

Graphical meta analysis: Estimation of covariance matrices from multiple studies

Anders Ellern Bilgrau
abilgrau@math.aau.dk

Poul Svante Eriksen
svante@math.aau.dk

Martin Bøgsted
m_boegsted@dcm.aau.dk

September 24, 2014

Abstract

This paper proposes a model and estimators for a common covariance matrix in cases where multiple datasets are present and thus provide a basis for meta analysis of gaussian graphical models. Our approach is inspired by traditional meta analysis using random effects models. We derive the basic properties and estimators of the model and compare our estimators method to the straight-forward approaches of simple averages or “mega analysis” were the datasets are combined, reprocessed, and plugged into the usual estimators. Though only a modest improvement, explicitly accounting for the inter-study variance is superior to the alternative. The model is also shown to be applicable as an intermediate between linear discriminant analysis (LDA) and quadratic discriminant analysis (QDA).

Keywords: *meta analysis, covariance estimation, integrative analysis, network integration, structural meta analysis, gaussian graphical modeling, linear discriminant analysis, quadratic discriminant analysis*

Contents

1	Introduction	3
2	A graphical random effects model	3
2.1	The likelihood function	4
2.2	Moment estimate	5
2.3	Maximization of the likelihood	5
2.4	Maximization using the EM algorithm	6
2.5	Estimation procedure	7
3	Numerical results	7
4	Applications	9
4.1	Supervised classification	9
4.2	DLBCL graphical meta analysis	10
5	Concluding remarks	10
A	Marginalization of Sigma	12
B	Non-concavity of log-likelihood	12
B.1	log-convexity of the multivariate gamma function	13
B.2	One-dimensional case	14
C	Negative semi-definite hessian in stationary points	15
C.1	1. order derivatives	15
C.2	2. order derivatives	15
C.3	Negative semi-definiteness in stationary points	16
D	log-likelihood of the precision	17
E	Model and notation overview	17

1 Introduction

The fundamental problem in statistics of accurately and precisely estimating the covariance matrix (or its inverse) is notoriously difficult though computationally easy. The usual bias-corrected maximum likelihood estimator (MLE), the sample covariance matrix, have long been known to perform poorly in general due to high variability [3]. The sample covariance is ill-conditioned when the sample size n is less than the number of variables p . Because of its central statistical role the list of statistical methods and applications utilizing the estimated covariance matrix is extraordinarily long. Beside the many standard statistical methods such as principal component analysis (PCA), linear discriminant analysis (LDA), quadratic discriminant analysis (QDA), more direct applications include [examples in medicine, biology, genetics, finance, forensics, physics, engineering, signal processing]. Among this expanding list of applications is also an increasing number of high-dimensional applications and datasets publicly available at online repositories.

In high-dimensional datasets the number of features p far exceed the number of samples n . Since the number of parameters increase quadratically in p and the sample covariance become singular when $p > n$ a plethora of shrinkage and regularization estimators have been proposed to combat the accompanying problems by effectively increasing the degrees of freedom. Instead of attempting to derive still more sophisticated estimators with diminishing improvements we attempt to alleviate the problem from a different angle by effectively increasing n and using more available data.

We are motivated by gene-gene interaction networks in diffuse large B-cell lymphoma (DLBCL) where the covariance matrix contain all information about the conditional dependencies of the genes. As with many other cancers, a large number of DLBCL studies are now available online and hence we wish to use these studies in combination with our own data to arrive at a good estimate of the covariance matrix and the inter-study variation.

2 A graphical random effects model

In a vein similar to regular effect-based meta analysis of Choi et al. [2], we think of the the different studies as related but perturbed experiments. The graphical analog to Choi et al. [2] is the following relatively simple graphical random effects model (GREM) of the observations. Let p be the number of features and k the number of studies. We model the each sample for the i th study as a p -dimensional zero-mean multivariate gaussian vector with the covariance matrix realized from a inverse Wishart distribution, i.e. the hierarchical model,

$$\begin{aligned}\Sigma_i | \Psi, \nu &\sim \mathcal{W}^{-1}(\Psi, \nu), \\ \mathbf{x} | \Sigma_i &\sim \mathcal{N}_p(\mathbf{0}_p, \Sigma_i), \quad i = 1, \dots, k,\end{aligned}\tag{1}$$

where $\mathcal{W}^{-1}(\Psi, \nu)$ denotes an inverse Wishart distribution with probability density function (pdf),

$$f(\Sigma_i | \Psi, \nu) = \frac{|\Psi|^{\frac{\nu}{2}}}{2^{\frac{\nu p}{2}} \Gamma_p(\frac{\nu}{2})} |\Sigma_i|^{-\frac{\nu+p+1}{2}} e^{-\frac{1}{2} \text{tr}(\Psi \Sigma_i^{-1})},$$

Ψ and Σ_i are positive semi-definite, $\nu > p - 1$, and $\mathcal{N}_p(\mathbf{0}_p, \Sigma_i)$ is a multivariate gaussian distribution with pdf

$$f(\mathbf{x} | \Sigma_i) = (2\pi)^{-\frac{p}{2}} \det(\Sigma_i)^{-\frac{1}{2}} e^{-\frac{1}{2} \mathbf{x}^\top \Sigma_i^{-1} \mathbf{x}}.$$

Throughout this paper, we use generic notation $f(\cdot | \cdot)$ and $f(\cdot)$ for the conditional and unconditional pdf of random variables. In the model, ν encodes the inter-study variation where the Σ_i 's concentrate around Ψ for $\nu \rightarrow \infty$. I.e. the inter-study variation goes to zero for larger ν . We wish to infer the parameters Ψ and ν from the observed data.

