E(\hat{R}_{ij}|R_{ij}))^2}{\text{var}(\hat{R}_{ij}|R_{ij})}$$

Using the expressions for conditional mean/variances from Lemma 2 (see below), in the above expression, we get:

$$\sum_{ij} \left(\hat{R}_{ij} - (R_{ij} + R_{ij}(1 - R_{ij}^2)\eta^2(n_{ij})) \right)^2 / ((1 - R_{ij}^2)^2\eta^2(n_{ij})).$$

We set $\hat{R}_{ij} = \hat{\Sigma}_{ij}/(\hat{\sigma}_i \hat{\sigma}_j)$ above, and obtain

$$\sum_{ij} \left\{ \frac{(\hat{\sigma}_i \hat{\sigma}_j R_{ij} + \hat{\sigma}_i \hat{\sigma}_j R_{ij}(1 - R_{ij}^2)\eta^2(n_{ij}) - \hat{\Sigma}_{ij})^2}{\hat{\sigma}_i \hat{\sigma}_j(1 - R_{ij}^2)\eta(n_{ij})} \right\}.$$

The loss function above is a highly nonconvex function in R_{ij} or Σ_{ij} . To this end, we approximate the above by replacing some unknown population quantities by their sample analogues. This results in a loss function:

$$\mathcal{L}(\Sigma) = \sum_{ij} \left\{ \frac{(\Sigma_{ij} + \Sigma_{ij}(1 - \hat{R}_{ij}^2)\eta^2(n_{ij}) - \hat{\Sigma}_{ij})^2}{\hat{\sigma}_i \hat{\sigma}_j(1 - \hat{R}_{ij}^2)\eta(n_{ij})} \right\}, \quad (33)$$

which is convex in Σ . In words, $\mathcal{L}(\Sigma)$ above, is a measure of how close Σ_{ij} s are to the pairwise covariance terms $\hat{\Sigma}_{ij}$ s—this critically depends upon the number of observed samples n_{ij} for every pair (i, j) .

We now present Lemma 2 and its proof:

Lemma 2. Assume that all conditions of Lemma 1 hold. If n_{ij} is large so that Cn_{ij}^{-4} is negligible for a constant C , we have:

$$E(\hat{R}_{ij}|R_{ij}) \approx R_{ij} + R_{ij}(1 - R_{ij}^2)\eta^2(n_{ij}) \quad (34)$$

and

$$\text{var}(\hat{R}_{ij}|R_{ij}) \approx (1 - R_{ij}^2)^2\eta^2(n_{ij}) \quad (35)$$

where $\eta(n_{ij})$ is as described in (15).

Proof of Lemma 2

We re-write \hat{R}_{ij} as a function of the Fisher Z-score

$$\hat{R}_{ij} = \frac{\exp(2\hat{Z}_{ij}) - 1}{\exp(2\hat{Z}_{ij}) + 1} \quad (36)$$

We then expand \hat{R}_{ij} as a function of \hat{Z}_{ij} around the population Fisher Z-score Z_{ij} using the 2nd order Taylor series expansion as follows:

$$\begin{aligned} \hat{R}_{ij} &\approx \frac{\exp(2Z_{ij}) - 1}{\exp(2Z_{ij}) + 1} + \frac{4\exp(2Z_{ij})}{\exp(2Z_{ij}) + 1}(\hat{Z}_{ij} - Z_{ij}) + \\ &\quad \frac{4\exp(2Z_{ij})(\exp(2Z_{ij}) - 1)}{(\exp(2Z_{ij}) + 1)^3}(\hat{Z}_{ij} - Z_{ij})^2 \\ &= R_{ij} + (1 - R_{ij}^2)(\hat{Z}_{ij} - Z_{ij}) + R_{ij}(1 - R_{ij}^2)(\hat{Z}_{ij} - Z_{ij})^2 \end{aligned} \quad (37)$$

Using the fact that $E(\hat{Z}_{ij}|R_{ij}) = E(\hat{Z}_{ij}|Z_{ij}) = Z_{ij}$, we get from (37)

$$\begin{aligned} E(\hat{R}_{ij}|R_{ij}) &\approx R_{ij} + R_{ij}(1 - R_{ij}^2)E((\hat{Z}_{ij} - Z_{ij})^2|R_{ij}) \\ &= R_{ij} + R_{ij}(1 - R_{ij}^2)\eta_{ij}^2 \end{aligned} \quad (38)$$

and

$$\begin{aligned} \text{var}(\hat{R}_{ij}|R_{ij}) &\approx (1 - R_{ij}^2)^2\eta^2(n_{ij}) + Cn_{ij}^{-4} \\ &\approx (1 - R_{ij}^2)^2\eta^2(n_{ij}), \end{aligned} \quad (39)$$

where (39) makes use of the fact that Cn_{ij}^{-4} is negligible as per the condition of Lemma 2; and the cross (covariance) term vanishes as it is the third moment of a Gaussian with mean zero.

Derivation of D in (7)

Here we discuss how we derive the analytical form of D in (7) in the optimization framework in (6).

Let $\tilde{\Sigma}$ be the sample covariance matrix of X^f (i.e., the fully observed version of X) We implicitly assume that the perturbation amount Δ is such that $\hat{\Sigma} + \Delta$ is a good approximation to the unobserved $\tilde{\Sigma}$. That is,

$$|\Delta_{ij}| \approx |\hat{\Sigma}_{ij} - \tilde{\Sigma}_{ij}| \leq D_{ij} \quad (40)$$

We can write

$$|\hat{\Sigma}_{ij} - \tilde{\Sigma}_{ij}| \leq |\hat{\Sigma}_{ij} - \Sigma_{ij}| + |\Sigma_{ij} - \tilde{\Sigma}_{ij}|. \quad (41)$$

We propose bounds on each of the two terms on the right using our results from the Robocov covariance matrix section. We know that the first term would be bounded by C_{ij} from Corollary 1. Note that $\tilde{\Sigma}_{ij}$ is an instance of $\hat{\Sigma}_{ij}$ when $n_{ij} = N$ —i.e., all samples are observed. Hence, the bound will be similar to C_{ij} but with n_{ij} replaced by N . We therefore define

$$Q_{ij} := \hat{\sigma}_i \hat{\sigma}_j \min \left(2, \eta(N) \left\{ 3(1 - \hat{R}_{ij}^2) + 2\sqrt{3}\eta(N) \right\} \right) \quad (42)$$

where \hat{R} is the correlation matrix corresponding to $\hat{\Sigma}$.

