

RESEARCH

Open Access



# Tisslet tissues-based learning estimation for transcriptomics

Ahmed Miloudi<sup>1†</sup>, Aisha Al-Qahtani<sup>2\*†</sup>, Thamanna Hashir<sup>3</sup>, Mohamed Chikri<sup>1\*</sup> and Halima Bensmail<sup>2\*</sup>

<sup>†</sup>Ahmed Miloudi and Aisha Al-Qahtani These authors equally contributed to this work.

\*Correspondence:  
aialqahtani@hbku.edu.qa;  
mohamed.chikri@usmba.ac.ma;  
hbensmail@hbku.edu.qa

<sup>1</sup> Faculty of Medicine  
and Pharmacy-FUSMBA, Fes,  
Morocco

<sup>2</sup> Qatar Computing Research  
Institute, HBKU, Doha, Qatar

<sup>3</sup> Carnegie Mellon University-  
Qatar, Doha, Qatar

## Abstract

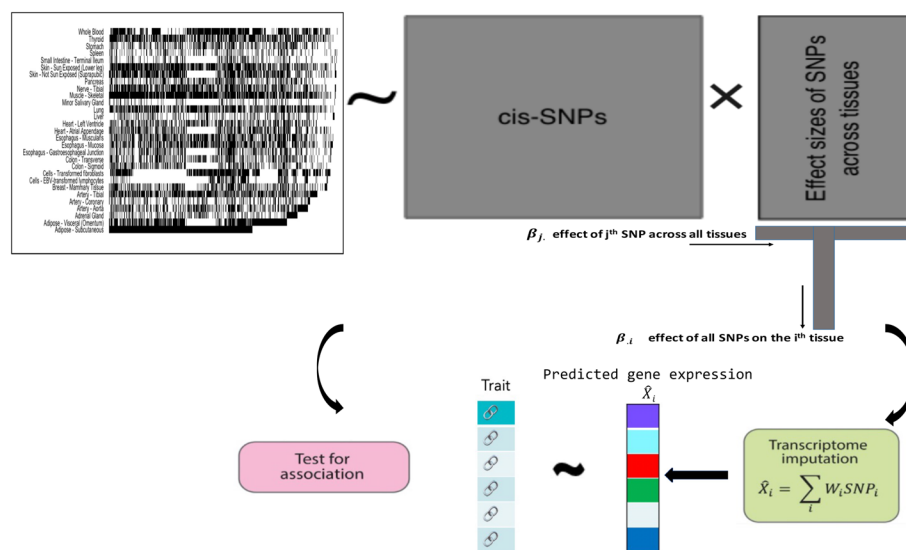
In the context of multi-omics data analytics for various diseases, transcriptome-wide association studies leveraging genetically predicted gene expression hold promise for identifying novel regions linked to complex traits. However, existing methods for multi-tissue gene expression prediction often fail to account for tissue-tissue expression interactions, limiting their accuracy and effectiveness. This research addresses the challenge of predicting gene expression across multiple tissues by incorporating tissue-tissue expression correlations based on a nonlinear multivariate model. Our findings demonstrate that this model excels in estimating tissue-tissue interactions and accurately predicting missing data. These results have significant implications for multi-omics data analytics and transcriptome-wide association studies, suggesting a novel approach for identifying regions associated with complex traits.

**Keywords:** Transcriptomics, Multiple-tissues, Machine learning, Sparse covariance matrix, EQTL, Likelihood estimator

## Introduction

eQTL (expression Quantitative Trait Loci) studies explore the relationship between genetic variants, such as SNPs (Single Nucleotide Polymorphisms), and gene expression levels. By identifying how specific SNPs influence gene expression, researchers can gain insights into the genetic basis of complex traits and diseases. In the context of transcriptomics, which involves the comprehensive analysis of RNA transcripts, eQTL mapping is particularly valuable (Fig. 1). When considering multiple tissues, incorporating tissue-tissue correlations becomes essential, as gene expression can be influenced by interactions across different tissues. This approach allows for a more holistic understanding of gene regulation and the identification of novel regions associated with complex traits, enhancing the precision and effectiveness of transcriptome-wide association studies. Transcriptome-wide association studies (TWAS) have therefore gained prominence in the field of genomics and genetics for their potential to uncover genetic variants associated with complex traits. These studies leverage gene expression data to bridge the gap between genetic variations and phenotypic traits. Although TWAS holds great promise, the accuracy of gene expression





**Fig. 1** Overview of TISSLET method. The inputs required for TISSLET is a matrix of gene expression for several tissues. We also use the subjects genotypes with sample matched measured expression. TISSLET's weights output are calculated based on the CEM algorithm using measured gene expression and provide the weights and covariance structure of tissues. For full details, see the Materials and Methods section

prediction across multiple tissues remains a challenge. Previous research addressed the challenge of multi-tissue gene expression (Grinberg and Wallace, 2021) [1]

Machine learning algorithms have been used to predict eQTL regulation of a gene by analyzing adjacent genetic variants [2, 3]. Researchers have used genetically predicted gene expression models to conduct transcriptome-wide association studies (TWAS) and identify novel areas linked to complex traits. However, many of these regions lack a GWAS relationship within 1Mb.

There are various benefits to such analyses: Leveraging gene expression enriches possible trait-associated SNPs, whereas joint eQTL modeling improves overall association strength and reduces the number of tests from millions to roughly 20,000 genes.

Leveraging shared eQTLs across tissues enhances eQTL discovery and gene expression imputation accuracy, leading to more powerful transcriptome-wide association analyses [4–6]. Hu et al. [7] and Molstad et al. [8] developed a penalized regression strategy for collaborative modeling of eQTLs. The penalty encourages shared eQTLs between tissues. Using genotype and expression data from the Genotype-Tissue Expression (GTEx) project, multi-tissue eQTL models significantly increase gene association identification and imputation accuracy compared to single-tissue techniques.

While Hu et al.'s technique [7] does not account for tissue-specific gene expression correlation in combined eQTL modeling, Molstad et al. [8] uses cross-tissue imputation with EM algorithm. Both approach uses linear model to derive their algorithm. Recent research suggests that accounting for metric based on skewness improves variable selection and prediction accuracy and can identify regulators or genes in large patient cohorts. A recent study showed that significant correlation was detected between expression skewness and the top 500 genes corresponding to the most significant differential DNA methylation occurring in the promotor regions in TCGA [9]

Recent research suggests that accounting for tissue-tissue correlation in high-dimensional penalized multivariate response linear regression improves variable selection and prediction accuracy. This phenomena may be explained by a seemingly unrelated regression interpretation of high-dimensional sparse multivariate response linear regression. Moreover, certain tissue types are more difficult to obtain due to biological and cost constraints. Using tissue-tissue correlation with skewness assumption enhances gene expression prediction accuracy, particularly for small data numbers.

In this paper, we propose an approach (TISSLET) which not only impute missing gene expression using cross information from several tissue, but also estimate tissue-tissue correlation. We calculate a joint eQTL weights while imputing missing gene expression using a skewed normal modeling. The technique by which our approach works is straightforward. Measuring expression in one tissue can provide a reliable estimate of expression in the other, especially if their expression is substantially correlated. Ignoring tissue-tissue correlation can significantly reduce gene expression prediction accuracy. Our methodology offers several advantages over previous approaches for multi-tissue joint eQTL mapping, including:

1. Incorporation of a full tissue-tissue correlation matrix in the model, rather than assuming a diagonal matrix, which reveals cross-tissue expression dependencies that eQTLs alone cannot explain.
2. Efficient estimation of eQTL weights by modeling cross-tissue associations.
3. Relaxation of the normality assumption by allowing for a heavy-tailed error distribution, assuming that errors follow a multivariate skewed distribution.

Figure 1 gives an overview of TISSLET method. This figure demonstrates the training of the imputation model. The inputs required for TISSLET is a matrix of gene expression for several tissues. We also use the subjects genotypes with sample matched measured expression. TISSLET's weights output are calculated based on the CEM algorithm using measured gene expression and provide the weights and covariance structure of tissues. For full details, see the Materials and Methods section.

## Materials and methods.

Regression models are commonly used to map the relationship between SNPs and gene expression levels in eQTL studies. Traditional models often assume normality of errors, which may not capture the true distribution of gene expression data. Introducing a skewed model with a cross-tissue expression based on SNP genotypes can provide a more accurate representation by allowing for heavy-tailed error distributions, improving the precision of eQTL mapping and the identification of genetic influences on gene expression.

Let  $\mathbf{x}_i \in \mathbb{R}^p$  represent the genotypes of  $p$  SNPs (both centered and normalized) and let  $\mathbf{y}_i \in \mathbb{R}^q$  represent the vector of centered and normalized measured expression in  $q$  tissues for the  $i^{th}$  subject for a specific gene within a certain distance (e.g. 500 kb) or less away from the gene of interest. We assume that gene expression represents a realization of the random vector for the  $i^{th}$  subject:

$$\mathbf{y}_i \sim \boldsymbol{\beta}^t \mathbf{x}_i + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\epsilon}_i \sim \text{SN}(\mathbf{0}, \boldsymbol{\Sigma}, \boldsymbol{\Lambda}). \quad (2.1)$$

where  $\text{SN}_q$  denotes the  $q$ -dimensional multivariate skewed normal distribution, where  $\boldsymbol{\beta} \in \mathbb{R}^{p \times q}$  is the unknown regression coefficient matrix (i.e., eQTL weights), and  $\boldsymbol{\Sigma}^{-1} \in \mathbb{R}^{q \times q}$  is the cross-tissue error precision (inverse covariance) matrix. We further assume that  $\boldsymbol{\epsilon}_i$  is independent of  $\boldsymbol{\epsilon}_j$  for all  $i \neq j$ .