2.1 The likelihood function

Suppose we have n_i i.i.d. observations, $\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i}$, from the $i = 1, \dots, k$ independent studies from the model given in (1). Let $\mathbf{X}_i = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i})^\top$ be the $n_i \times p$ (transposed) expression matrix of the i 'th study where rows correspond to samples and columns to variables. By the independence assumptions, the log-likelihood is given by

$$\begin{aligned} \ell(\Psi, \nu | \mathbf{X}_1, \dots, \mathbf{X}_k) &= \log f(\mathbf{X}_1, \dots, \mathbf{X}_k | \Psi, \nu) \\ &= \log \int f(\mathbf{X}_1, \dots, \mathbf{X}_k | \Sigma_1, \dots, \Sigma_k, \Psi, \nu) f(\Sigma_1, \dots, \Sigma_k | \Psi, \nu) d\Sigma_1 \cdots d\Sigma_k \\ &= \log \int \prod_{i=1}^k f(\mathbf{X}_i | \Sigma_i) f(\Sigma_i | \Psi, \nu) d\Sigma_1 \cdots d\Sigma_k \\ &= \sum_{i=1}^k \log \int f(\mathbf{X}_i | \Sigma_i) f(\Sigma_i | \Psi, \nu) d\Sigma_i. \end{aligned}$$

Since the inverse Wishart distribution is conjugate to the multivariate gaussian distribution the integral, of which the integrand forms a gaussian-inverse-Wishart distribution, can be evaluated. Hence Σ_i can be marginalized out, cf. Appendix A, and we arrive at the following expression for the log-likelihood,

$$\begin{aligned} \ell(\Psi, \nu | \mathbf{X}_1, \dots, \mathbf{X}_k) &= \sum_{i=1}^k \log \frac{|\Psi|^{\frac{\nu}{2}} \Gamma_p(\frac{\nu+n_i}{2})}{2^{\frac{n_i p}{2}} |\Psi + \mathbf{X}_i^\top \mathbf{X}_i|^{\frac{\nu+n_i}{2}} \Gamma_p(\frac{\nu}{2})} \\ &= c + \sum_{i=1}^k \frac{\nu}{2} \log |\Psi| + \log \Gamma_p\left(\frac{\nu+n_i}{2}\right) - \frac{\nu+n_i}{2} \log |\Psi + \mathbf{X}_i^\top \mathbf{X}_i| - \log \Gamma_p\left(\frac{\nu}{2}\right), \quad (2) \end{aligned}$$

where the constant terms are abbreviated by c and Γ_p is the multivariate generalization for the gamma function Γ , which can be seen in equation (10) of

Appendix B.1. As should be expected, the scatter matrix $\mathbf{S}_i = \mathbf{X}_i^\top \mathbf{X}_i$ and n_i are sufficient statistics for each study.

While we are not able to shown general log-concavity of the log-likelihood, we show that it is log-concave in ν for fixed $\boldsymbol{\Sigma}$, and hence there exists a unique maxima for the marginal $\ell(\nu)$, cf. Appendix B. Reversely, for fixed ν we can show, that the hessian in any stationary point (where $\frac{\partial \ell}{\partial \boldsymbol{\Psi}} = 0$) is negative semi-definite and hence a local maxima, cf. Appendix C. This combined with the observation that $\ell(\boldsymbol{\Psi}) \rightarrow -\infty$ whenever an eigenvalue $\lambda_i(\boldsymbol{\Psi}) \rightarrow 0$, we conjecture that any stationary point indeed a global maximum.

In the following, estimation of $\boldsymbol{\Psi}$ for fixed ν is considered.

2.2 Moment estimate

The moment estimate of $\boldsymbol{\Psi}$ is readily available. It is clear, by construction that the expectation of a single scatter matrix is

$$\mathbb{E}[\mathbf{S}_i] = \mathbb{E}[\mathbf{X}_i^\top \mathbf{X}_i] = \mathbb{E}[\mathbb{E}[\mathbf{X}_i^\top \mathbf{X}_i | \boldsymbol{\Sigma}_i]] = \mathbb{E}[n_i \boldsymbol{\Sigma}_i] = n_i \mathbb{E}[\boldsymbol{\Sigma}_i] = n_i \frac{\boldsymbol{\Psi}}{\nu - p - 1}.$$

Thus the estimator

$$\hat{\boldsymbol{\Psi}} = \frac{1}{k} \sum_{i=1}^k \frac{\nu - p - 1}{n_i} \mathbf{X}_i^\top \mathbf{X}_i = (\nu - p - 1) \frac{1}{k} \sum_{i=1}^k \frac{1}{n_i} \mathbf{S}_i \quad (3)$$

of $\boldsymbol{\Psi}$ is unbiased. This estimate is a scaled average of the empirical covariance matrices.

2.3 Maximization of the likelihood

To find the maximizing parameters we differentiate (2) w.r.t. $\boldsymbol{\Psi}$ and equate to zero while assuming ν known and constant. Assuming $\boldsymbol{\Psi}$ unstructured, then $\nabla_{\mathbf{Z}} \log |\mathbf{Z}| = \mathbf{Z}^{-1}$ and

$$\begin{aligned} \mathbf{0} &= \frac{k\nu}{2} \boldsymbol{\Psi}^{-1} - \sum_{i=1}^k \frac{\nu + n_i}{2} (\boldsymbol{\Psi} + \mathbf{S}_i)^{-1} \\ &= \frac{k\nu}{2} \boldsymbol{\Psi}^{-1} - \sum_{i=1}^k \frac{\nu + n_i}{2} (\mathbf{I} + \boldsymbol{\Psi}^{-1} \mathbf{S}_i)^{-1} \boldsymbol{\Psi}^{-1}. \end{aligned} \quad (4)$$

While assuming $\boldsymbol{\Psi}$ unstructured is not strictly correct it (in this case) leads to a valid equation. The proper 1. order derivative can in alternative notation be seen in equation (13) and can easily be seen to be equivalent to the above. Equation (4) is equivalent to

$$k\nu \mathbf{I} - \sum_{i=1}^k (\nu + n_i) (\mathbf{I} - (-\boldsymbol{\Psi}^{-1} \mathbf{S}_i))^{-1} = \mathbf{0}.$$

which can be rewritten as

$$k\nu\mathbf{I} - \sum_{i=1}^k(\nu + n_i) \sum_{l=0}^{\infty} (-\Psi^{-1}\mathbf{S}_i)^l = \mathbf{0}.$$

by the Neumann series $\left((\mathbf{I} + \mathbf{A})^{-1} = \sum_{l=0}^{\infty} \mathbf{A}^l\right)$ if $\lim_{l \rightarrow \infty} (\mathbf{I} - \Psi^{-1}\mathbf{S}_i)^l = \mathbf{0}$ for all i , i.e. the eigenvalues of $\Psi^{-1}\mathbf{S}_i$ should be less than 1.