When N is reasonably large, $|\eta(N)\hat{R}_{ij}^2 - \eta(N)\hat{R}_{ij}^2|$ is very small since both \hat{R}_{ij}^2 and \hat{R}_{ij}^2 are bounded between 0 and 1 and $\eta(N) \rightarrow 0$ as $N \rightarrow \infty$.

Therefore we can effectively replace Q_{ij} by C'_{ij} defined as:

$$C'_{ij} := \hat{\sigma}_i \hat{\sigma}_j \min \left(2, \eta(N) \left\{ 3(1 - \hat{R}_{ij}^2) + 2\sqrt{3}\eta(N) \right\} \right) \quad (43)$$

This provides a justification for the choice of D appearing in (7).

Arriving at the Robocov inverse covariance estimator in Section 2.2

Note that Problem (6) involves minimization of a pointwise maximum (over Δ) of convex functions $\Omega \mapsto L(\Omega; \hat{\Sigma} + \Delta) + \lambda \sum_{ij} |\Omega_{ij}|$. Hence, Problem (6) is convex [?] in Ω .

Here we explain how the min-max optimization problem in (6) leads to the optimization problem in (8).

To this end, note that:

$$\begin{aligned} & \max_{\Delta: |\Delta_{ij}| \leq D_{ij}, \forall i,j} \left\{ -\log \det(\Omega) + \langle \Omega, \hat{\Sigma} + \Delta \rangle \right\} \\ &= \max_{\Delta: |\Delta_{ij}| \leq D_{ij}, \forall i,j} \left\{ -\log \det(\Omega) + \langle \Omega, \hat{\Sigma} \rangle + \langle \Omega, \Delta \rangle \right\} \\ &= -\log \det(\Omega) + \langle \Omega, \hat{\Sigma} \rangle + \max_{\Delta: |\Delta_{ij}| \leq D_{ij}, \forall i,j} \langle \Omega, \Delta \rangle \\ &= -\log \det(\Omega) + \langle \Omega, \hat{\Sigma} \rangle + \sum_{i,j} D_{ij} |\Omega_{ij}| \end{aligned} \quad (44)$$

where, the last line follows by noting that

$$\langle \Omega, \Delta \rangle = \sum_{ij} \Omega_{ij} \Delta_{ij} \leq \sum_{ij} |\Omega_{ij}| \cdot |\Delta_{ij}| \leq \sum_{ij} |\Omega_{ij}| D_{ij}$$

and an equality above holds when $\Delta_{ij} = \text{sign}(\Omega_{ij}) |D_{ij}|$ for all i, j ,

Using (44), Problem (6) becomes:

$$\begin{aligned} & \min_{\Omega \succeq 0} \left\{ \max_{\Delta: |\Delta_{ij}| \leq D_{ij}, \forall i,j} \left\{ -\log \det(\Omega) + \langle \Omega, \hat{\Sigma} + \Delta \rangle \right\} \right. \\ & \quad \left. + \lambda \sum_{ij} |\Omega_{ij}| \right\} \\ &= \min_{\Omega \succeq 0} \left\{ -\log \det(\Omega) + \langle \Omega, \hat{\Sigma} \rangle + \sum_{i,j} D_{ij} |\Omega_{ij}| \right. \\ & \quad \left. + \lambda \sum_{ij} |\Omega_{ij}| \right\} \end{aligned}$$

which is the formulation appearing in (8).

Simulation settings

The parameter models for the simulated population models in Figure 1 are as follows.

- **Hub:** The hub matrix population model for both Figure 1 and Table 1 comprised of correlation blocks of size 5. Each block had all off-diagonal entries equal to 0.7.
- **Toeplitz:** The Toeplitz matrix population model A in Figure 1 had entries of the form $A_{ij} = \max \{0, 1 - 0.1 * |i - j|\}$.
- **1-band precision:** The 1-band precision matrix population model in Figure 1 is of the form $A_{i,i+1} = 0.5$ and $A_{i,j} = 0$ for $j \neq i, i + 1$ for each feature i .

Performance metrics

Three performance metrics were used to compare different correlation and partial correlation estimators for different simulation settings (Table 1). They include

- **FP2 : False Positive 2-norm:** Euclidean distance of the estimated correlation or partial correlation values for feature pairs with population correlation or partial correlation equal to 0.
- **FPR: False Positive Rate:** The proportion of feature pairs with population correlation (partial correlation) equal to 0 that have estimated correlation (partial correlation) greater than 0.1.
- **FNR: False Negative Rate:** The proportion of feature pairs with population correlation (partial correlation) greater than 0.1 that have estimated correlation (partial correlation) less than 0.01.

Stratified LD-score regression

Stratified LD score regression (S-LDSC) is a method that assesses the contribution of a genomic annotation to disease and complex trait heritability[? ?]. S-LDSC assumes that the per-SNP heritability or variance of effect size (of standardized genotype on trait) of each SNP is equal to the linear contribution of each annotation

$$\text{var}(\beta_j) := \sum_c a_{cj} \tau_c, \quad (45)$$

where a_{cj} is the value of annotation c for SNP j , where a_{cj} is binary in our case, and τ_c is the contribution of annotation c to per-SNP heritability conditioned on

other annotations. S-LDSC estimates the τ_c for each annotation using the following equation

$$E \left[\chi_j^2 \right] = N \sum_c l(j, c) \tau_c + 1, \quad (46)$$

where $l(j, c) = \sum_k a_{ck} r_{jk}^2$ is the *stratified LD score* of SNP j with respect to annotation c and r_{jk} is the genotypic correlation between SNPs j and k computed using data from 1000 Genomes Project[?] (see URLs); N is the GWAS sample size.

We assess the informativeness of an annotation c using two metrics. The first metric is enrichment (E_c), defined as follows (for binary and probabilistic annotations only):

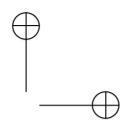
$$E_c = \frac{\frac{h_g^2(c)}{h_g^2}}{\frac{\sum_j a_{cj}}{M}}, \quad (47)$$

where $h_g^2(c)$ is the heritability explained by the SNPs in annotation c , weighted by the annotation values.