The skewed distribution  $\text{SN}$  can reveal the asymmetric information when the observations such as gene expression are skewed [10]. Asymmetrical distributions have an additional shape parameter  $\boldsymbol{\Lambda} \in \mathbb{R}^{q \times q}$  to represent the direction of the asymmetry of the density. If the skewness in observations is ignored, inferences with symmetric distributions may result in biased or even misleading conclusions.

### Penalized skew-normal log-likelihood

In appendix (7.2) we outline in details the derivation of the model parameters (weights and cross-correlation matrices) under the normality assumption. This approach uses the ECM algorithm for parameter estimation. Algorithm 1 gives the pseudocode of the approach using normality assumption.

**Algorithm 1** Regularized ECM Algorithm with normality

---

```

1: Initialize  $\boldsymbol{\beta}^{(1)} \in \mathbb{R}^{p \times q}$  and  $\boldsymbol{\Omega}^{(1)} \in \mathbb{R}^{q \times q}$ .
2: Set  $k = 1$ 
3: while not converged do
4:    $\boldsymbol{\Omega}^{(k+1)} \leftarrow \arg \min_{\boldsymbol{\Omega} \in \mathbb{R}^{q \times q}} \left[ \frac{1}{n} \sum_{i=1}^n \text{tr} \{ \mathbf{S}_i(\boldsymbol{\beta}^{(k)}, \boldsymbol{\Sigma}^{(k)}) \boldsymbol{\Omega} \} - \log |\boldsymbol{\Omega}| + \mathcal{R}(\boldsymbol{\Omega}) \right]$ 
5:    $\mathbf{A} \leftarrow \left( \tilde{\mathbf{y}}_i^{(k)} - \boldsymbol{\beta}^T \mathbf{x}_i \right)^T \left( \tilde{\mathbf{y}}_i^{(k)} - \boldsymbol{\beta}^T \mathbf{x}_i \right)$ 
6:    $\boldsymbol{\beta}^{(k+1)} \leftarrow \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{p \times q}} \left[ \frac{1}{n} \sum_{i=1}^n \mathbf{A}^T \boldsymbol{\Omega}^{(k+1)} \mathbf{A} + \mathcal{R}(\boldsymbol{\beta}) \right]$ 
7:   Update  $\mu_i^{(k+1)}, V_i^{(k+1)}$  for  $i \in [n]$  as in (7.1) to (7.5) in the appendix
8:   Construct  $\mathbf{S}(\boldsymbol{\beta}^{(k+1)}, \boldsymbol{\Sigma}^{(k+1)})$ 
9:   Set  $k \leftarrow k + 1$ 
10: end while

```

---

Similarly, using equation (2.1) and assuming that the errors  $\boldsymbol{\epsilon}$ 's  $\stackrel{\text{i.i.d.}}{\sim} \text{SN}(\mathbf{0}, \boldsymbol{\Sigma}, \boldsymbol{\Lambda})$ , for  $1 \leq i \leq n$ , and assuming that gene expression represents a realization of the random vector for the  $i^{\text{th}}$  subject:

$$(\mathbf{y}_i | \boldsymbol{\beta}, \boldsymbol{\Lambda}, \boldsymbol{\Sigma}) \sim \mathbb{N}_p(\boldsymbol{\beta}^t \mathbf{x}_i + \boldsymbol{\Lambda} \boldsymbol{\tau}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\tau} \sim \text{TN}_p(\mathbf{0}, \mathbb{I}_p)$$

where  $\mathbf{x}_i = (\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,p})$  is the vector of covariates (here  $\mathbf{x}_i$  is the  $i^{\text{th}}$  individus genotype),  $\boldsymbol{\beta}^t = (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_p)^t$ , is an unknown matrix of mean regression coefficient and  $\text{TN}$  is a truncated normal distribution.

On the other hand, since gene expression matrix has random missing values, we consider that  $\mathbf{y}_i$ 's are partially observed with an arbitrary missing pattern. In order to set up estimating equations for multivariate data with possible missing values, we separate  $\mathbf{y}_i (q \times 1)$  into two components  $(\mathbf{y}_i^o, \mathbf{y}_i^m)$  accordingly, where  $\mathbf{y}_i^o (q_i^o \times 1)$  is the observed component and  $\mathbf{y}_i^m ((q - q_i^o) \times 1)$  is the missing component. Further, we introduce two missingness indicator matrices, denoted by  $\mathbf{O}_i$  and  $\mathbf{M}_i$  henceforth, corresponding to  $\mathbf{y}_i$  such that  $\mathbf{y}_i^o = \mathbf{O}_i \mathbf{y}_i$  and  $\mathbf{y}_i^m = \mathbf{M}_i \mathbf{y}_i$ , respectively. More specifically,  $\mathbf{O}_i (q_i^o \times q)$  and  $\mathbf{M}_i ((q - q_i^o) \times q)$  are sub-matrices extracted from the rows of an identity matrix of order  $q$ ,  $\mathbf{I}_q$ , corresponding to row positions of  $\mathbf{y}_i^o$  and  $\mathbf{y}_i^m$  in  $\mathbf{y}_i$ , respectively. When  $\mathbf{y}_i$  is fully observed,  $\mathbf{O}_i = \mathbf{I}_q$  and  $\mathbf{M}_i$  is null. Meanwhile, it is easy to verify that:  $\mathbf{y}_i = \mathbf{O}_i^T \mathbf{y}_i^o + \mathbf{M}_i^T \mathbf{y}_i^m$  and  $\mathbf{O}_i^T \mathbf{O}_i + \mathbf{M}_i^T \mathbf{M}_i = \mathbf{I}_q$ . Using appendix (7.3) and Bayes Theorem the negative log-likelihood for the observations  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n$

$$\begin{aligned} \mathcal{L}(\Theta | \mathbf{y}^o, \mathbf{y}^m, \tau) &= \log p(\mathbf{y}^m | \mathbf{y}^o, \tau) p(\tau | \mathbf{y}^o) = \\ &= \frac{1}{2} \sum_{i=1}^n \{ \log |\mathbf{\Omega}| + (\mathbf{y}_i - \beta^t \mathbf{x}_i - \mathbf{\Lambda} \tau) \mathbf{\Omega} (\mathbf{y}_i - \beta^t \mathbf{x}_i - \mathbf{\Lambda} \tau) + \tau_i \tau_i^t \} \end{aligned} \quad (2.2)$$

Then the expected conditional log likelihood  $\mathbf{Q}(\beta, \mathbf{\Omega}, \mathbf{\Lambda})$  is expressed as

$$\mathbb{E}[\mathcal{L}(\beta, \mathbf{\Omega}, \mathbf{\Lambda} | \mathbf{y}^o, \mathbf{y}^m, \tau)] \propto \frac{1}{n} \sum_{i=1}^n \text{tr}\{\mathbf{R}_i(\beta, \mathbf{\Sigma}) \mathbf{\Omega}\} - \log |\mathbf{\Omega}|$$

where

$$\mathbf{R}_i(\beta, \mathbf{\Sigma}, \mathbf{\Lambda}) = \mathbb{E} \left[ (\mathbf{y}_i - \xi - \mathbf{\Lambda} \tau_i) (\mathbf{y}_i - \xi - \mathbf{\Lambda} \tau_i)^t | \mathbf{y}_i^o, \Theta \right] \quad (2.3)$$

### Regularization for precision matrix $\mathbf{\Omega}$

We regularize the entries of  $\mathbf{\Omega}$  using  $\ell_1$  penalties, then the regularized conditional log likelihood  $\mathbf{Q}$  is expressed as:

$$\mathbf{Q}(\beta, \mathbf{\Omega}, \mathbf{\Lambda}) + \mathcal{R}(\mathbf{\Omega}) = \mathbf{Q}(\beta, \mathbf{\Omega}, \mathbf{\Lambda}) + \lambda_1 \sum_{j=1}^q \sum_{k=1}^p |\mathbf{\Omega}_{j,k}|$$

For sufficiently large tuning parameter values  $\lambda_1$ , the penalty results in precision matrix estimations  $\mathbf{\Omega}$  with all off-diagonal entries equal to zero. The penalty assumes that some entries are equal to zero. In multi-tissue joint eQTL mapping, a zero in the  $(j,k)^{th}$  entry of  $\mathbf{\Omega}$  indicates independent expression in the  $j^{th}$  and  $k^{th}$  tissues, given expression in all other tissues and  $p$  SNP genotypes (1). We do this for two reasons: first, when  $p > n$  (which is the case for practically every gene we test, more SNPs than observations), without penalizing the diagonals, a perfect fit can occur in the M-step and not in E step (see Algorithm 2): Recent literature, has shown that precision matrix estimators employed in predictive models exhibit this behavior [11]. Next we summarize the main ECM Algorithm 2. Algorithm 3 gives a detailed version of Algorithm in appendix 7.3):