Thus, we can approximate the solution by using the first order expansion ($l = 1$) and solve for Ψ

$$\begin{aligned} \mathbf{0} &= k\nu\mathbf{I} - \sum_{i=1}^k(\nu + n_i)(\mathbf{I} - \Psi^{-1}\mathbf{S}_i) \\ &= k\nu\mathbf{I} - k\nu\mathbf{I} - n_{\bullet}\mathbf{I} + \Psi^{-1} \sum_{i=1}^k \nu\mathbf{S}_i + \Psi^{-1} \sum_{i=1}^k n_i\mathbf{S}_i \\ &= -n_{\bullet}\mathbf{I} + \Psi^{-1} \sum_{i=1}^k (\nu + n_i)\mathbf{S}_i \end{aligned}$$

where $n_{\bullet} = \sum_{i=1}^k n_i$ is the total number of observations. This implies

$$\Psi^{-1} \sum_{i=1}^k (\nu + n_i)\mathbf{S}_i = n_{\bullet}\mathbf{I}$$

which is equivalent to

$$\hat{\Psi} = \frac{\sum_{i=1}^k (\nu + n_i)\mathbf{S}_i}{n_{\bullet}} = \frac{\sum_{i=1}^k (\nu + n_i)\mathbf{X}_i^{\top} \mathbf{X}_i}{\sum_{i=1}^k n_i}, \quad (5)$$

corresponding to a weighted sum of the scatter matrices.

2.4 Maximization using the EM algorithm

We now derive the updating scheme of the expectation-maximization (EM) algorithm for fixed ν . Suppose $\hat{\Psi}_{(t)}$ is the current estimate of Ψ and that $\mathbf{S}_i = \mathbf{X}_i^{\top} \mathbf{X}_i = \sum_{j=1}^{n_i} \mathbf{x}_{ij} \mathbf{x}_{ij}^{\top}$ is the empirical scatter matrix where $\mathbf{X}_i = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i})^{\top}$. We now compute the expectation step of the EM-algorithm.

From (1) we have that,

$$\begin{aligned} \mathbf{S}_i | \Sigma_i &\sim \mathcal{W}_p(\Sigma_i, n_i), \quad i = 1, \dots, k \\ \Sigma_i &\sim \mathcal{W}_p^{-1}(\Psi, \nu) \end{aligned}$$

Let $\Delta_i = \Sigma_i^{-1}$ be the precision matrix (or, concentration matrix) and let $\Theta = \Psi^{-1}$, then

$$\begin{aligned} \mathbf{S}_i | \Delta_i &\sim \mathcal{W}_p(\Delta_i^{-1}, n_i) &\Leftrightarrow & f(\mathbf{S}_i | \Delta_i) \propto |\Delta_i|^{\frac{1}{2}} e^{-\frac{1}{2} \text{tr}(\Delta_i \mathbf{S}_i)}, \\ \Delta_i &\sim \mathcal{W}_p(\Theta, \nu) &\Leftrightarrow & f(\Delta_i) \propto |\Theta|^{-\frac{\nu}{2}} e^{-\frac{1}{2} \text{tr}(\Theta^{-1} \Delta_i)}. \end{aligned} \quad (6)$$

From the conjugacy of the inverse Wishart and the Wishart distribution, we have the posterior distribution for the precision matrix,

$$\mathbf{\Delta}_i | \mathbf{S}_i \sim \mathcal{W}_p((\mathbf{\Theta}^{-1} + \mathbf{S}_i)^{-1}, n_i + \nu)$$

Hence, the expectation step is given by

$$\mathbb{E}[\mathbf{\Delta}_i | \mathbf{S}_i] = (n_i + \nu)(\mathbf{\Theta}^{-1} + \mathbf{S}_i)^{-1}.$$

The maximization step, in which the log-likelihood $\ell(\mathbf{\Theta} | \mathbf{\Delta}_1, \dots, \mathbf{\Delta}_k)$ is maximized, yields the estimate

$$\hat{\mathbf{\Theta}} = \frac{1}{k\nu} \sum_{i=1}^k \mathbf{\Delta}_i,$$

which is the mean of the scaled precision matrices $\frac{1}{\nu} \mathbf{\Delta}_i$. The derivation of this estimate can be seen in Appendix D. The above yield the updating scheme

$$\hat{\mathbf{\Theta}}_{(t+1)} = \frac{1}{k\nu} \sum_{i=1}^k (n_i + \nu) \left((\hat{\mathbf{\Theta}}_{(t)})^{-1} + \mathbf{S}_i \right)^{-1} \quad (7)$$

for $\mathbf{\Theta}$. The connection to the maximum likelihood estimate is immediately seen through equation (4).

2.5 Estimation procedure

We propose an alternating procedure between estimating ν and $\mathbf{\Psi}$ while keeping the other fixed. Given parameters $\hat{\nu}_{(t)}$ and $\hat{\mathbf{\Psi}}_{(t)}$ at iteration t , we estimate $\hat{\mathbf{\Psi}}_{(t+1)}$ using the fixed $\hat{\nu}_{(t+1)}$. Subsequently, find $\hat{\nu}_{(t+1)}$ by standard one-dimensional numerical optimization procedure keeping $\hat{\mathbf{\Psi}}_{(t)}$ fixed. This coordinate ascent-like approach is repeated until convergence. More precisely, in pseudo-code, we propose the algorithm seen in Algorithm 1. The update function U in the algorithm is defined by the derived estimators above. That is, equations (3), (5), and (7) define U to be the moment, approximate MLE, and EM estimate respectively.