The second metric is standardized effect size (τ^*) defined as follows (for binary, probabilistic, and continuous-valued annotations):

$$\tau_c^* = \frac{\tau_c s d_c}{\frac{h_g^2}{M}}, \quad (48)$$

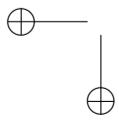
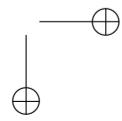
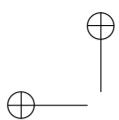
where $s d_c$ is the standard error of annotation c , h_g^2 the total SNP heritability and M is the total number of SNPs on which this heritability is computed (equal to 5, 961, 159 in our analyses). τ_c^* represents the proportionate change in per-SNP heritability associated to a 1 standard deviation increase in the value of the annotation.



a

1

Supplementary Figures



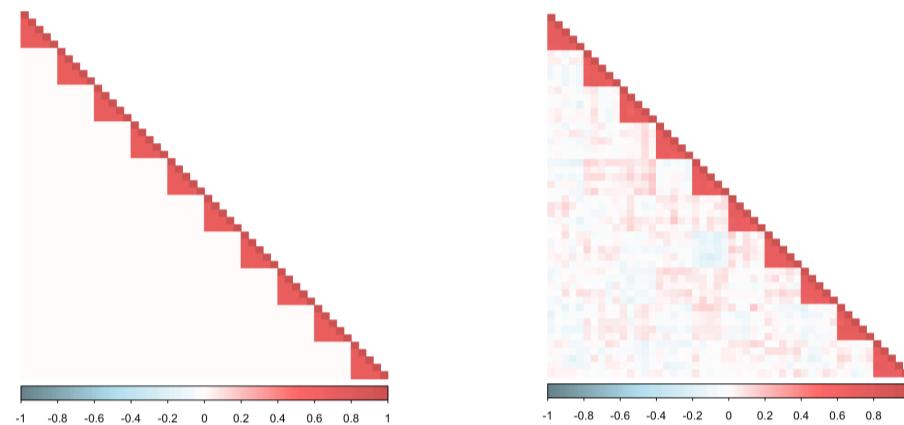
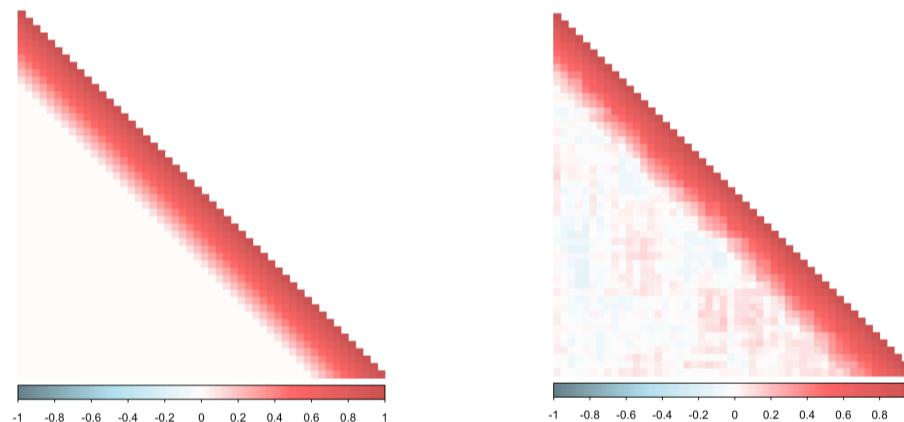
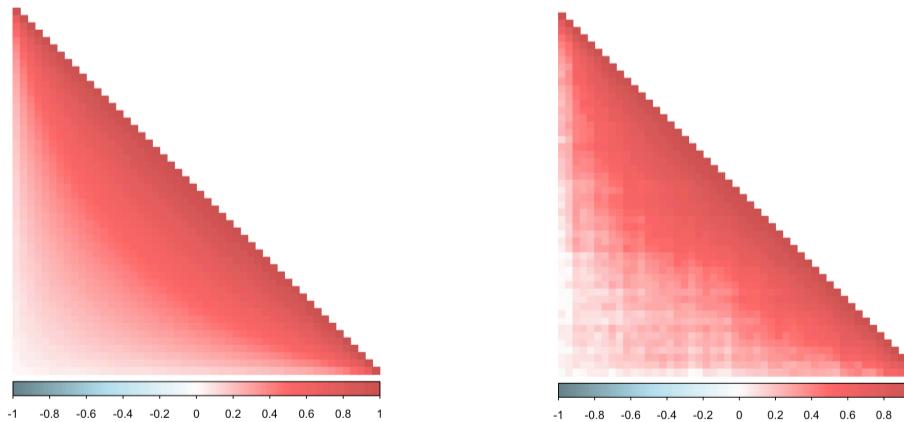
Population correlation**Standard Pairwise Correlation****Hub****Toeplitz****1-band precision**

Fig. S1. We applied standard pairwise correlation estimator on data generated from the simulation models from Figure 1—this comprises of Hub, Toeplitz or 1-band precision matrix-based population models with $N = 500$ samples, $P = 50$ features and $\pi = 50\%$ proportion of missing data.

