

Supplementary Material for “HDTD: Analyzing multi-tissue gene expression data”

Anestis Touloumis^{1,2*}, John C. Marioni^{1,3} and Simon Tavaré¹

¹ CRUK Cambridge Institute, University of Cambridge, Li Ka Shing Centre, Robinson Way, Cambridge CB2 0RE, United Kingdom.

² Computing, Engineering and Mathematics, University of Brighton, Brighton, BN2 4GJ, United Kingdom.

³ EMBL-European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, United Kingdom.

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Associate Editor: XXXXXXXX

STATISTICAL METHODS

A full description of the hypothesis testing procedures for the mean matrix and the covariance matrices can be found in Touloumis *et al.* (2015) and Touloumis *et al.* (2014), respectively. Herein, we derive the exact formulae for the shrinkage estimators of the two covariance matrices Σ_C and Σ_R . To accomplish this, we extend the Stein-type covariance matrix estimation approach employed in Touloumis (2015) for estimating a single covariance matrix with vector-valued random variables to simultaneous estimation of two covariance matrices with matrix-valued random variables. For identifiability reasons, we impose the restriction $\text{tr}(\Sigma_R) = r$ (see Touloumis *et al.*, 2014), where $\text{tr}(\mathbf{A})$ denotes the trace of the square matrix \mathbf{A} , that is the sum of the diagonal elements of \mathbf{A} . It can be readily shown that the gene-wise variance scaling restriction $\text{tr}(\Sigma_R) = r$ does not affect the gene- and column-wise correlation pattern.

To estimate the column covariance matrix Σ_C , first let

$$\Sigma_C^* = (1 - \lambda_C) \frac{1}{Nr} \sum_{i=1}^N \mathbf{Y}_i^T \mathbf{Y}_i + \lambda_C \mu_C \mathbf{I}_c \quad (1)$$

where $\mathbf{Y}_i = \mathbf{X}_i - \widehat{\mathbf{M}}$, $\widehat{\mathbf{M}}$ is the sample mean matrix, and where \mathbf{I}_c is the identity matrix of size c . The optimal shrinkage intensity λ_C and the parameter μ_C are defined so as to minimise the expectation of $\text{tr}[(\Sigma_C^* - \Sigma_C)^2]$. This strategy leads to the closed form solutions

$$\lambda_C = \frac{\frac{\text{tr}(\Sigma_R^2)}{r^2(N-1)} [\text{tr}^2(\Sigma_C) + \text{tr}(\Sigma_C^2)]}{\frac{\text{tr}(\Sigma_R^2)}{r^2(N-1)} [\text{tr}^2(\Sigma_C) + \text{tr}(\Sigma_C^2)] + \text{tr}(\Sigma_C^2) - \text{tr}^2(\Sigma_C)/c}$$

and $\mu_C = \text{tr}(\Sigma_C)/c$. Both λ_C and μ_C depend on the unknown parameters $\text{tr}(\Sigma_C)$, $\text{tr}(\Sigma_C^2)$ and $\text{tr}(\Sigma_R^2)$ which can be estimated

by

$$T_{1N} = \frac{1}{rN} \sum_{i=1}^N \text{tr}(\mathbf{X}_i^T \mathbf{X}_i) - \frac{1}{cN(N-1)} \sum_{i \neq j} \text{tr}(\mathbf{X}_i^T \mathbf{X}_j),$$

$$T_{2N} = \frac{1}{r^2 N(N-1)} \sum_{i \neq j} \text{tr}(\mathbf{X}_i^T \mathbf{X}_i \mathbf{X}_j^T \mathbf{X}_j) - \frac{2}{r^2 N(N-1)(N-2)} \sum_{i \neq j \neq k} \text{tr}(\mathbf{X}_i^T \mathbf{X}_i \mathbf{X}_j^T \mathbf{X}_k) + \frac{1}{r^2 N(N-1)(N-2)(N-3)} \sum_{i \neq j \neq k \neq l} \text{tr}(\mathbf{X}_i^T \mathbf{X}_j \mathbf{X}_k^T \mathbf{X}_l),$$

and

$$T_{3N} = \frac{T_{2N}^{-1}}{N(N-1)} \sum_{i \neq j} (\mathbf{R}_i^T \mathbf{R}_j)^2 - \frac{2T_{2N}^{-1}}{N(N-1)(N-2)} \sum_{i \neq j \neq k} \mathbf{R}_i^T \mathbf{R}_j \mathbf{R}_i^T \mathbf{R}_k + \frac{T_{2N}^{-1}}{N(N-1)(N-2)} \sum_{i \neq j \neq k \neq l} \mathbf{R}_i^T \mathbf{R}_j \mathbf{R}_k^T \mathbf{R}_l,$$