3 Numerical results

To assess the model and the stability of the estimation procedure we generated data from (1) for $p = 10$ variables in $k = 3$ studies each with an equal number of observations, $n = n_1 = n_2 = n_3$. We chose the parameters $\nu = 15$ and

$$\mathbf{\Psi} = \begin{bmatrix} 1 & 0.5 & 0.5 \\ 0.5 & 1 & 0.5 \\ 0.5 & 0.5 & 1 \end{bmatrix}.$$

Algorithm 1 Pseudo-code for the GREM estimation procedure

```

1: procedure GREM COORDINATE ASCENT
2:   Input:
3:   Sufficient data:  $(\mathbf{S}_1, n_1), \dots, (\mathbf{S}_k, n_k)$ 
4:   Initial parameters:  $\hat{\Psi}_{(0)}, \hat{\nu}_{(0)}$ 
5:   Convergence criterion:  $\varepsilon > 0$ 
6:   Output:
7:   Parameters:  $\hat{\Psi}_{(t')}, \hat{\nu}_{(t')}$ 
8:
9:   Initialize:  $l_{(0)} \leftarrow \ell(\hat{\Psi}_{(0)}, \hat{\nu}_{(0)})$ 
10:  for  $t = 1, 2, 3, \dots$  do
11:     $\hat{\Psi}_{(t)} \leftarrow U(\hat{\Psi}_{(t-1)}, \hat{\nu}_{(t-1)})$ 
12:     $\nu_{(t)} \leftarrow \arg \max_{\nu} \ell(\hat{\Psi}_{(t)}, \nu)$ 
13:     $l_{(t)} \leftarrow \ell(\hat{\Psi}_{(t)}, \hat{\nu}_{(t)})$ 
14:    if  $l_{(t)} - l_{(t-1)} < \varepsilon$  then
15:      return  $(\hat{\Psi}_{(t)}, \nu_{(t)})$ 

```

The number of observations in each study, n_i , was varied range [5, 14].

We measure the precision of the estimates values against the expected covariance matrix given by

$$\Sigma = \mathbb{E}[\Sigma_i] = \frac{1}{\nu - p - 1} \Psi.$$

Let $\hat{\Psi}$ and $\hat{\nu}$ be the estimates obtained in the model described and let

$$\hat{\Sigma}_{\text{GREM}} = \frac{1}{\hat{\nu} - p - 1} \hat{\Psi}$$

be an estimate of Σ . We benchmark this against the pooled covariance matrix,

$$\hat{\Sigma}_{\text{pool}} = \frac{1}{k} \sum_{i=1}^k \frac{1}{n_i} S_i,$$

as a simple alternative estimate of Σ . We benchmark the two against each other using the following sum of squared errors,

$$\text{SSE}(\hat{\Sigma}) = \sum_{i \leq j} \frac{(\hat{\Sigma}_{ij} - \Sigma_{ij})^2}{\text{var}(\Sigma_{ij})}$$

where

$$\text{var}(\Sigma_{ij}) = n(\Psi_{ij}^2 + \Psi_{ii}\Psi_{jj}).$$

For each n from 5 to 14, the squared sum of errors for each estimator, $\text{SSE}(\hat{\Sigma}_{\text{GREM}})$, and $\text{SSE}(\hat{\Sigma}_{\text{pool}})$ where computed for 500 datasets and the average of these values seen in figure 1 as function of n . Unsurprisingly, the model outperforms the simple pool average covariance matrix.

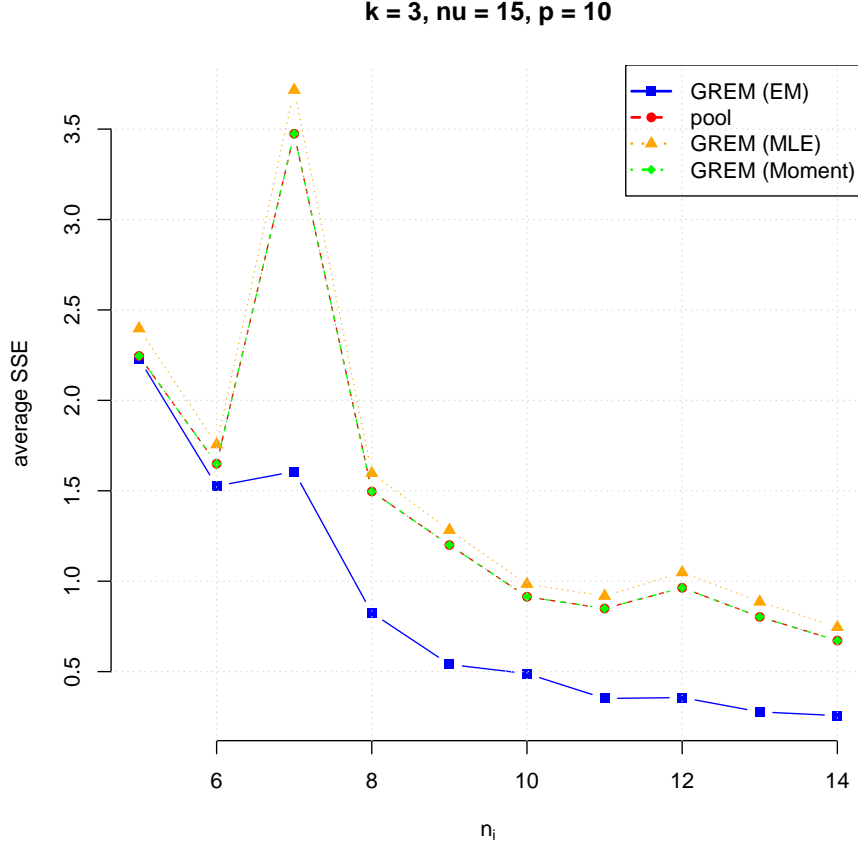


Figure 1: The average sum or squared errors (SSE), of 500 simulations, as a function of the number of samples n_i in each study.

4 Applications

4.1 Supervised classification

The estimate obtained from the model can be utilized in supervised learning as a intermediate case of linear discriminant analysis (LDA) and quadratic dis-

criminant analysis (QDA) to regularized linear discriminant analysis (RDA).