a

3

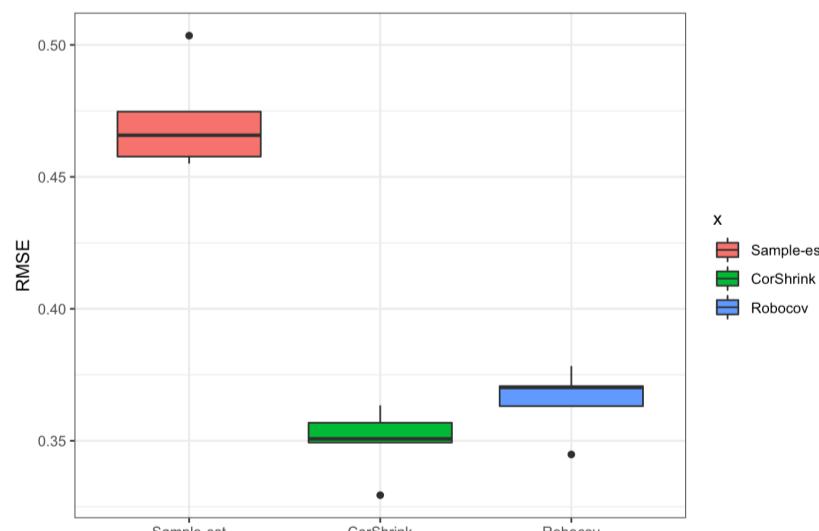
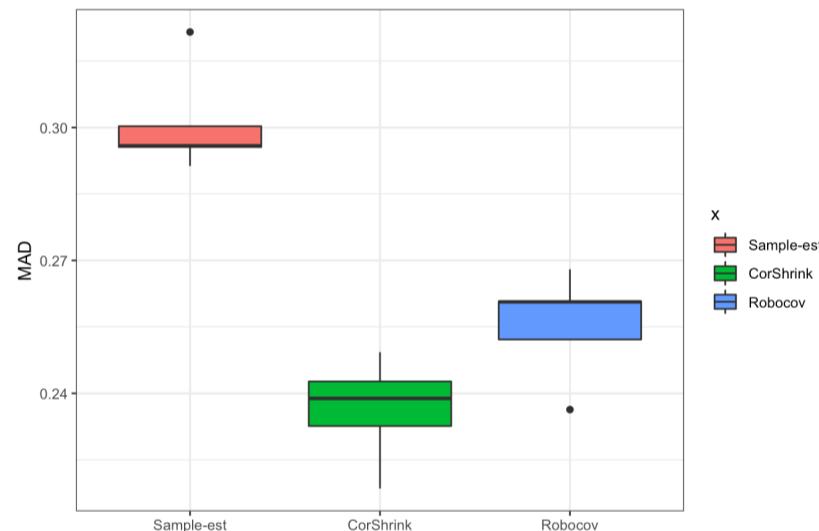
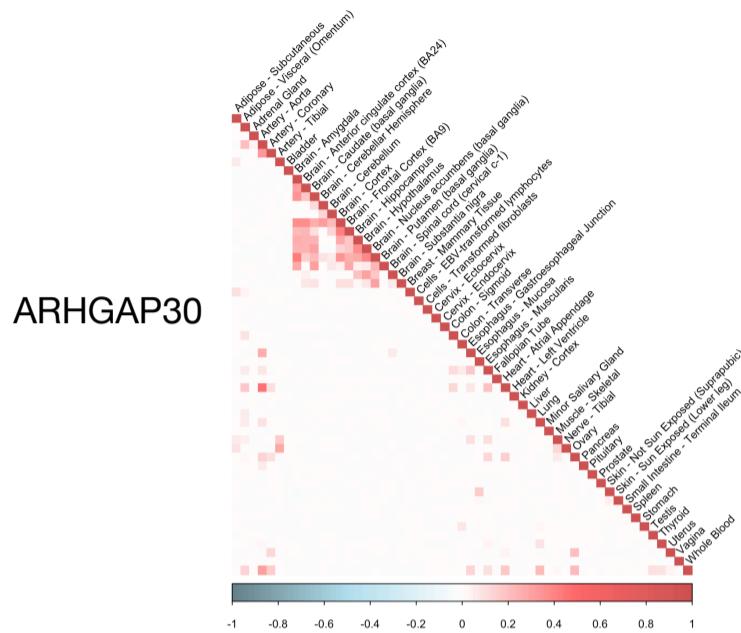


Fig. S2. (Top panel) Robcov correlation estimate of the ARHGAP30 gene. (Middle and bottom panels) We split the data matrix randomly into 2 equal groups. We compare the Robcov, CorShrink and pairwise sample correlation estimators from one half of the data with the pairwise sample correlation matrix on the other half. We use Median Absolute Deviation (MAD) (middle panel) and Root Mean Squared Error (RMSE) (lower panel) metrics. The results are averaged over 50 such random splits. See Table S1 for a numerical summary.

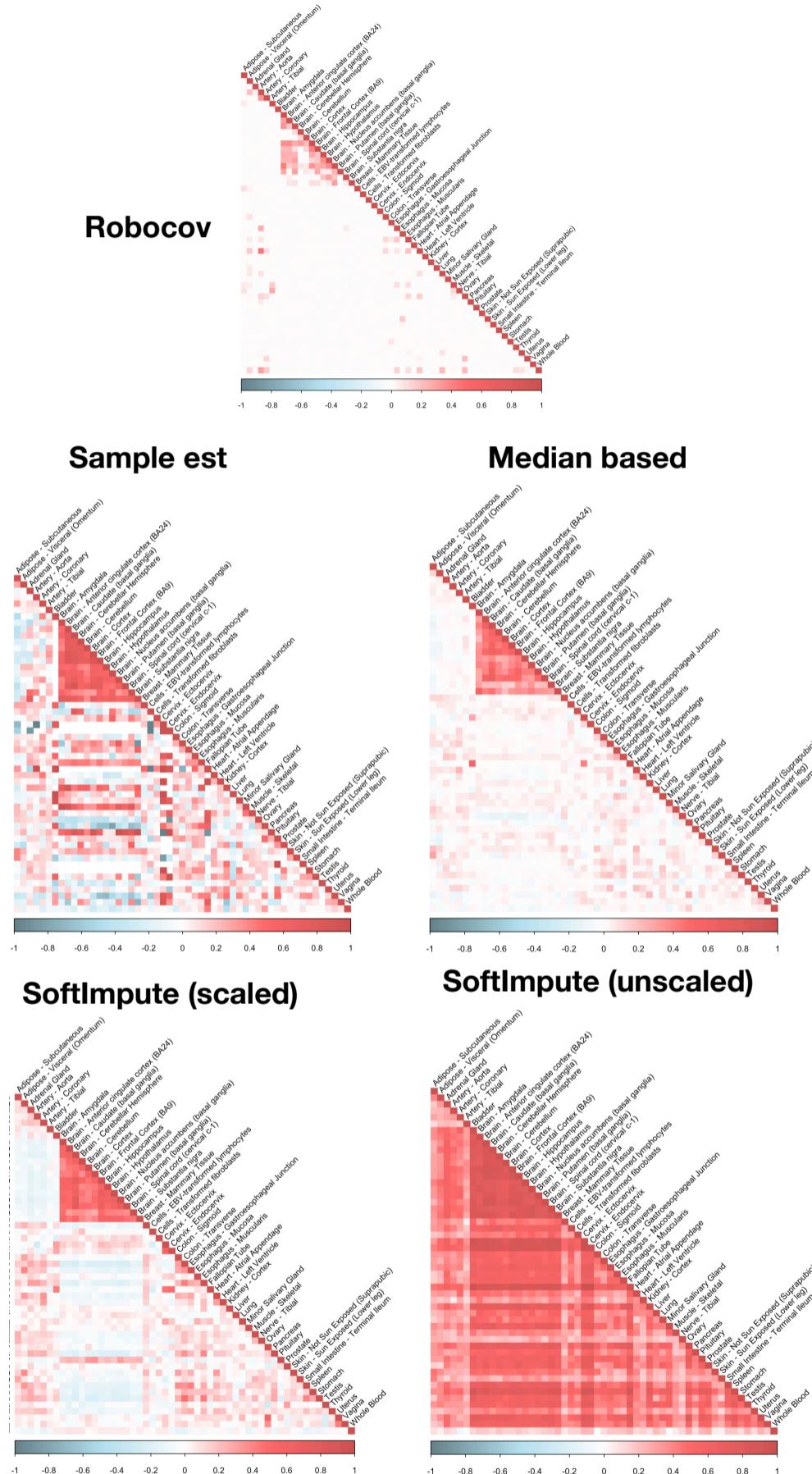


Fig. S3. We compare the Robocov correlation estimator for the ARHGAP30 gene with four other estimators. They include the standard pairwise sample correlation estimator, the sample correlation matrix computed over data imputed by either a median-based approach (missing entries of a feature replaced by the median of observed entries), the scaled SoftImpute[?] approach; and an unscaled SoftImpute[?] approach.