**Algorithm 2** Regularized ECM Algorithm with non-normality

---

```

1: Initialize  $\beta^{(1)} \in \mathbb{R}^{p \times q}$  and  $\Omega^{(1)} \in \mathbb{R}^{q \times q}$ .
2: Set  $k = 1$ .
3: while not converged do
4:   Compute  $\mathbf{Q}(\beta, \Omega, \Lambda | \beta^{(k)}, \Omega^{(k)}, \Lambda^{(k)})$ 
5:    $\Omega^{(k+1)} = \arg \min_{\Omega \in \mathbb{R}^{q \times q}} \left[ \text{tr}(\mathbf{R}^{(k)} \Omega) - \log |\Omega| + \mathcal{R}(\Omega) \right]$ 
6:    $\xi^{(k+1)} = \arg \min_{\xi \in \mathbb{R}^{p \times q}} \left[ \text{tr}(\mathbf{R}^{(k)} \Omega) \right]$ 
7:    $\hat{\beta}^{(k+1)} = \hat{\xi}^{(k+1)} \mathbf{x}_i / (\mathbf{x}_i \mathbf{x}_i^t)$ 
8:   If the previous objective value Q has not converged, update  $k = k + 1$  and
      return to Step 3.
9: end while

```

---

**Remark 1:** Step 4 of Algorithm 2 can be expressed

$$\Omega^{(k+1)} = \arg \min_{\Omega \in \mathbb{R}^{q \times q}} \left[ \text{tr}(\mathbf{R}^{(k)} \Omega) - \log |\Omega| + \lambda_1 \sum_{j=1}^q \sum_{k=1}^p |\Omega_{j,k}| \right] \quad (2.4)$$

Conveniently, (2.4) is exactly the optimization problem for computing the  $\ell_1$ -penalized normal log-likelihood precision matrix estimator with input sample covariance matrix  $\mathbf{R}(\cdot)$ . Many efficient algorithms and software packages exist for computing (2.4). We can use our R software BiGQUIC to solve (2.4) which is available at [BiGQUIC](#) [12]

**Remark 2.** Step 5 of Algorithm 2 can be expressed as:

$$\xi^{(k+1)} = \arg \min_{\xi \in \mathbb{R}^{p \times q}} \frac{1}{n} \sum_{i=1}^n \left\{ (\mathbf{y}_i^{(k)} - \xi - \Lambda \tau_i) \Omega^{(k+1)} (\mathbf{y}_i^{(k)} - \xi - \Lambda \tau_i)^t \right\}$$

Using appendix (7.3), we have

$$\hat{\xi}^{(k+1)} = \frac{1}{n} \left( \sum_{i=1}^n \hat{\mathbf{y}}_i^{(k)} - \hat{\Lambda}^{(k)} \sum_{i=1}^n \hat{\eta}_i^{(k)} \right)$$

where

$$\hat{\mathbf{y}}_i = \Omega^{-1} \mathbf{S}_i^{oo} \mathbf{y}_i + \left( \mathbf{I} - \Omega^{-1} \mathbf{S}_i^{oo} \right) \left( \hat{\xi} + \hat{\Lambda} \hat{\eta}_i \right)$$

## Illustrations

### Simulation study

In this section, we utilize 80 replications of data generated from multivariate regression, where the sample size is  $n = 50, 150$ ,  $p = 22$ , and  $q = 24$ . The choice of  $q$  is made to align with the dimension of the regression models applied to the GTEx data. In each

replication, a sparse matrix  $\mathbf{B}$  is generated using the element-wise product of three matrices:  $\mathbf{B} = \mathbf{W} \odot \mathbf{K} \odot \mathbf{Q}$  where  $\odot$  is the elementwise product,  $(\mathbf{W})_{ij} \sim N(0, 1)$  and  $(\mathbf{K})_{ij} \sim \text{Bernoulli}(s_1)$  and each row of  $\mathbf{Q}$  consists of either all 1's or all 0's with a success probability of 1's equal to  $s_2$ .

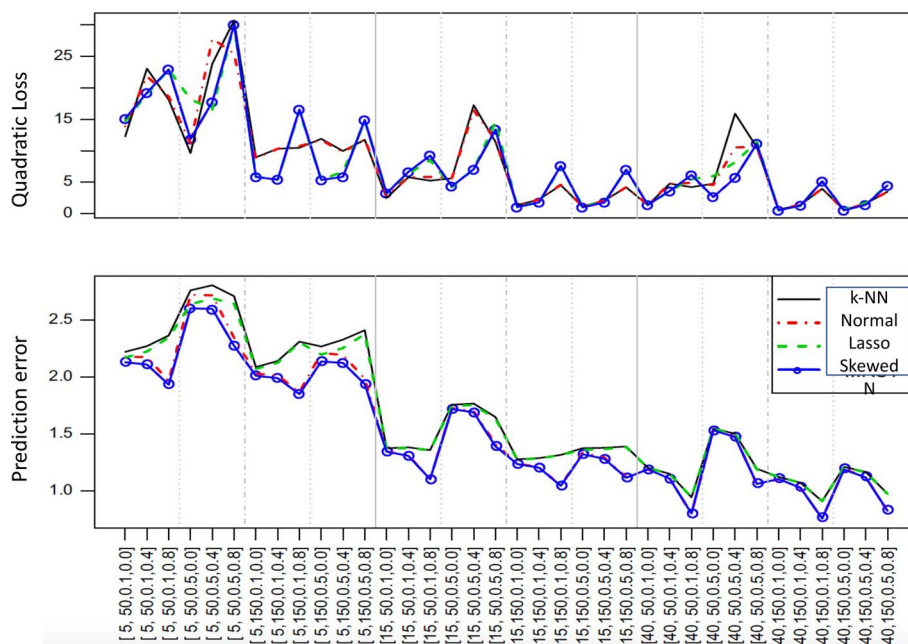
By generating  $\mathbf{B}$  in this manner, we anticipate that  $(1 - s_2)p$  predictors will be irrelevant for all  $q$  responses, and each predictor will be relevant for  $s_1q$  of all the response variables. An  $n \times p$  predictor matrix  $\mathbf{X}$  is also generated with rows drawn independently from  $N(\mathbf{0}, \Sigma_X)$ , where  $(\Sigma_X)_{ij} = 0.7^{|i-j|}$ , following the approach of Yuan and Lin (2007) [13] and Peng et al. (2010) [14]. We consider an AR(1) covariance structure for the scale matrix of the errors, which is  $\Sigma = \rho^{|i-j|}$ .

Lastly, every row of the error matrix  $\mathbf{E}$  is independently sampled from a multivariate skew-Normal distribution  $\text{SN}(\mathbf{0}, \Sigma, \Lambda)$ , and the response matrix  $\mathbf{Y}$  is formed as  $\mathbf{Y} = \mathbf{XB} + \mathbf{E}$ . To reduce computation time, we independently generate validation data (sample size  $n = 50$ ) within each replication to estimate the prediction error for the algorithms, akin to performing a K-fold cross-validation for the algorithm, as described in Rothman et al. (2010) [15].

We consider 36 different combinations of  $\Lambda$ ;  $\rho$ ;  $s_1$  and  $s_2$  from the following ranges: (1)  $\rho = \{0, 0.4, 0.8\}$ , (2)  $\Lambda = \text{diag}(-1, 1, -1, \dots, 1)$  or  $\mathbb{1}_q$ , where  $\mathbb{1}_q$  is a column vector of ones, (3)  $s_1 = 0.1, 0.5$ , and (4)  $s_2 = 1$ . Tuning parameters is selected from the set  $\{2^a, a \in 0; \pm 1, \dots, \pm 5\}$ .