Let Y be a random variable denoting the class $1, \dots, k$ and suppose \mathbf{x} is a random vector of the explanatory variables. Recall, that QDA (and LDA) finds the class y maximizing

$$P(Y = y | \mathbf{X} = \mathbf{x}) = \frac{\pi_y f(\mathbf{x} | Y = y)}{\sum_{y'=1}^k \pi_{y'} f(\mathbf{x} | Y = y')}$$

where $\mathbf{x} | Y = y$ is assumed multivariate normal, i.e.

$$(\mathbf{x} | Y = y) \sim \mathcal{N}_p(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y).$$

LDA differs from QDA only by the additional assumption that $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_y$ for all classes y . An intermediate classifier of LDA and QDA can thus be constructed by assuming the $\boldsymbol{\Sigma}_y$'s inversely Wishart distributed as in (1), i.e. $\boldsymbol{\Sigma}_y \sim \mathcal{W}^{-1}(\boldsymbol{\Psi}, \nu)$. This hierarchical discriminant analysis (HDA) is thus straight-forward to implement given that

$$\begin{aligned} f(\mathbf{x} | Y = y) &= \int f(\mathbf{x} | \boldsymbol{\Sigma}_y, Y = y) f(\boldsymbol{\Sigma}_y | Y = y) d\boldsymbol{\Sigma}_y \\ &= \frac{|\boldsymbol{\Psi}|^{\frac{\nu}{2}} \Gamma_p\left(\frac{\nu+1}{2}\right)}{\pi^{-\frac{n}{2}} |\boldsymbol{\Psi} + (\mathbf{x} - \boldsymbol{\mu}_k)(\mathbf{x} - \boldsymbol{\mu}_k)^\top| \frac{\nu+1}{2} \Gamma_p\left(\frac{\nu}{2}\right)}, \end{aligned}$$

analogous to the computation done in appendix A, using the matrix determinant lemma, $|\mathbf{A} + \mathbf{u}\mathbf{v}^\top| = (1 + \mathbf{v}^\top \mathbf{A}\mathbf{u})|\mathbf{A}|$, to simplify and speed up the computations \square .

4.2 DLBCL graphical meta analysis

5 Concluding remarks

While the improvements are modest, the results above demonstrate an advantage of modelling the inter-study variance as a hierarchical random effects model. However, the virtue of such a model is not from improvement in accuracy alone. Also desirable is the explicit quantification of the inter-study variance. If $\hat{\nu}$ is estimated to be large, the studies exhibit a largely common covariance structure, and vice-versa when $\hat{\nu}$ is small.

The generalization of the model to $n \gg p$ is extremely interesting though out of scope for this article.

References

- [1] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, 2004. ISBN 9780511804441. doi: 10.1017/CBO9780511804441.
- [2] J. K. Choi, U. Yu, S. Kim, and O. J. Yoo. Combining multiple microarray studies and modeling interstudy variation. *Bioinformatics*, 19(Suppl 1):i84–i90, July 2003. ISSN 1367-4803. doi: 10.1093/bioinformatics/btg1010.
- [3] AP Dempster. Covariance Selection. *Biometrics*, 28(1):157–175, 1972.
- [4] KB Petersen and MS Pedersen. The Matrix Cookbook. *Technical University of Denmark*, 2008. URL http://www2.imm.dtu.dk/pubdb/views/publication_details.php?id=3274.

A Marginalization of Σ

This sections shows the marginalization out of Σ in (2). For ease of notation we drop the subscript i used in Σ_i , \mathbf{X}_i , $\mathbf{S}_i = \mathbf{X}_i \mathbf{X}_i^\top$, and n_i in the above text. We do the computation straight-forwardly by the assumptions of the model,

$$\begin{aligned}
f(\mathbf{X}|\Psi, \nu) &= \int f(\mathbf{X}|\Sigma) f(\Sigma|\Psi, \nu) d\Sigma \\
&= \int \left[\prod_{j=1}^n (2\pi)^{-\frac{p}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2} \text{tr}(\mathbf{x}_j \mathbf{x}_j^\top \Sigma^{-1})} \right] \frac{|\Psi|^{-\frac{p}{2}}}{2^{\frac{\nu p}{2}} \Gamma_p(\frac{\nu}{2})} |\Sigma|^{-\frac{\nu+p+1}{2}} e^{-\frac{1}{2} \text{tr}(\Psi \Sigma^{-1})} d\Sigma \\
&= (2\pi)^{-\frac{np}{2}} \frac{|\Psi|^{-\frac{p}{2}}}{2^{\frac{\nu p}{2}} \Gamma_p(\frac{\nu}{2})} \int |\Sigma|^{-\frac{n}{2}} e^{-\frac{1}{2} \text{tr}(\mathbf{S} \Sigma^{-1})} |\Sigma|^{-\frac{\nu+p+1}{2}} e^{-\frac{1}{2} \text{tr}(\Psi \Sigma^{-1})} d\Sigma \\
&= \frac{|\Psi|^{-\frac{p}{2}}}{\pi^{-\frac{np}{2}} 2^{\frac{(\nu+n)p}{2}} \Gamma_p(\frac{\nu}{2})} \int |\Sigma|^{-\frac{(\nu+n)+p+1}{2}} e^{-\frac{1}{2} \text{tr}((\Psi+\mathbf{S}) \Sigma^{-1})} d\Sigma.
\end{aligned}$$

The integrand here can be recognized as a unnormalized inverse Wishart pdf, $\mathcal{W}^{-1}(\Psi + \mathbf{S}, \nu + n)$, and so integral evaluates to the reciprocal value of the normalizing constant in that density. Thus,

$$\begin{aligned}
f(\mathbf{X}|\Psi, \nu) &= \frac{|\Psi|^{-\frac{p}{2}}}{\pi^{-\frac{np}{2}} 2^{\frac{(\nu+n)p}{2}} \Gamma_p(\frac{\nu}{2})} \frac{2^{\frac{(\nu+n)p}{2}} \Gamma_p(\frac{\nu+n}{2})}{|\Psi + \mathbf{S}|^{\frac{\nu+n}{2}}} \\
&= \frac{|\Psi|^{-\frac{p}{2}} \Gamma_p(\frac{\nu+n}{2})}{\pi^{-\frac{np}{2}} |\Psi + \mathbf{S}|^{\frac{\nu+n}{2}} \Gamma_p(\frac{\nu}{2})},
\end{aligned}$$

which was what was wanted.