a

5

Eigenvalue trend for ARHGAP30 gene

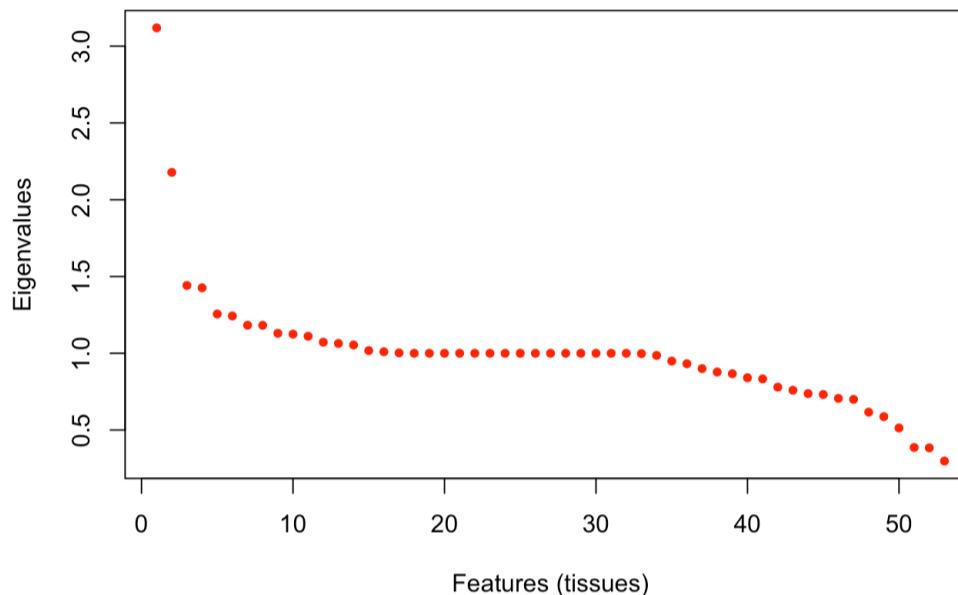


Fig. S4. Plot of eigenvalues sorted from highest to lowest in magnitude for tissue-tissue pairwise correlation matrix for a particular gene (ARHGAP30). The eigenvalues do not show any sharp drop close to 0 as one would expect if the matrix allowed a low rank (+noise) structure. This suggests relatively high dimensional structure in the GTEx gene expression data which may explain why a low rank imputation method such as SoftImpute[?] performs poorly in S3.

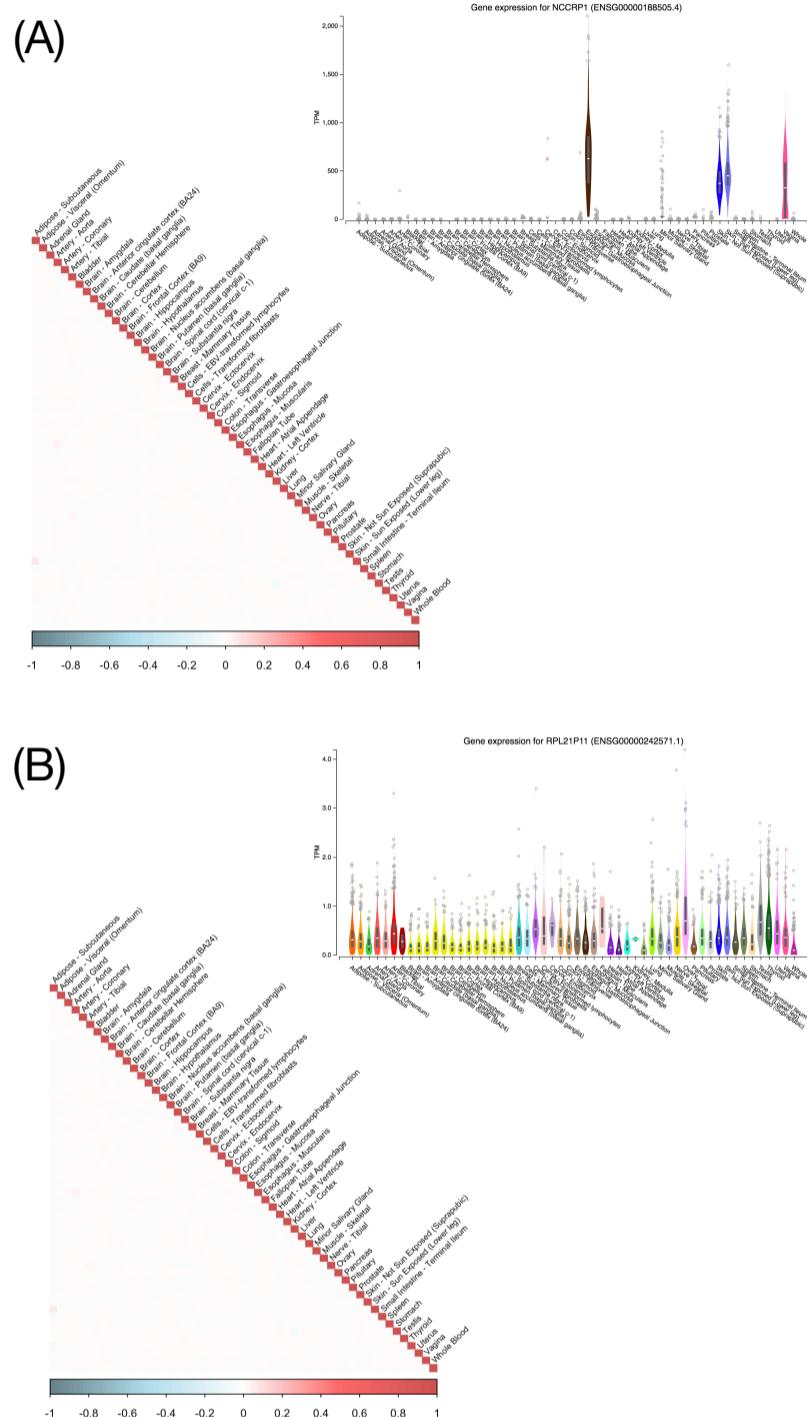


Fig. S5. Examples of two genes, NCCRP1 (top) and RPL21P11 (bottom), both of which have close to 0 average correlation in expression across tissue-pairs but having very distinctive expression profiles. NCCRP1 has high expression in a few specific tissues including Whole Blood, while RPL21P11 has uniformly low expression across all tissues. The expression profile plots for the genes have been fetched from the GTEx Portal (<https://gtexportal.org/home/>).