respectively. Here, \mathbf{R}_i denotes the vectorized form of \mathbf{X}_i obtained by stacking the columns of \mathbf{X}_i on top of one another. Following results in Touloumis *et al.* (2014), it can be shown that T_{1N} , T_{2N} and T_{3N} are consistent estimators of $\text{tr}(\Sigma_C)$, $\text{tr}(\Sigma_C^2)$ and $\text{tr}(\Sigma_R^2)$, respectively. Thus, we estimate λ_C with the consistent estimator

$$\widehat{\lambda}_C = \frac{\frac{T_{3N}}{r^2(N-1)} (T_{1N}^2 + T_{2N})}{\frac{T_{3N}}{r^2(N-1)} (T_{1N}^2 + T_{2N}) + T_{2N} - T_{1N}^2/c}$$

and we estimate μ_C with the consistent estimator

$$\widehat{\mu}_C = T_{1N}/c.$$

*to whom correspondence should be addressed

Finally, the shrinkage estimator of the column covariance matrix Σ_C is the plug-in estimator of (1), that is

$$\widehat{\Sigma}_C = (1 - \widehat{\lambda}_C) \frac{1}{Nr} \sum_{i=1}^N \mathbf{Y}_i^T \mathbf{Y}_i + \widehat{\lambda}_C \widehat{\mu}_C \mathbf{I}_C.$$

Next, to estimate the row covariance matrix Σ_R , let

$$\Sigma_R^* = (1 - \lambda_R) \frac{1}{N \text{tr}(\widehat{\Sigma}_C)} \sum_{i=1}^N \mathbf{Y}_i \mathbf{Y}_i^T + \lambda_R \mathbf{I}_r$$

where

$$\lambda_R = \frac{\frac{1}{N-1} \frac{\text{tr}(\widehat{\Sigma}_C^2)}{\text{tr}^2(\widehat{\Sigma}_C)} [r^2 + \text{tr}(\Sigma_R^2)]}{\frac{1}{N-1} \frac{\text{tr}(\widehat{\Sigma}_C^2)}{\text{tr}^2(\widehat{\Sigma}_C)} [r^2 + \text{tr}(\Sigma_R^2)] + \text{tr}(\Sigma_R^2) - r}$$

minimizes the expected value of $\text{tr}[(\Sigma_R^* - \Sigma_R)^2]$. In a similar manner as with Σ_C^* , we can estimate λ_R with the consistent estimator

$$\widehat{\lambda}_R = \frac{\frac{1}{N-1} \frac{T_{2N}}{T_{1N}^2} (r^2 + T_{3N})}{\frac{1}{N-1} \frac{T_{2N}}{T_{1N}^2} (r^2 + T_{3N}) + T_{3N} - r}$$

and then use

$$\widehat{\Sigma}_R = (1 - \widehat{\lambda}_R) \frac{1}{N \text{tr}(\widehat{\Sigma}_C)} \sum_{i=1}^N \mathbf{Y}_i \mathbf{Y}_i^T + \widehat{\lambda}_R \mathbf{I}_r$$

as a shrinkage estimator of the gene covariance matrix Σ_R .

Derivation of the expectations of $\text{tr}[(\Sigma_C^* - \Sigma_C)^2]$ and $\text{tr}[(\Sigma_R^* - \Sigma_R)^2]$ relies on a matrix-variate normal model (Allen and Tibshirani, 2010) assumption for the matrix-valued random variables $\mathbf{X}_1, \dots, \mathbf{X}_N$. However, under mild moment assumptions outlined in Touloumis *et al.* (2014) and by carrying out similar arguments as in Touloumis (2015), it can be argued that $\widehat{\Sigma}_R$ and $\widehat{\Sigma}_C$ are robust to departures from the normality assumption. Therefore, *HDTD* implements nonparametric estimation and testing procedures that can be used beyond the scope of the matrix-variate normal model.

MULTIPLE TISSUE EXAMPLE

The ‘‘GTEX Analysis V4’’ gene-expression RPKM values were available for downloading from <http://www.gtexportal.org/home> (last assessed on March 19, 2016). The subsamples were then identified using the accompanying annotation files.

The nine tissue samples in each subject-specific data matrix were ordered as skin (sun exposed), nerve, adipose (subcutaneous), artery (tibial), lung, skeletal muscle, heart (left ventricle), blood and thyroid. The estimated tissue covariance matrix $\widehat{\Sigma}_C$ is given in Table 1 and the resulting tissue-wise correlation matrix calculated from $\widehat{\Sigma}_C$ is given in Table 2 (both Tables are presented using a two-letter abbreviation for the tissues). From these two matrices, we can infer the tissue-specific variability by inspecting the diagonal elements of $\widehat{\Sigma}_C$ and the correlation of all possible tissue pairs. For example, the estimated variances for the lung and the skeletal muscle are 18244.2 and 27828.3 respectively, while

their correlation is 0.13, that is the element (5,6) in Table 2. By comparing the diagonal elements of $\widehat{\Sigma}_C$, we can deduce that blood was the most variable tissue ($SE = \sqrt{757511.9} = 870.4$).