### Results on simulation study

In the simulation study, we measure the overall performance of various methods in terms of the mean squared prediction error (PE). We have computed the prediction errors (PE) values of the entropy loss functions of the estimators of  $\Omega$  for the 80 simulated datasets



**Fig. 2** Prediction errors (PE) values of the entropy loss functions of the estimators of  $\Omega$  for the 80 simulated datasets using the k-NN, Normal, the G-Lasso and the TISLET (skewed) algorithm



**Table 1** TPR/TNR for the matrix  $\Omega$  averaged over 80 replications with  $s_1 = 0.1$ ,  $s_2 = 1$  and  $\lambda = (1; 1; 1; \dots; 1)^T$ 

Sample size	Correlation	Normal	Skewed Normal
n=50	$\rho = 0.0$	78.24/78.64	80.43/80.20
	$\rho = 0.4$	79.26/79.39	82.10/78.60
	$\rho = 0.8$	86.88/76.79	88.27/76.04
n=150	$\rho = 0.0$	86.73/77.68	88.55/76.72
	$\rho = 0.4$	88.41/74.60	89.66/74.93
	$\rho = 0.8$	92.16/74.27	93.35/73.43

**Table 2** TPR/TNR for the matrix  $\Omega$  averaged over 80 replications with  $s_1 = 0.5$ ,  $s_2 = 1$  and  $\lambda = (1; 1; 1; \dots; 1)^T$ 

Sample size	Correlation	Normal	Skewed Normal
n=50	$\rho = 0.0$	88.16/37.44	88.89/35.43
	$\rho = 0.4$	87.99/38.92	88.52/37.42
	$\rho = 0.8$	90.16/39.12	91.36/34.34
n=150	$\rho = 0.0$	92.88/43.17	92.97/43.58
	$\rho = 0.4$	93.21/42.20	93.27/41.60
	$\rho = 0.8$	95.76/36.93	95.97/36.81

using the k-NN, Normal, the G-lasso and the TISSLET algorithm (Fig. 2). To assess the significance of the difference in prediction error between the Normal and Tisslet algorithms, we conducted a Wilcoxon signed-rank test ( $R_i = \text{Sign}(D_i) - \text{rank}(|D_i|)$ ) where  $D_i$  is the difference of both predictions. The results indicated a statistically significant difference, with a p-value of 0.0124 at the 5% significance level ( $\alpha = 0.05$ ).

While the  $\ell_1$  loss is used to evaluate the performance of  $\Omega$ , the sparsity recognition performance of  $\Omega$  is measured by the true positive rate (TPR) as well as the true negative rate (TNR) defined as

$$\text{TPR}(\hat{\Omega}, \Omega) = \frac{\#\{(i, j); \hat{\Omega}_{ij} \neq 0 \text{ and } \Omega_{ij} \neq 0\}}{\#\{(i, j); \Omega_{ij} \neq 0\}}$$

$$\text{TNR}(\hat{\Omega}, \Omega) = \frac{\#\{(i, j); \hat{\Omega}_{ij} = 0 \text{ and } \Omega_{ij} = 0\}}{\#\{(i, j); \Omega_{ij} = 0\}}$$

The corresponding true positive rates (TPR) and true negative rates (TNR) for  $\Omega$  are reported in Tables 1 and 2. From these tables it is evident that a slightly higher TPR is accompanied by a lower TNR for our algorithm Tisslet algorithm versus Molstad et al. algorithm (normal assumption). We have also compared the numerical performance of the algorithms and the G-Lasso method obtained from Friedman et al. (2022) [16] by setting  $\Lambda$  to be  $\mathbf{0}$  for the 36 combinations. In general, it turns out that for higher error correlations ( $\rho = 0.8$ ) the mean square error of these methods is somewhat lower compared to the K-NN method, but the estimators of  $\Omega$  are improved considerably. Finally, the computational times of the two algorithms are computed.

The computational times for the four algorithms were evaluated and compared. On average, the CPU time ratios of Normal, K-NN, and Lasso to our algorithm,



Tisslet, were 3.79, 4.72 and 2.17 respectively. These computations were performed on the Panther cluster, equipped with 2.5 GHz processors. The differences in computational time are attributed to their computational complexities, which are generally  $O(np)$ , for all except for K-NN, which has a complexity of  $O(kn^2p)$ . Lasso exhibits a lower ratio compared to the other algorithms because it penalizes small values in the covariance matrix, effectively reducing the time spent on its computation.

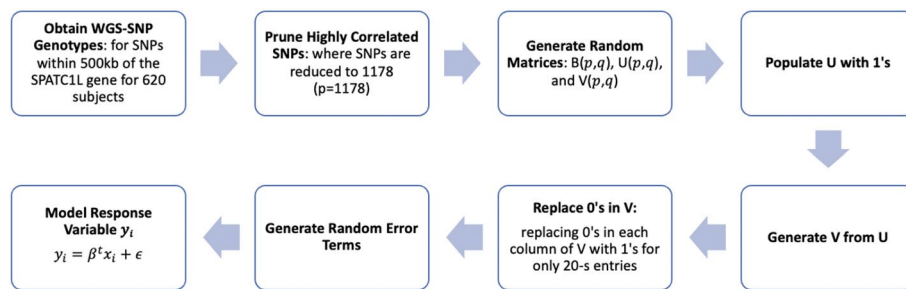
### Data generating study for SPATC1L gene

We conducted extensive numerical experiments to examine how the number of shared eQTLs, population  $R^2$  (also known as narrow-sense heritability), and how tissue-tissue correlation structure affect the performance of various methods for estimating eQTL weights across multiple tissues. To closely replicate the conditions of joint eQTL mapping in GTEx data, we obtained whole-genome sequencing SNP genotype data for all SNPs within 500kb of the SPATC1L gene for 620 subjects from the GTEx dataset [17]. After removing highly correlated SNPs (see Data Preparation section), we were left with  $p = 1178$  SNP genotypes. For each replication, we then generated  $n = 620$  subjects' expressions in  $q = 29$  tissues. Denoting  $\mathbf{x}_i \in \mathbb{R}^p$  as the SNP genotypes for the  $i^{th}$  subject, we generated  $\mathbf{y}_i \in \mathbb{R}^q$  as a realization of the random vector  $\beta_*^t \mathbf{x}_i + \epsilon$  for  $i = [n]$ , where  $\beta \in \mathbb{R}^{p \times q}$  are the eQTL weights and errors are independent and identically distributed as  $\text{SN}_q(0, \mathbf{A}, \mathbf{\Omega}_*^{-1})$  (Same protocol as [8] and [7]). For five hundred independent replications in each scenario, we randomly divided the  $n = 620$  subjects into a training set of size  $n_{\text{train}} = 400$ , a validation set of size 110, and a testing set of size 110.

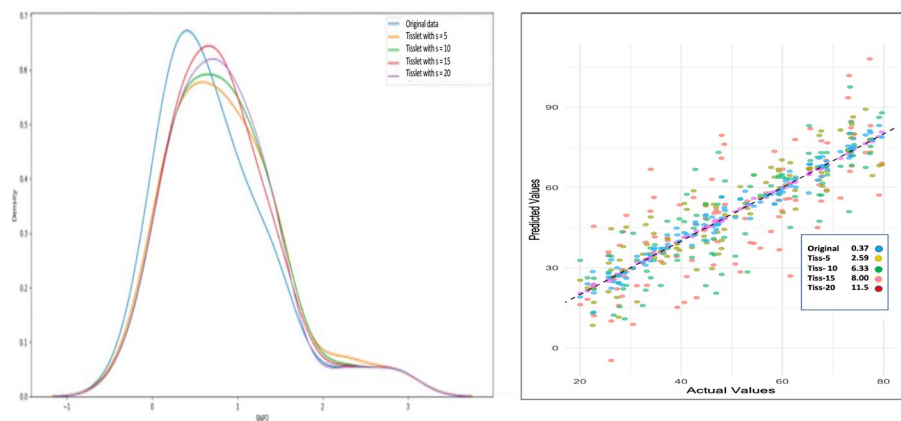
To create missingness, in each independent replication, we generated  $\mathbf{y}$  as follows: first, we created a matrix  $\mathbf{B} \in \mathbb{R}^{p \times q}$  with entries that were independent  $\mathcal{N}(0, 1)$ . Then, we generated a matrix  $\mathbf{S} \in \mathbb{R}^{p \times q}$  to be a matrix whose rows are either all zero or all one: we randomly selected  $s$  rows to be nonzero, where  $s \in [20]$ . Given  $\mathbf{S}$ , we then generated a matrix  $\mathbf{U} \in \mathbb{R}^{p \times q}$  so that each of the  $q$  columns has  $20 - s$  randomly selected entries equal to one only from entries which are zero in  $\mathbf{S}$  and all others equal to zero. We calculated  $\mathbf{y} = \mathbf{B} \odot \mathbf{S} + \mathbf{B} \odot \mathbf{U}$ , where  $\odot$  denotes the element-wise product. This construction ensured that each tissue has twenty total eQTLs,  $s$  of which are shared across all tissue types. We considered  $s = \{5, 10, 15, 18, 20\}$  in the simulations presented in this section. It is important to note that due to high linkage disequilibrium, many SNPs are highly correlated, resulting in a larger number of SNPs associated with gene expression. Figure 3 gives a summary of the protocol for generating data.

In addition, we randomly introduced missing values to both the training and validation set responses with a missing probability of 0.55, which corresponds to the missing rate in the GTEx gene expression data. For each method, we trained the model on the training data, chose tuning parameters using the validation data, and evaluated prediction and variable selection accuracy on the testing data.

To construct the covariance matrix  $\mathbf{\Omega}$ , we build it to have a block-diagonal structure and to control the  $R^2$ . Specifically, we consider a covariance structure for the scale matrix of the errors, which is  $\Sigma_{ij} = \rho^{|i-j|}$ , for  $i \in [20, 20]$  and  $\Sigma_{ij} = (\rho + 0.2)^{|i-j|}$ , for  $i \in [10, 10]$ .