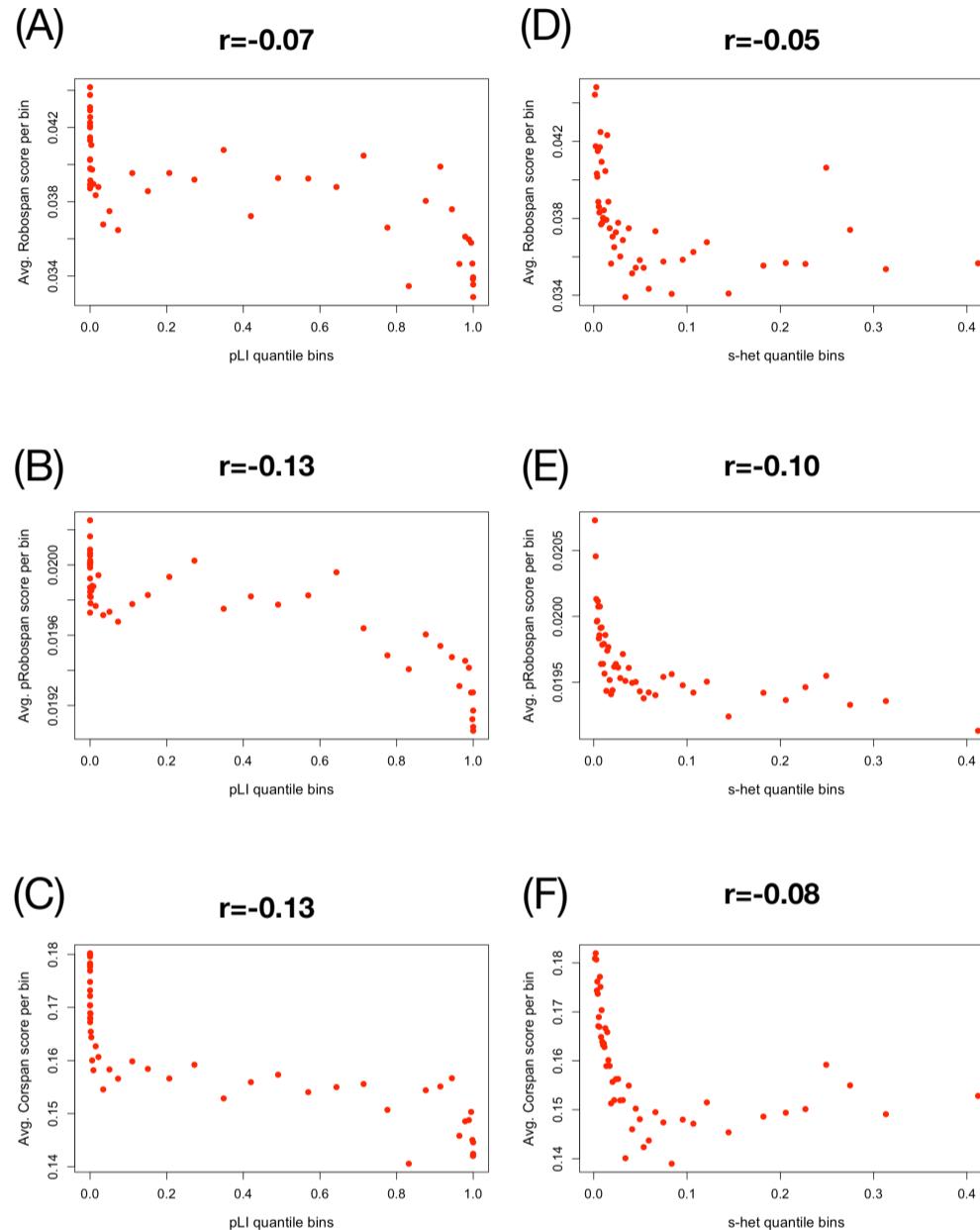


Fig. S6. Comparison of pLI gene score with (A) Robospan-score, (B) pRobospan-score and (C) Corspan-score for all genes (See Results section for details). Comparison of s_het gene score with (A) Robospan-score, (B) pRobospan-score and (C) Corspan-score for all genes (See Results section for details).