Table 3 displays the adjusted p -values that are greater than 0.05 when testing the sphericity hypothesis for all possible pairs of tissues. Although we rejected the sphericity hypothesis for the majority of tissue pairs, it is worth noting that the skin, the adipose and the lung tissue appear to form a network of mutually uncorrelated tissues with equal variance (p -value= 0.082).

Table 4 displays the adjusted p -values when testing the statistical significance of the tissue-specific gene-lists provided by Melé *et al.* (2015). For a given tissue, small adjusted p -values across the remaining eight tissues suggest that the differential expression of the gene-list under study is tissue-specific. From our analysis, the tissue-specificity of the artery-, adipose- and thyroid-specific gene lists provided by Melé *et al.* (2015) was not supported.

The analysis in the main text and the Supplementary Materials can be reproduced by executing the commands in the R script.

REFERENCES

- Allen, G.I. and Tibshirani, R. (2010). Transposable regularized covariance models with an application to missing data imputation. *The Annals of Applied Statistics*, **4**, 764–790.
- Melé, M., *et al.* (2015). The human transcriptome across tissues and individuals. *Science*, **348**, 660–665.
- Touloumis, A., Marioni, J.C. and Tavaré, S. (2014). Hypothesis testing for the covariance matrix in high-dimensional transposable data with Kronecker product dependence structure. *Submitted. arXiv:1404.7684v2*.
- Touloumis, A., Tavaré, S. and Marioni, J.C. (2015). Testing the mean matrix in high-dimensional transposable data. *Biometrics*, **71**, 157–166.
- Touloumis, A. (2015). Nonparametric Stein-type shrinkage covariance matrix estimators in high-dimensional settings. *Computational Statistics and Data Analysis*, **83**, 251–261.

Table 1. The estimated tissue-wise covariance matrix $\hat{\Sigma}_C$.

	SK	NE	AD	AR	LU	SM	HE	BL	TH
SK	19122.9	576.1	263.6	519.7	543.9	1176.7	2995.3	208.3	1159.8
NE	576.1	16882.4	554.0	860.1	751.6	1847.0	2723.1	138.3	1050.0
AD	263.6	554.0	17523.0	461.8	416.6	741.5	1286.6	159.2	617.4
AR	519.7	860.1	461.8	17070.6	796.2	1904.7	2442.6	300.5	1190.7
LU	543.9	751.6	416.6	796.2	18244.2	2824.8	5641.3	3497.2	1963.9
SM	1176.7	1847.0	741.5	1904.7	2824.8	27828.3	13027.8	348.6	3589.7
HE	2995.3	2723.1	1286.6	2442.6	5641.3	13027.8	53880.8	1231.6	9565.9
BL	208.3	138.3	159.2	300.5	3497.2	348.6	1231.6	757511.9	4730.4
TH	1159.8	1050.0	617.4	1190.7	1963.9	3589.7	9565.9	4730.4	20444.9

Table 2. The estimated tissue-wise correlation matrix based on $\hat{\Sigma}_C$.

	SK	NE	AD	AR	LU	SM	HE	BL	TH
SK	1.00	0.03	0.01	0.03	0.03	0.05	0.09	0.00	0.06
NE	0.03	1.00	0.03	0.05	0.04	0.09	0.09	0.00	0.06
AD	0.01	0.03	1.00	0.03	0.02	0.03	0.04	0.00	0.03
AR	0.03	0.05	0.03	1.00	0.05	0.09	0.08	0.00	0.06
LU	0.03	0.04	0.02	0.05	1.00	0.13	0.18	0.03	0.10
SM	0.05	0.09	0.03	0.09	0.13	1.00	0.34	0.00	0.15
HE	0.09	0.09	0.04	0.08	0.18	0.34	1.00	0.01	0.29
BL	0.00	0.00	0.00	0.00	0.03	0.00	0.01	1.00	0.04
TH	0.06	0.06	0.03	0.06	0.10	0.15	0.29	0.04	1.00

Table 3. List of adjusted p -values > 0.05 from applying the sphericity test to all possible tissue pairs.

Tissue 1	Tissue 2	Adjusted p -value
skin	adipose	0.2837
skin	lung	0.5914
skin	thyroid	0.1099
nerve	adipose	0.0588
adipose	artery	0.7992
adipose	lung	0.2724

Table 4. The adjusted p -values from testing the significance of tissue-specific gene lists (rows) to each one of the remaining tissues (columns).

Tissue-Specific Gene List	SK	NE	AD	AR	LU	SM	HE	BL	TH
SK									
NE	0.000								
AD	0.105	0.043							
AR	0.668	0.716	0.007						
LU	0.037	0.000	0.000	0.001					
SM	0.000	0.000	0.000	0.000	0.000				
HE	0.000	0.000	0.000	0.000	0.000	0.000			
BL	0.000	0.000	0.000	0.000	0.000	0.000	0.000		
TH	0.000	0.000	0.000	0.000	0.000	0.781	0.000	0.000	