**Fig. 3** Overview of the Data Pre-processing Steps. Input WGS SNP data then prune highly correlated SNPs. Create U and V and generate B. Generate error and fit the model

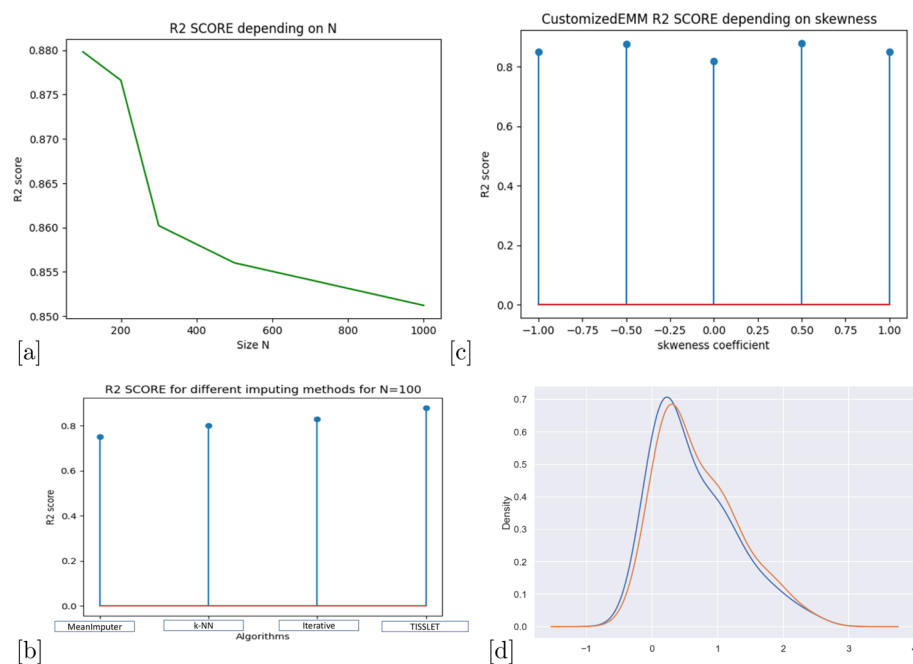


**Fig. 4** Left: comparing distribution of original versus Tislet imputed  $SNP_2$  using several choice of  $s$  (5, 10, 15, 20) distributions. Right: providing the Mean absolute error (MAE)

### Results on SPATC1L gene study

**Comparative Analysis of Covariance Approaches:** We evaluated the  $s$ -TISSLET imputation method's performance at different values of  $s$  for imputing missing data generated in section (3.3). Figure 4 shows that despite the varying  $s$  values, the imputed distributions closely align with the original distribution. This indicates that the imputation method maintains the overall structure of the gene expression data, which is critical for the reliability of subsequent analyses.

**Dataset size and prediction accuracy.** In the assessment of the TISSLET framework's performance, a nuanced relationship between dataset size ( $N$ ) and  $R^2$  scores became apparent (Fig. 5a). An increase in  $N$  was correlated with a decline in  $R^2$  scores, underscoring the challenges in sustaining prediction accuracy with the expansion of data volume inherent in multi-omics studies. Despite this, the TISSLET framework displayed a remarkable resilience in  $R^2$  scores across varying levels of skewness in gene expression data (Fig. 5b), attesting to its robustness and adaptability to the diverse data distributions encountered in multi-omics research. Further illustrating its efficacy, the TISSLET method demonstrated superior predictive accuracy when compared with other imputation methods such as MEANimputer,  $k$ -NN and iterative



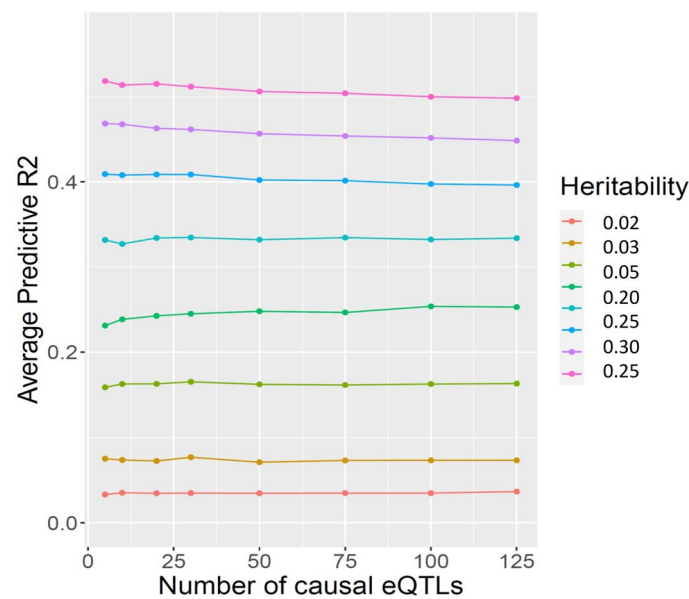
**Fig. 5** **a** Relationship between dataset size and  $R^2$  scores. TISSLET (CEM)  $R^2$  scores across skewness coefficients. **b** A bar graph displaying the TISSLET framework's  $R^2$  scores, indicating consistent predictive accuracy despite varying levels of skewness  $\{-1.0; -0.5; 0; 0.5; 1.0\}$  in gene expression data. **c**  $R^2$  score comparison among imputing methods for  $N = 100$ . **d** Bar chart illustrating  $R^2$  scores for different imputing methods, highlighting the superior performance of the TISSLET framework over MEANimputer, k-Nearest Neighbor ( $k=2$ ), and iterative methods

imputation for multi-tissue gene expression prediction (Fig. 5c), reinforcing its suitability for complex genomic analyses.

**Impact of the number of causal variants in simulation study:** We extended our simulation study to assess the impact of a larger number of causal variants, ranging from 5 to 125, across a wide range of heritability levels. We found that increasing the number of causal variants had very little effect on the predictive performance of TISSLET. We believe this is because the expected imputation accuracy primarily depends on the total heritability explained by the causal SNPs (Fig. 6).

**Performance of fitting the skewed normal versus normal on predicted gene:** Finally the density plot provides a visual comparison between the distributions of gene expression predictions using two [8] and our method of imputation. The curves are closely aligned and have a similar shape, peaking around a value of 1, which suggests that both methods yield relatively similar predictions in terms of their distribution. However, one curve is slightly shifted towards the right indicating that our approach may predict slightly higher expression values compared to the one that uses Normality assumption (Fig. 5d).

**Results of for SPATC1L gene study:** In this section, we present results with tissue-tissue correlation for several errors, where  $\rho \in \{0; 0.1; 0.3; 0.5; 0.7\}$  varying. In this setting, we observe that our method, TISSLET performs better than all realistic competitors: only G-Lasso, outperforms slightly our method when  $\rho$  is less or equal than 0.015. As one would expect, when expression is nearly uncorrelated ( $\rho = 0$ ), our method



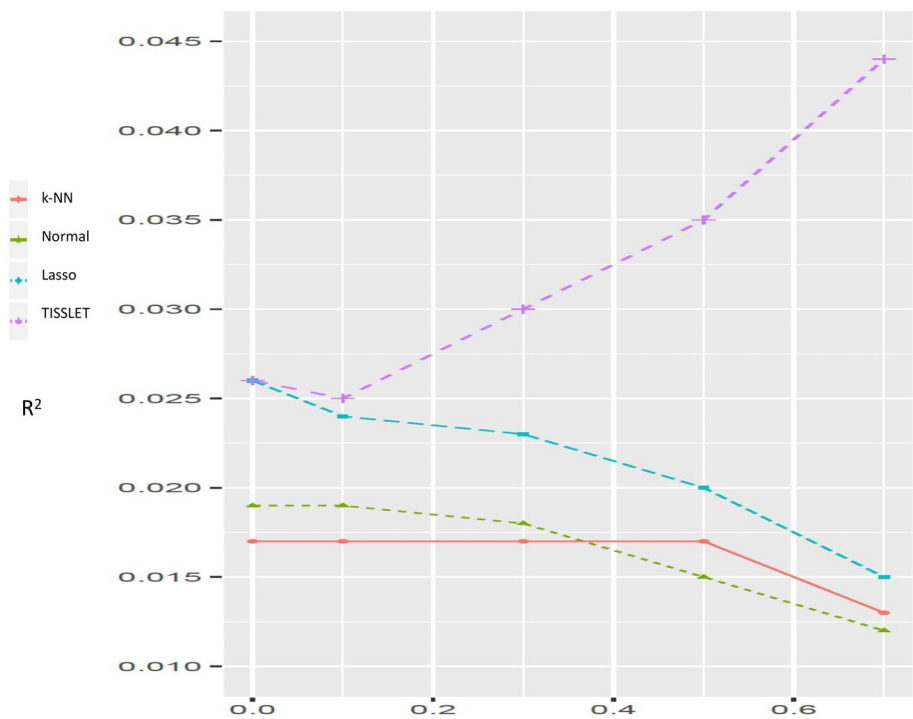
**Fig. 6** Impact of number of causal variants in simulation study

TISSLET performs better than all the others 3 methods that does implicitly assumes no tissue-tissue correlation. Remarkably, when  $\rho$  is greater than or equal to 0.3, TISSLET outperforms even the normal method which assume normality and missingness. In fact, the prediction accuracy of TISSLET increases as  $\rho$  increases, whereas all other methods, which do not explicitly model tissue-tissue correlation, have prediction accuracy remaining constant or slightly decreasing as  $\rho$  increases. This demonstrates the benefit of not only accounting for tissue-tissue correlation in multi-tissue joint eQTL mapping but also assume the correct distribution of the gene expression when expression across tissue types can be reasonably assumed to be conditionally dependent (Fig. 7).