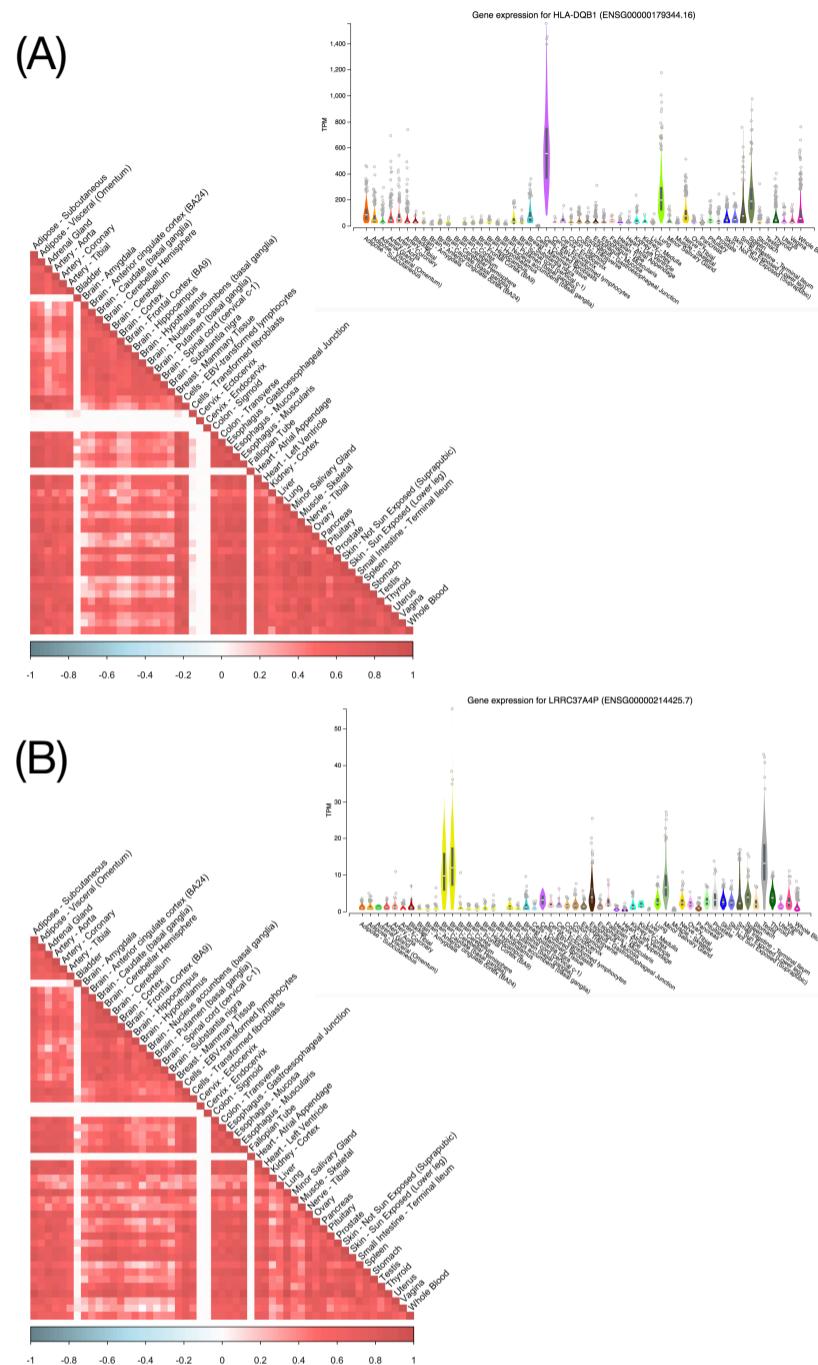


Fig. S7. Examples of genes high Robospan-score but not specifically expressed in blood or uniformly expressed across tissues. The two genes are HLA-DQB1 (top) and LRRC37A4P (bottom). HLA-DQB1 is specifically expressed in lymphocyte cell line which is related to blood. LRRC37A4P has highest expression in brain cerebellum and testis. The expression profile plots for the genes have been fetched from the GTEx Portal (<https://gtexportal.org/home/>).

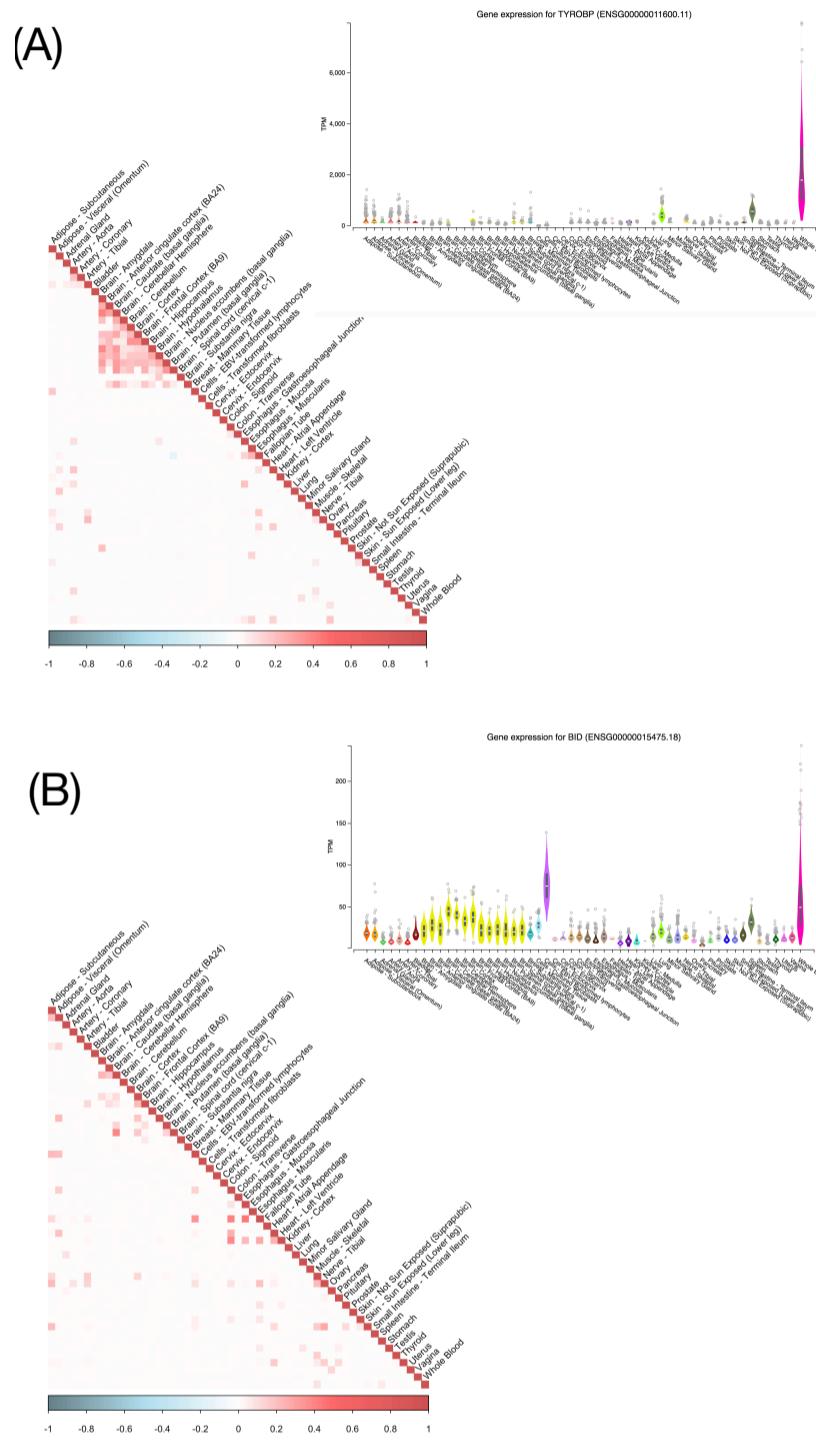


Fig. S8. Examples of genes that are specifically expressed in Whole Blood but do not show high Robospan-score. The two genes are TYROBP (top) and BID (bottom). They show tissue specific expression in Whole blood and are in top 10% specifically expressed genes (SEG) as per ref.[?], but do not show consistently high expression correlation across tissue pairs and hence, do not have high Robospan score. The expression profile plots for the genes have been fetched from the GTEx Portal (<https://gtexportal.org/home/>).