B Non-concavity of log-likelihood

The log-likelihood is not concave. This section analyse the (non)-concavity of the log-likelihood function,

$$\begin{aligned}
\ell(\Psi, \nu | \mathbf{X}_1, \dots, \mathbf{X}_k) \\
= c + \sum_{i=1}^k \frac{\nu}{2} \log |\Psi| + \log \Gamma_p\left(\frac{\nu + n_i}{2}\right) - \frac{\nu + n_i}{2} \log |\Psi + \mathbf{S}_i| - \log \Gamma_p\left(\frac{\nu}{2}\right). \quad (8)
\end{aligned}$$

We first analyze the terms involving Ψ in (8). Not counting the constant, the first term, $\frac{\nu}{2} \log |\Psi|$, is concave by the well known result that the log-determinant of positive semi-definite matrix is concave [See e.g. 1, pp. 73-74]. The third term in (8), $-\frac{\nu+n_i}{2} \log |\Psi + \mathbf{S}_i|$, is convex by the same argument however negated by the sign. The sum $\Psi + \mathbf{S}_i$ is positive semi-definite since both summands are.

Next, we look at the terms involving ν . Clearly, the mixed terms involving both ν and Ψ are log-linear in ν and hence log-concave. We thus restrict our attention to the remaining terms not dependent on Ψ . By themselves,

the second term in (8), $\log \Gamma_p\left(\frac{\nu+n_i}{2}\right)$ is convex since the multivariate gamma function is log-convex, cf. section B.1. In the same manner, the fourth term, $-\log \Gamma_p\left(\frac{\nu}{2}\right)$, is concave by the negative sign. The sum however of these terms involving Γ_p are concave in ν , since

$$\log \Gamma_p\left(\frac{\nu+n_i}{2}\right) - \log \Gamma_p\left(\frac{\nu}{2}\right) = \log \frac{\Gamma_p\left(\frac{\nu+n_i}{2}\right)}{\Gamma_p\left(\frac{\nu}{2}\right)} = \log \prod_{j=1}^p \frac{\Gamma\left(\frac{\nu+1-j}{2} + \frac{n_i}{2}\right)}{\Gamma\left(\frac{\nu+1-j}{2}\right)}.$$

which can be seen to be concave since $n_i \geq 2$ for all i and

$$x \mapsto \log \left(\frac{\Gamma(x+a)}{\Gamma(x)} \right) \quad (9)$$

is concave for all $x > 0$ and $a > 0$. While we do not provide a proof of this the mapping is plotted in figure 2. Hence, the log-likelihood is log-concave in ν .

```
logGammaRatio <- function(x, a) {
  log(gamma(x + a)/gamma(x))
}
xs <- seq(0.01, 2, by = 0.01)
par(mfrow = c(1,2), mar = c(2,2,0,0) + 0.5)
plot(xs, logGammaRatio(xs, a = 1e-3), type = "l", xlab = "", ylab = "")
```

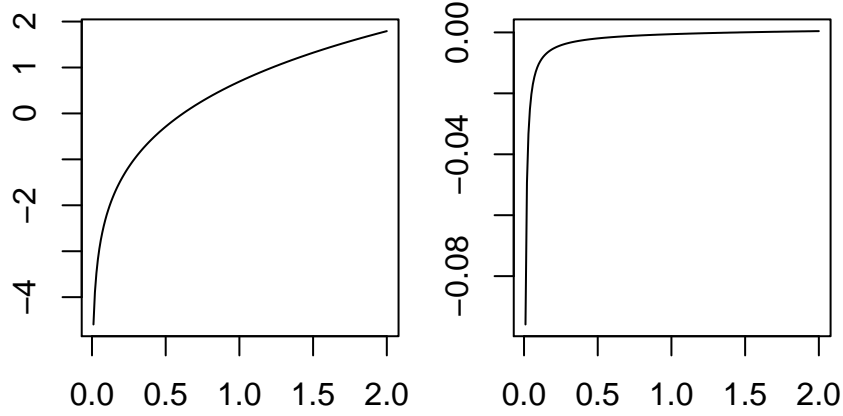


Figure 2: Plots of the mapping given in (9) for different values of a .

B.1 log-convexity of the multivariate gamma function

The log-convexity of the multivariate gamma function can be seen using the following characterization of Γ_p ,

$$\Gamma_p(t) = \pi^{\frac{1}{2}\binom{p}{2}} \prod_{j=1}^p \Gamma\left(t + \frac{1-j}{2}\right) \quad \text{where } \Gamma(t) = \int_0^\infty x^{t-1} e^{-x} dx. \quad (10)$$

From this

$$\log \Gamma_p(t) = \frac{1}{2} \binom{p}{2} \log \pi + \sum_{j=1}^p \Gamma \left(t + \frac{1-j}{2} \right), \quad (11)$$

which is convex since Γ is log-convex and sums of convex functions are convex. Hence Γ_p is log-convex.