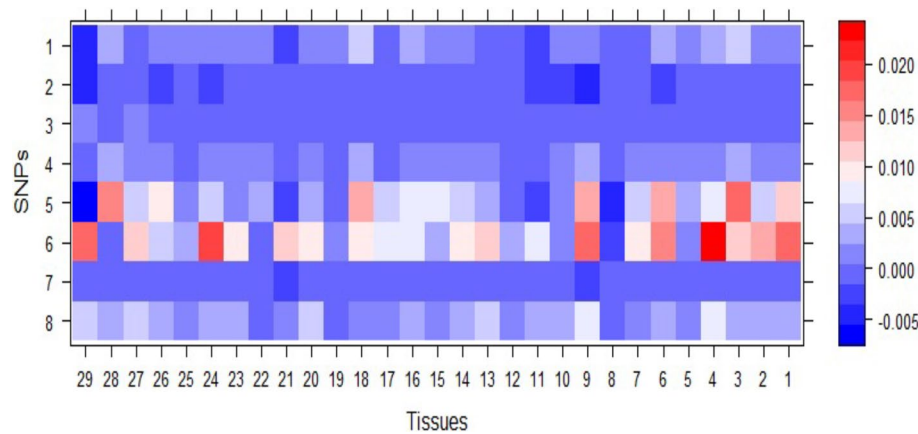
Furthermore, the heatmap of gene expressions presented in Fig. 8 illustrates a detailed example of a specific gene. In this case, out of the eight eQTLs (expression quantitative trait loci) that have been identified for this gene, seven of these eQTLs are consistently shared across all 29 examined tissues, demonstrating a uniform regulatory effect. Additionally, one of the eQTLs is shared across 28 of the 29 tissues, indicating a high level of commonality with a single tissue exception. This comprehensive visualization underscores the widespread influence of these eQTLs on gene expression across a broad range of tissue types, highlighting their significant role in the regulatory network.

## Conclusion and discussion

Our results suggest that the proposed TISSLET framework can robustly handle missing data in multi-tissue gene expression studies. The framework's ability to retain accuracy across different imputation parameters, such as 'k', and its superior performance compared to other imputation methods, underscores its potential in improving



**Fig. 7** Average test  $R^2$  for 4 competing methods where  $\rho$  the correlation of the errors varies from 0.0 to 0.8



**Fig. 8** A heatmap depicting the expression of SNPs per tissue, with darker blue shades indicating weaker relationships

complex trait predictions in multi-omics studies. The consistency of  $R^2$  scores regardless of skewness and dataset size provides evidence of the model's stability. However, the observed decline in prediction accuracy with increasing dataset size suggests that further optimization of the framework may be necessary to manage the complexities of large-scale multi-omics data.

In conclusion, our findings advocate for the TISSLET framework's application in multi-omics studies to improve the prediction of complex traits. Its ability to accurately account for tissue-tissue correlations and withstand data distribution asymmetries positions it as a powerful tool in the advancement of personalized medicine. As the complexity of multi-omics studies grows, the development of robust statistical machine learning like TISSLET is crucial for unveiling the genetic foundations of complex traits and advancing our understanding of gene expression dynamics.

While the TISSLET framework marks a significant advancement in gene expression prediction across multiple tissues, this study demonstrates its precision and stability, showcasing its high accuracy in predictions. However, the framework is still in its foundational stages and requires further refinement. The next steps involve optimizing the model to enhance its performance, particularly for large-scale datasets, thereby fully realizing its potential in complex multi-omics research and personalized medicine.

## Appendix

### Definition and notation

- $\beta_{.,m_i} = \mathbf{M}_i \beta_{.,i}^t$        $\beta_{.,o_i} = \mathbf{O}_i \beta_{.,i}^t$
- $\Sigma_{o_i} = \mathbf{O}_i \Sigma \mathbf{O}_i^t$        $\Sigma_{o_i} \in \mathbf{R}^{(q_0 \times q_0)}$
- $\Sigma_{m_i} = \mathbf{M}_i \Sigma \mathbf{M}_i^t$        $\Sigma_{m_i} \in \mathbf{R}^{(q-q_0) \times (q-q_0)}$
- $\Sigma_{o_i, m_i} = \mathbf{O}_i \Sigma \mathbf{M}_i^t$        $\Sigma_{o_i, m_i} \in \mathbf{R}^{(q_0 \times (q-q_0))}$
- $\Sigma_{m_i, o_i} = \mathbf{M}_i \Sigma \mathbf{O}_i^t$        $\Sigma_{m_i, o_i} \in \mathbf{R}^{(q-q_0) \times q_0}$
- $\mathbf{y}_i^{m_i} \sim \mathbf{N}_{q-q_0}(\beta_{.,m_i}^T \mathbf{x}_i, \Sigma_{m_i})$
- $\mathbf{y}_i^{o_i} \sim \mathbf{N}_{q_0}(\beta_{.,o_i}^T \mathbf{x}_i, \Sigma_{o_i})$

### Penalized normal log-likelihood

Using Aaron et al. (2020), then can summarize the derivation of EM for estimating missing gene expression  $\hat{y}_i$ , eQTL weights  $\beta$  and cross-tissue correlation  $\Omega$ .

$$\mathbf{y}_i^{m_i} | \mathbf{y}_i^{o_i}, \Theta \sim \mathbf{N}(\mu_i, \mathbf{V}_i)$$

where

$$\mu_i = \beta_{.,m_i}^T \mathbf{x}_i + \Sigma_{m_i, o_i} \Sigma_{o_i}^{-1} (\mathbf{y}_i^{o_i} - \beta_{.,o_i}^T \mathbf{x}_i) \quad (5.1)$$

$$\mathbf{V}_i = \Sigma_{m_i} - \Sigma_{m_i, o_i} \Sigma_{o_i}^{-1} \Sigma_{o_i, m_i} \quad (5.2)$$

Then the negative log-likelihood of the **conditional** multivariate normal distribution for  $\mathbf{y}_{i, m_i} | \mathbf{y}_{i, o_i}$  is proportional to

$$\frac{1}{n} \sum_i \left( \mathbf{y}_i^{o_i} - \beta_{.,o_i}^T \mathbf{x}_i \right)^T \left( \mathbf{y}_i^{m_i} - \beta_{.,m_i}^T \mathbf{x}_i \right) \Sigma^{-1} \left( \mathbf{y}_i^{o_i} - \beta_{.,o_i}^T \mathbf{x}_i \right)^T \left( \mathbf{y}_i^{m_i} - \beta_{.,m_i}^T \mathbf{x}_i \right)^T + \log |\Sigma|$$

The negative log likelihood is useful in the EM steps because we needed to calculate the expectation of the log-likelihood  $Q(\Theta)$ :

$$Q(\Theta) = \text{tr}\{S_{(\beta, \Sigma)}\Omega\} - \log \det \Omega$$

We regularize the entries of  $\beta$  and  $\Omega$  using  $l_1$  penalties, and estimate them by minimizing the penalties likelihood

$$Q(\Theta) = \text{tr}\{S_{(\beta, \Sigma)}\Omega\} - \log \det \Omega + \lambda_1 \|\Omega\|_1 + \lambda_2 \sum_{jk} |\beta_{jk}|$$

where the empirical covariance matrix

$$S_{(\beta, \Sigma)} = \frac{1}{n} \sum_{i=1}^n \Gamma_i \quad (5.3)$$

$$\Gamma_i = \begin{pmatrix} \Gamma_{o_i, o_i} & \Gamma_{o_i, m_i} \\ \Gamma_{m_i, o_i} & \Gamma_{m_i, m_i} \end{pmatrix}$$

$$\Gamma_{o_i, o_i} = (y_{i, o_i}^T - x_i^T \beta_{\cdot, o_i})^T (y_{i, o_i}^T - x_i^T \beta_{\cdot, o_i})$$

$$\Gamma_{m_i, m_i} = (\mu_{i, m_i}^T - x_i^T \beta_{\cdot, m_i})^T (\mu_{i, m_i}^T - x_i^T \beta_{\cdot, m_i}) + V_{m_i} \quad (5.4)$$

$$\Gamma_{m_i, o_i} = (y_{i, o_i}^T - x_i^T \beta_{\cdot, o_i})^T (\mu_{i, m_i}^T - x_i^T \beta_{\cdot, m_i}) \quad (5.5)$$

Then the algorithm is summarized as the following:

#### Penalized skew-normal log-likelihood

$$\mathcal{L}(\Theta | y^o, y^m, \tau) = \log p(y^m | y^o, \tau) p(\tau | y^o) \quad (5.6)$$

$$\mathcal{L}(\Theta | y) = \frac{1}{2} \sum_{i=1}^n \{ \log(\det \Omega^{-1}) + (y_i - \beta^t x_i - \Lambda \tau) \Omega (y_i - \beta^t x_i - \Lambda \tau) + \tau_i \tau_i^t \} \quad (5.7)$$

where  $\Theta = (\beta, \Omega, \Lambda)$  represents all unknown parameters. In the E-step of ECM, we need to calculate the Q-function, which is the conditional expectation of the complete data log-likelihood function (7) given the observed data  $y^o$  and the current estimate  $\hat{\Theta}^k = (\hat{\beta}, \hat{\Sigma}, \hat{\Lambda})$ . Herein, the term  $-\frac{1}{2} E(\tau_j^t \tau_j | y_j^o, \Theta^k)$  can be omitted because it does not include any parameters. Therefore, we have

$$Q(\Theta | \Theta^{(k)}) = \frac{1}{2} \sum_{i=1}^n \left[ \log |\Omega| - \mathbb{E}(Z_i Z_i^t) \right] \quad (5.8)$$

Using the formula  $\mathbb{E}(ZZ^t) = \text{Var}(Z) + \mathbb{E}(Z)\mathbb{E}(Z)^t$ , we have:



$$\mathbb{E}(\mathbf{Z}\mathbf{Z}^t) = \text{Var} \left[ \left( \mathbf{y} - \xi - \mathbf{\Lambda} \boldsymbol{\tau} \right) \boldsymbol{\Omega}^{1/2} \right] + \left( \mathbb{E}(\mathbf{y}|\mathbf{y}^o) - \xi - \mathbf{\Lambda} \hat{\boldsymbol{\eta}} \right) \left( \mathbb{E}(\mathbf{y}|\mathbf{y}^o) - \xi - \mathbf{\Lambda} \hat{\boldsymbol{\eta}} \right) \boldsymbol{\Omega} \quad (5.9)$$

$$= \text{Var} \left[ \left( \hat{\mathbf{y}} - \xi - \mathbf{\Lambda} \boldsymbol{\tau} \right) \boldsymbol{\Omega}^{1/2} \right] + \text{tr} \left( \hat{\mathbf{y}} - \xi - \mathbf{\Lambda} \hat{\boldsymbol{\eta}} \right) \left( \hat{\mathbf{y}} - \xi - \mathbf{\Lambda} \hat{\boldsymbol{\eta}} \right)^t \boldsymbol{\Omega}$$

Using iterative expectation and  $\mathbf{O}_i^t \mathbf{O}_i (\mathbf{I}_p - \boldsymbol{\Sigma} \mathbf{S}_i^{oo}) = 0$ , we have

$$\mathbb{E}(\mathbf{y}_i^m | \mathbf{y}_i^o) = \mathbf{M}_i (\xi + \mathbf{\Lambda} \boldsymbol{\eta}_i + \boldsymbol{\Sigma} \mathbf{S}_i^{oo} (\mathbf{y}_i - \xi - \mathbf{\Lambda} \boldsymbol{\eta}_i))$$

$$\text{cov}(\mathbf{y}_i, \mathbf{y}_i) = \mathbf{M}_i (\mathbf{I}_p - \boldsymbol{\Sigma} \mathbf{S}_i^{oo}) (\boldsymbol{\Sigma} + \mathbf{\Lambda} (\boldsymbol{\Psi}_i - \boldsymbol{\eta}_i \boldsymbol{\eta}_i^t) \mathbf{\Lambda} (\mathbf{I}_p - \mathbf{S}_i^{oo} \boldsymbol{\Sigma})) \mathbf{M}_i^t$$

$$\mathbb{E}(\mathbf{y}_i \boldsymbol{\tau}^t | \mathbf{y}_i^o) = \mathbf{M}_i (\mathbf{I}_p - \boldsymbol{\Sigma} \mathbf{S}_i^{oo}) (\xi \boldsymbol{\eta}_i^t + \mathbf{\Lambda} \boldsymbol{\Psi}_i) + \boldsymbol{\Sigma} \mathbf{S}_i^{oo} \mathbf{y}_i \boldsymbol{\eta}_i^t$$

After simplification we have

$$\mathbf{Q}(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}}) = \frac{1}{2} \sum_i \left\{ \text{tr} \left( \mathbf{R}_{(\xi, \mathbf{\Lambda})} \boldsymbol{\Omega}^{-1} \right) - \log \det (\boldsymbol{\Omega}) \right\} + \lambda_1 \|\boldsymbol{\Omega}\|_1 + \lambda_2 \sum_{jk} |\boldsymbol{\beta}_{jk}| \quad (5.10)$$

$$\hat{\mathbf{R}}_i(\xi, \mathbf{\Lambda}) = \left( (\mathbf{I}_p - \hat{\boldsymbol{\Omega}}^{-1} \hat{\mathbf{S}}_i^{oo} \hat{\mathbf{\Lambda}} - \mathbf{\Lambda}) \left( \hat{\boldsymbol{\Psi}}_i - \hat{\boldsymbol{\eta}}_i \hat{\boldsymbol{\eta}}_i^t \right) \times \left( (\mathbf{I} - \hat{\boldsymbol{\Omega}}^{-1} \hat{\mathbf{S}}_i^{oo} \hat{\mathbf{\Lambda}} - \mathbf{\Lambda}) \right.$$

$$\left. + \left( \mathbf{I} - \hat{\boldsymbol{\Omega}}^{-1} \hat{\mathbf{S}}_i^{oo} \right) \hat{\boldsymbol{\Omega}}^{-1} + \left( \hat{\mathbf{y}}_i - \xi - \mathbf{\Lambda} \hat{\boldsymbol{\eta}}_i \right) \left( \hat{\mathbf{y}}_i - \xi - \mathbf{\Lambda} \hat{\boldsymbol{\eta}}_i \right) \right)$$

$$\hat{\mathbf{S}}_i^{oo} = \mathbf{O}_i^t \left( \mathbf{O}_i \hat{\boldsymbol{\Omega}}^{-1} \mathbf{O}_i \right)^{-1} \mathbf{O}_i$$

$$\hat{\boldsymbol{\eta}}_i = \mathbb{E} \left( \boldsymbol{\tau}_i | \mathbf{y}_i^o, \boldsymbol{\Theta} \right)$$

$$\hat{\boldsymbol{\Psi}}_i = \mathbb{E} \left( \boldsymbol{\tau}_i \boldsymbol{\tau}_i^t | \mathbf{y}_i^o, \boldsymbol{\Theta} \right)$$

The predicted gene expression  $\hat{\mathbf{y}}_i$  is given by

$$\hat{\mathbf{y}}_i = \mathbb{E}(\mathbf{y}_i | \mathbf{y}_i^o, \boldsymbol{\Theta}) = \boldsymbol{\Omega}^{-1} \mathbf{S}_i^{oo} \mathbf{y}_i + \left( \mathbf{I} - \boldsymbol{\Omega}^{-1} \mathbf{S}_i^{oo} \right) \left( \hat{\xi} + \hat{\mathbf{\Lambda}} \hat{\boldsymbol{\eta}}_i \right)$$

Using those derivation, Algorithm becomes:

**Algorithm 3** Regularized EM Algorithm with skewed-normality

---

```

1: Initialize:  $\Omega^0, \xi^0, \Lambda^0$  and  $\eta^0$ .
2: Set  $k = 1$ .
3: E-Step:
4: while not converged do
5:    $\Upsilon_i^k \leftarrow I - \hat{\Omega}^{-1(k)} \hat{S}_i^{oo(k)}$ 
6:    $\hat{\mathbf{y}}_i^{(k+1)} \leftarrow \hat{\Omega}^{(k)} \hat{S}^{oo(k)} \mathbf{y}_i^{(k)} + \Upsilon^k (\hat{\xi}^{(k)} + \hat{\Lambda}^{(k)} \hat{\eta}_i^{(k)})$ 
7:    $\mathbf{C}_i = \mathbf{O}_i^t \Omega^{oo} \mathbf{O}_i$ 
8:    $\tau^{(k+1)} | \mathbf{y}_i^{(k)} \leftarrow TN_q \left( \hat{\Lambda}^{(k)} \mathbf{C}_i^{oo(k)} (\hat{\mathbf{y}}_i^{(k)} - \hat{\xi}^{(k)}), (I - \hat{\Lambda}^{(k)} \mathbf{C}_i^{oo(k)} \hat{\Lambda}^{(k)}), \mathbb{R}_+^q \right)$ 
9:    $\hat{\eta}^{(k+1)} \leftarrow \mathbb{E}(\tau | \mathbf{y}_i)$ 
10:   $\hat{\psi}^{(k+1)} \leftarrow \text{var}(\tau | \mathbf{y}_i) + \hat{\eta}^{(k)} \hat{\eta}^{(k)}$ 
11: end while
12: CM Step 1:
13:  $\hat{\xi}^{(k+1)} \leftarrow \frac{1}{n} \left( \sum_{i=1}^n \hat{\mathbf{y}}_i^{(k)} - \hat{\Lambda}^{(k)} \sum_{i=1}^n \hat{\eta}_i^{(k)} \right)$ 
14: CM-Step 2:
15:  $\hat{\Omega}^{(k+1)} \leftarrow \frac{1}{n} \sum_{i=1}^n \mathbf{R}_i^{(k)}$ 
16: For  $i := 1$  to  $n$ 
17:   $\mathbf{R}_i^{(k+1)} \leftarrow \left( \Upsilon_i^k \hat{\Lambda} - \Lambda \right) \left( \hat{\Psi}_i - \hat{\eta}_i \hat{\eta}_i^t \right) \left( \Upsilon_i^k \hat{\Lambda} - \Lambda \right)^t + \Upsilon_i^k \hat{\Omega}^{-1} + (\hat{\mathbf{y}}_i - \xi - \Lambda \hat{\eta}_i)(\hat{\mathbf{y}}_i - \xi - \Lambda \hat{\eta}_i)$ 
18: End
20: CM-Step 3:
21:  $\hat{\Lambda}^{(k+1)} \leftarrow \text{diag} \left\{ (\hat{\Omega}^{(k)} \odot \sum_{i=1}^n \hat{\Psi}^{(k)})^{-1} (\hat{\Omega}^{(k+1)} \odot \sum_{i=1}^n \hat{\Phi}_i^{(k)}) \mathbf{1}_p \right\}$ 
22:
23: For  $i := 1$  to  $n$   $\hat{\Phi}_i^{(k+1)} \leftarrow (\hat{\Psi}^{(k)} - \hat{\eta}^{(k)} \hat{\eta}^{(k)t}) \hat{\Lambda}^{(k)} (\mathbf{I} - \hat{\Omega}^{(k)-1} \hat{S}^{oo(k)}) + \hat{\eta}^{(k)} (\hat{\mathbf{y}}_i - \hat{\xi}^{(k)})$ 
24: End