Supplementary Tables

Table S1. Predictive comparison of CorShrink, Robocov and sample correlation estimators for a GTEx gene. We report Mean Absolute Deviation (MAD) and Root Mean Squared Deviation (RMSD) metrics between an estimator (e.g., sample correlation matrix, CorShrink and Robocov) computed on the gene expression data (GTEx project) for half of the individuals (training set) and the sample correlation matrix computed from other half (testing set) of all individuals. Results are averaged over 30 such different training/testing data-splits with the standard errors reported in brackets.

Method	MAD	RMSE
Sample-Est	0.30 (0.01)	0.47 (0.02)
CorShrink	0.24 (0.01)	0.35 (0.01)
Robocov	0.25 (0.01)	0.36 (0.01)

Table S2. Pathway enrichment analysis of Robospan, pRobospan and Corspan genes. Pathway enrichment is performed using the ConsensusPathDB database[?]. Only the top 5 non-redundant and statistically significant ($q\text{-value} < 0.05$) pathways for a gene set are reported.

Gene Set	Top pathways
Robospan	Interferon signaling (1.1e-18), Immune system (3.1e-08), HSF1 activation (1.1e-07), Antigen processing and presentation (2.4e-07), Allograft rejection (1.5e-06)
pRobospan	Immune system (2.7e-21), Interferon signaling (3.4e-15), Innate immune system (5.1e-12), TNF signaling pathway (7.9e-11), Neutrophil degranulation (2.0e-10)
Corspan	Interferon signaling (1.1e-17), Immune system (1.6e-07), Antigen processing and presentation (1.2e-05), HSF1 activation (4.1e-05), Neutrophil degranulation (1e-04),

Table S3. List of 11 blood and autoimmune traits (5 blood traits and 6 autoimmune traits) analyzed in this paper.

Annotation	Traits
Blood traits	Red blood Cell Distribution Width (UKBB[?]), Red blood Cell Count (UKBB[?]), White blood Cell Count (UKBB[?]), Platelet Count (UKBB[?]), Eosinophil Count (UKBB[?])
Immune traits	Ulcerative Colitis[?], Rheumatoid Arthritis[?], Celiac[?], Lupus[?], Crohn’s disease[?], Auto Immune and Inflammatory traits

Table S4. S-LDSC results for SNP annotations corresponding to Robospan, pRobospan, Corspan and SEG-Blood gene sets for blood and autoimmune traits. Standardized Effect sizes (τ^*) and Enrichment (E) of 8 SNP annotations corresponding to 4 gene sets (Robospan, pRobospan, Corspan and SEG-Blood[?]) and 2 SNP annotations corresponding to 5kb and 100kb window based SNP-to-gene linking strategies for each gene set. Results for all annotations are conditional on 86 baselineLD-v2.1 annotations. Reports are meta-analyzed across 11 Blood and Autoimmune traits.

Robospan						
	τ^*	se(τ^*)	p(τ^*)	E	se(E)	p(E)
5kb (2.6%)	0.086	0.024	0.00048	2.7	0.16	1.5e-07
100kb (10%)	0.12	0.03	7.9e-05	2.3	0.12	2e-09
pRobospan						
	τ^*	se(τ^*)	p(τ^*)	E	se(E)	p(E)
5kb (2.3%)	0.096	0.028	0.00057	3.2	0.22	9.3e-08
100kb (9.9%)	0.11	0.034	0.0011	2.4	0.12	5.5e-09
Corspan						
	τ^*	se(τ^*)	p(τ^*)	E	se(E)	p(E)
5kb (2.5%)	0.04	0.024	0.093	2.4	0.15	4.8e-07
100kb (9.8%)	0.038	0.02	0.059	2.1	0.1	1.7e-08
SEG-Blood						
	τ^*	se(τ^*)	p(τ^*)	E	se(E)	p(E)
5kb (2.7%)	0.24	0.036	7.6e-11	3.6	0.26	8.7e-10
100kb (10.1%)	0.21	0.029	1.3e-13	2.4	0.095	2.2e-10

Table S5. S-LDSC results for SNP annotations corresponding to Robospan, pRobospan, Corspan and SEG-Blood gene sets for 6 autoimmune traits. Standardized Effect sizes (τ^*) and Enrichment (E) of 8 SNP annotations corresponding to 4 gene sets (Robospan, pRobospan, Corspan and SEG-Blood[?]) and 2 SNP annotations corresponding to 5kb and 100kb window based SNP-to-gene linking strategies for each gene set. Results for all annotations are conditional on 86 baselineLD-v2.1 annotations. Reports are meta-analyzed across 6 Autoimmune traits.

Robospan						
	τ^*	se(τ^*)	p(τ^*)	E	se(E)	p(E)
5kb (2.6%)	0.12	0.036	7e-04	2.7	0.25	5e-05
100kb (10%)	0.14	0.049	0.0051	2.3	0.19	2e-06
pRobospan						
	τ^*	se(τ^*)	p(τ^*)	E	se(E)	p(E)
5kb (2.3%)	0.1	0.043	0.016	3.3	0.39	1e-04
100kb (9.9%)	0.11	0.059	0.061	2.5	0.24	1e-05
Corspan						
	τ^*	se(τ^*)	p(τ^*)	E	se(E)	p(E)
5kb (2.5%)	0.08	0.035	0.021	2.5	0.23	1e-04
100kb (9.8%)	0.037	0.035	0.28	2	0.15	8e-06
SEG-Blood						
	τ^*	se(τ^*)	p(τ^*)	E	se(E)	p(E)
5kb (2.7%)	0.33	0.042	5e-15	4.2	0.25	8e-06
100kb (10.1%)	0.3	0.036	9e-17	2.7	0.11	7e-06

Table S6. Joint S-LDSC results for annotations corresponding to Robospan, pRobospan, Corspan and SEG-Blood gene sets.: Standardized Effect sizes (τ^*) and Enrichment (E) of SNP annotations that are significant when all annotations from Table S4 are modeled jointly and subjected to forward stepwise elimination. 2 annotations survive in the resulting joint model. The analysis is conditional on 86 baselineLD-v2.1 annotations. Reports are meta-analyzed across 11 Blood and Autoimmune traits.

	τ^*	se(τ^*)	p(τ^*)	E	se(E)	p(E)
Robospan(100kb)	0.1	0.03	1e-04	2.3	0.1	2e-09
SEG-Blood(100kb)	0.2	0.03	2e-13	2.4	0.09	2e-10