B.2 One-dimensional case

In the one-dimensional case, we have

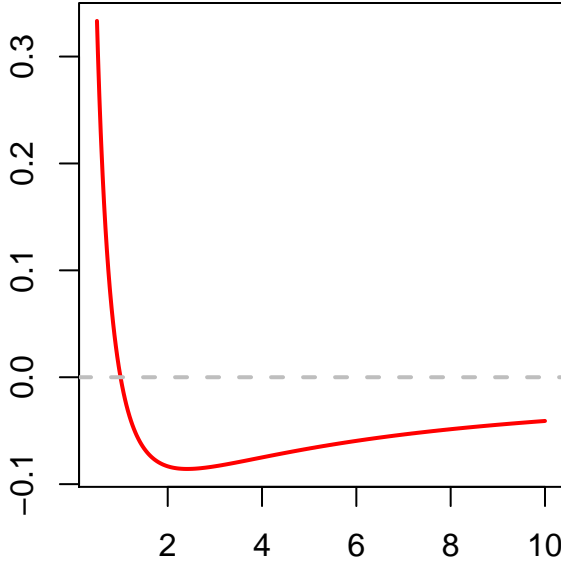
$$\ell'(\phi) = \frac{k\nu}{2} \frac{1}{\phi} + \sum_{i=1}^k \frac{\nu + n_i}{2} \frac{1}{\phi + x_i^2}$$

which is not clearly convex. We see, that

$$\lim_{\phi \rightarrow 0} \ell'(\phi) = \infty \text{ and } \lim_{\phi \rightarrow \infty} \ell'(\phi) = 0$$

We draw this derivative with $k = 1$ and other appropriately chosen parameters:

```
dl <- function(phi, k = 1, nu = 1, ni = 1, xi = 1) {
  k*nu/2*1/phi - (nu + ni)/2 * 1/(phi + xi^2)
}
phi <- seq(0.5, 10, by = 0.01)
plot(phi, dl(phi), type = "l", col = "red", lwd = 2)
abline(h = 0, col = "grey", lty = 2, lwd = 2)
```



Which is be seen to have a unique root. However, the log-likelihood ℓ is not log-concave since a differentiable function of one variable is concave if and only if its derivative is monotonically non-increasing.

C Negative semi-definite hessian in stationary points

This section proves that the hessian is negative semi-definite in all stationary points. The log-likelihood in (2), assuming ν fixed, obey

$$2\ell(\Psi) = k\nu \log |\Psi| - \sum_{a=1}^k (n_a + \nu) \log |\Psi + S_a| \quad (12)$$

up to an addition of a constant. The matrix cookbook by Petersen and Pedersen [4] has been a useful reference here¹.

C.1 1. order derivatives

From the log-likelihood expression, we compute the 1. order derivative $\nabla_{\Psi} 2\ell(\Psi)$ which is the matrix-valued function where each entry is given by

$$\frac{\partial 2\ell}{\partial \Psi_{ij}} = k\nu \operatorname{tr}(\mathbf{E}^{ij} \Psi^{-1}) - \sum_{a=1}^k (\nu + n_a) \operatorname{tr}(\mathbf{E}^{ij} (\Psi + S_a)^{-1}). \quad (13)$$

where \mathbf{E}^{ij} is a matrix with ones at entries (i, j) and (j, i) and zeros elsewhere. This \mathbf{E}^{ij} is introduced as the derivative is not straight-forward because of the symmetric structure of Ψ . Had Ψ been unstructured, then $\frac{\partial}{\partial \Psi} \log |\Psi| = \Psi^{-1}$. However, when Ψ is symmetric we have that $\frac{\partial}{\partial \Psi_{ij}} \log |\Psi| = \operatorname{tr}(\mathbf{E}^{ij} \Psi^{-1})$ which is to say $\frac{\partial}{\partial \Psi} \log |\Psi| = 2\Psi^{-1} - \Psi^{-1} \circ \mathbf{I}$ where \circ denotes the Hadamard product [4, eq. (43) and (141)].

The first order derivative lives in a $\binom{p+1}{2}$ -dimensional vector space indexed by (i, j) , $i \leq j$, with basis vectors \mathbf{E}^{ij} .

C.2 2. order derivatives

We proceed with the second order derivative $\nabla_{\Psi}^2 2\ell(\Psi)$ with entries given by

$$\begin{aligned} \frac{\partial^2 2\ell}{\partial \Psi_{kl} \partial \Psi_{ij}} &= -k\nu \operatorname{tr}(\mathbf{E}^{ij} \Psi^{-1} \mathbf{E}^{kl} \Psi^{-1}) \\ &+ \sum_{a=1}^k (\nu + n_a) \operatorname{tr}(\mathbf{E}^{ij} (\Psi + S_a)^{-1} \mathbf{E}^{kl} (\Psi + S_a)^{-1}), \end{aligned}$$

obtained by differentiation of (13) combined with $\frac{\partial}{\partial \Psi_{ij}} \Psi^{-1} = -\Psi^{-1} \mathbf{E}^{ij} \Psi^{-1}$ for symmetric matrices [4, eq. (40)] and the linearity of the trace operator.

The second order derivative is a $\binom{p+1}{2} \times \binom{p+1}{2}$ -dimensional matrix indexed by (i, j) and (k, l) , $i \leq j$, $k \leq l$.

¹See equation (41, p. 8) and (59, p. 9) and pages 14 and 52-53 in http://www2.imm.dtu.dk/pubdb/views/publication_details.php?id=3274

C.3 Negative semi-definiteness in stationary points

With the above expressions we now show that the hessian is negative semi-definite in all stationary points (or, extrema). Let $\mathbf{Y} = \sum_{(i,j)} y_{ij} E_{ij}$ be a symmetric matrix in the vector space where $\mathbf{Y} \neq \mathbf{0}$. The analog to $\mathbf{y}^\top \mathbf{A} \mathbf{y} = y_i \sum_{ij} A_{ij} y_j \leq 0$ in our vector space then becomes

$$\sum_{i \leq j, k \leq l} Y_{ij} (\nabla_{\Psi}^2 2\ell(\Psi))_{(i,j),(k,l)} Y_{kl} \leq 0$$

which amounts to showing that

$$-k\nu \operatorname{tr}(\mathbf{Y} \Psi^{-1} \mathbf{Y} \Psi^{-1}) + \sum_{a=1}^k (\nu + n_a) \operatorname{tr}(\mathbf{Y} (\Psi + \mathbf{S}_a)^{-1} \mathbf{Y} (\Psi + \mathbf{S}_a)^{-1}) \leq 0. \quad (14)$$