```

---

**Precision matrix estimator**

To estimate the inverse of the covariance matrix, several thresholding approaches have been proposed. The two most popular ones that we propose to use in this paper are GLASSO, available in R, that was proposed par Friedman et al. (2007) [16] and CLIME that was proposed by Cai et al. (2011) [18] If we use CLIME from Cai and Zhou approach, their equation emphasizes the fact that the solution may not be unique. Such nonuniqueness always occurs since  $\Sigma = \mathbf{X}\mathbf{X}^T$  where  $\mathbf{X} \in \mathbb{R}^{n \times p}$  is the sample data matrix and  $\text{rank}(\mathbf{X}) \leq n \ll p$  which is our case. However, there is no guarantee that the thresholding estimator is always positive definite [19]. Although the positive definite property is guaranteed in the asymptotic setting with high probability, the actual estimator can be an indefinite matrix, especially in real data analysis.

In the following we will describe quickly the argument and parameters of implementing BigQuiC in R.

**Previous ad-hoc process for association SNPs with tissues**

1. Create a copy of the genotype data matrix for each tissue. Only keep the samples that are matched with the corresponding gene expression and covariates data.
2. Remove SNPs with low minor allele frequency (MAF).

3. Remove genes with consistently low expression level.
4. For each tissue, load the tailored genotype data, gene expression data, covariates data, and gene/SNP location files into R.
5. To do a single tissue analysis, we will feed the data into MatrixEQTL to conduct single-tissue eQTL analysis and obtain summary statistics (i.e., t-statistics).
6. Using the cis window size to be 1 kb (i.e.,  $1 \times 10^5$ ) or 1 Mb (i.e.,  $1 \times 10^6$ ) and the p value threshold to be 1 in order to output summary statistics for all cis gene-SNP pairs then we convert each summary statistic  $t$  to a correlation  $r$  in each tissue  $r = \frac{1}{\sqrt{df+t^2}}$  where df is the degree of freedom (the number of samples minus the number of covariates minus two) in the corresponding tissue.
7. Further we convert the correlations to z-statistics using Fisher transformation  $z = \frac{1}{2} \sqrt{df-1} \log\left(\frac{1+r}{1-r}\right)$ .
8. Then we get the list of common gene-SNP pairs across all tissues and curate the obtained z-statistics into a matrix where each row corresponds to a gene-SNP pair in the common list and each column corresponds to a tissue. The z-statistics matrix is all we need for subsequent analyses.

#### Acknowledgements

We extend our gratitude to Dr. Molstad from the University of Minnesota for his invaluable assistance in providing access to his code and offering detailed explanations.

#### Author Contributions

The mathematical model was developed by H.B., writing the theoretical part and overview revision of the paper. M.C. Revised the paper. The author A.M. worked on running the mathematical model and producing the results. The author T.H. worked on producing the results in the earlier stage. The author A.A. worked on writing the results and discussion part with H.B. and A.M.; The corresponding author for submitting the article is A.A..

#### Funding

Not Applicable.

#### Availability of data and materials

Not Applicable.

#### Declarations

##### Ethics approval and consent to participate

Not Applicable.

##### Competing interests

The authors have no Conflict of interest to declare that are relevant to the content of this article.

Received: 19 September 2024 Accepted: 27 December 2024

Published online: 25 February 2025

#### References

1. Grinberg NF, Wallace C. Multi-tissue transcriptome-wide association studies. *Genetic Epidemiology*. 2021;45(3):324–37.
2. Gamazon ER, Wheeler HE, Shah KP, Mozaffari SV, Aquino-Michaels K, Carroll RJ, Eyler AE, Denny JC, Consortium G, Nicolae DL, et al. A gene-based association method for mapping traits using reference transcriptome data. *Nature genetics*. 2015;47(9):1091–8.
3. Manor O, Segal E. Genoexp: a web tool for predicting gene expression levels from single nucleotide polymorphisms. *Bioinformatics*. 2015;31(11):1848–50.
4. Lage K, Hansen NT, Karlberg EO, Eklund AC, Roque FS, Donahoe PK, Szallasi Z, Jensen TS, Brunak S. A large-scale analysis of tissue-specific pathology and gene expression of human disease genes and complexes. *Proceedings of the National Academy of Sciences*. 2008;105(52):20870–5.
5. Bai Y, Wang X, Hou J, Geng L, Liang X, Ruan Z, Guo H, Nan K, Jiang L. Identification of a five-gene signature for predicting survival in malignant pleural mesothelioma patients. *Frontiers in Genetics*. 2020;11:899.

6. Zhang N, Li P-c, Liu H, Huang T-c, Liu H, Kong Y, Dong Z-c, Yuan Y-h, Zhao L-l, Li J-h. Water and nitrogen in-situ imaging detection in live corn leaves using near-infrared camera and interference filter. *Plant Methods*. 2021;17:1–11.
7. Hu Y, Li M, Lu Q, Weng H, Wang J, Zekavat SM, Yu Z, Li B, Gu J, Muchnik S, et al. A statistical framework for cross-tissue transcriptome-wide association analysis. *Nature genetics*. 2019;51(3):568–76.
8. Molstad AJ, Sun W, Hsu L. A covariance-enhanced approach to multi-tissue joint eqtl mapping with application to transcriptome-wide association studies. *The annals of applied statistics*. 2021;15(2):998.
9. Church BV, Williams HT, Mar JC. Investigating skewness to understand gene expression heterogeneity in large patient cohorts. *BMC bioinformatics*. 2019;20:1–14.
10. Hosking JRM, Wallis JR, Wood EF. Estimation of the generalized extreme-value distribution by the method of probability-weighted moments. *Technometrics*. 1985;27(3):251–61.
11. Neal RM, Hinton GE. A view of the em algorithm that justifies incremental, sparse, and other variants. In: *Learning in Graphical Models*, 1998;pp. 355–368. Springer, ???
12. Hsieh C-J, Sustik MA, Dhillon IS, Ravikumar PK, Poldrack R. Big & quic: Sparse inverse covariance estimation for a million variables. *Advances in neural information processing systems* 2013;**26**
13. Yuan M, Lin Y. Model selection and estimation in the gaussian graphical model. *Biometrika*. 2007;94(1):19–35.
14. Peng C. Estimating and testing quantile-based process capability indices for processes with skewed distributions. *Journal of Data Science*. 2010;8(2):253–68.
15. Rothman AJ, Levina E, Zhu J. Sparse multivariate regression with covariance estimation. *Journal of Computational and Graphical Statistics*. 2010;19(4):947–62.
16. Friedman JTR, Hastie T: Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* 2007;**9**(3)
17. Morgan A, Vuckovic D, Krishnamoorthy N, Rubinato E, Ambrosetti U, Castorina P, Franzè A, Vozzi D, La Bianca M, Cappellani S, et al. Next-generation sequencing identified spact1l as a possible candidate gene for both early-onset and age-related hearing loss. *European Journal of Human Genetics*. 2019;27(1):70–9.
18. Cai W.L. Tony, Luo X. A constrained  $\ell_1$  minimization approach to sparse precision matrix estimation. *Journal of the American Statistical Association*. 2011;106(494):594–607.
19. Elanbari M, Rawi R, Ceccarelli M, Bouhali O, Bensmail H. Advanced computation of a sparse precision matrix hadap: A hadamard-dantzig estimation of a sparse precision matrix. In: *Proceedings of the Sixth International Conference on Computational Logics, Algebras, Programming, Tools, and Benchmarking* 2015. Citeseer

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.