Now, the positive-definiteness of Ψ , let

$$\begin{aligned} \mathbf{Y} &:= \Psi^{-\frac{1}{2}} \mathbf{Y} \Psi^{-\frac{1}{2}} \text{ and} \\ \mathbf{S}_a &:= \Psi^{-\frac{1}{2}} \mathbf{S}_a \Psi^{-\frac{1}{2}}. \end{aligned}$$

we can assume that $\Psi = \mathbf{I}$. Hence, the likelihood equation (12) equated to zero, becomes

$$k\nu \mathbf{I} = \sum_a (n_a + \nu) (\mathbf{I} + \mathbf{S}_a)$$

which implies (by multiplication by \mathbf{Y}^2)

$$\begin{aligned} k\nu \operatorname{tr}(\mathbf{Y}^2) &= \sum_a (n_a + \nu) \operatorname{tr}(\mathbf{Y}^2 (\mathbf{I} + \mathbf{S}_a)) \\ &= \sum_a (n_a + \nu) \operatorname{tr}(\mathbf{Y} (\mathbf{I} + \mathbf{S}_a) \mathbf{Y}). \end{aligned} \quad (15)$$

We substitute (15) into (14) to get

$$\begin{aligned} &\sum_a (n_a + \nu) \operatorname{tr}(\mathbf{Y} (\mathbf{I} + \mathbf{S}_a) \mathbf{Y} (\mathbf{I} + \mathbf{S}_a) - \mathbf{Y} (\mathbf{I} + \mathbf{S}_a) \mathbf{Y}) \\ &= \sum_a (n_a + \nu) \operatorname{tr}(\mathbf{Y} (\mathbf{I} + \mathbf{S}_a) \mathbf{Y} [(\mathbf{I} + \mathbf{S}_a) - \mathbf{I}]) \leq 0 \end{aligned}$$

Since \mathbf{S}_a is positive semi-definite, it can be diagonalized $\mathbf{S}_a = \mathbf{U}_a \mathbf{D}_a \mathbf{U}_a^\top$ by the spectral theorem where the diagonal matrix $\mathbf{D}_a \succeq 0$ (all non-zero entries are positive) and \mathbf{U}_a is orthonormal (i.e. $\mathbf{U}_a \mathbf{U}_a^\top = \mathbf{U}_a^\top \mathbf{U}_a = \mathbf{I}$). Using the diagonalization,

$$\sum_a (n_a + \nu) \operatorname{tr}(\underbrace{\mathbf{U}^\top \mathbf{Y} (\mathbf{I} + \mathbf{S}_a)^{-1} \mathbf{Y} \mathbf{U}}_{\text{symmetric}} [(\mathbf{I} + \mathbf{D}_a)^{-1} - \mathbf{I}]) \leq 0 \quad (16)$$

where the underbraced matrix is positive semi-definite and hence have non-negative diagonal elements. Furthermore, since \mathbf{D}_a is diagonal also with non-negative elements, the diagonal matrix $(\mathbf{I} + \mathbf{D}_a)^{-1} - \mathbf{I}$ clearly have non-positive entries and is thus negative semi-definite.

Since the trace of a matrix product is the sum of the element-wise products, the trace (and thus the sum) will always be non-positive and hence (16) will always hold.

D log-likelihood of the precision

Suppose we have k i.i.d. realizations, $\mathbf{\Delta}_1, \dots, \mathbf{\Delta}_k$, from the Wishart distribution given in equation (6). The corresponding log-likelihood can be computed straight-forwardly:

$$\begin{aligned} \ell(\mathbf{\Theta} | \mathbf{\Delta}_1, \dots, \mathbf{\Delta}_k) &= \sum_{i=1}^k \log f(\mathbf{\Delta}_i | \mathbf{\Theta}) \\ &= \sum_{i=1}^k \log \frac{|\mathbf{\Theta}|^{-\frac{\nu}{2}}}{2^{-\frac{\nu p}{2}} \Gamma_p\left(\frac{\nu}{2}\right)} |\mathbf{\Delta}_i|^{\frac{\nu-p-1}{2}} e^{-\frac{1}{2} \text{tr}(\mathbf{\Theta}^{-1} \mathbf{\Delta}_i)} \\ &= c + \sum_{i=1}^k -\frac{\nu}{2} \log |\mathbf{\Theta}| - \frac{1}{2} \text{tr}(\mathbf{\Theta}^{-1} \mathbf{\Delta}_i) \\ &= c + \frac{\nu k}{2} \left(\log |\mathbf{\Theta}| + \text{tr} \left(\mathbf{\Theta}^{-1} \frac{1}{\nu k} \sum_{i=1}^k \mathbf{\Delta}_i \right) \right). \end{aligned}$$

The last expression is to be maximized with respect to $\mathbf{\Theta}$ and can be recognized as the MLE problem in a multivariate Gaussian distribution. Hence,

$$\mathbf{\Theta} = \frac{1}{k\nu} \sum_{i=1}^k \mathbf{\Delta}_i,$$

is the MLE in this model.

E Model and notation overview

Let $\mathbf{\Delta}_i = \mathbf{\Sigma}_i^{-1}$, $\mathbf{\Theta} = \mathbf{\Psi}^{-1}$, and $\mathbf{S}_i = \mathbf{X}_i^\top \mathbf{X}_i$. Then the following equivalences hold.

$$\begin{array}{lll} \mathbf{\Sigma}_i \sim \mathcal{W}_p^{-1}(\mathbf{\Psi}, \nu) & \iff & \mathbf{\Delta}_i \sim \mathcal{W}_p(\mathbf{\Theta}, \nu) \\ \mathbf{X}_i | \mathbf{\Sigma}_i \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{\Sigma}_i) & \iff & \mathbf{X}_i | \mathbf{\Delta}_i \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{\Delta}_i^{-1}) \\ \Updownarrow & & \Updownarrow \\ \mathbf{\Sigma}_i | X_i \sim \mathcal{W}_p^{-1}(\mathbf{\Psi} + \mathbf{S}_i, \nu + n_i) & \iff & \mathbf{\Delta}_i | X_i \sim \mathcal{W}_p((\mathbf{\Theta}^{-1} + \mathbf{S}_i)^{-1}, \nu + n_i) \end{array}$$

List of Corrections