

# How to use DeCovarT: a toy example with two genes and two cell populations

true                      true                      true

## 1 Objectives

### 1.1 Rationale of the new generative model

As in most traditional deconvolution models, we assume that the overall measured gene expression can be reconstructed by summing the individual contributions of each cell population weighted by its frequency. Formally, let  $\mathbf{X} = (x_{gj}) \in \mathcal{M}_{\mathbb{R}^{G \times J}}$  the signature matrix representing the purified transcriptomic profiles of  $J$  cell populations and  $\mathbf{p} = (p_{ji}) \in ]0, 1[^{J \times N}$  the unknown relative proportions of cell populations in  $N$  samples, then the linear relation relating the bulk expression ( $\mathbf{y} = (y_{gi}) \in \mathbb{R}_+^{G \times N}$ ) to the individual cell expression profiles is given by the matrix product:  $\mathbf{y} = \mathbf{X} \times \mathbf{p}$ . In addition, we consider unit simplex constraint on the cellular ratios (Eq. (1)):

$$\begin{cases} \sum_{j=1}^J p_j = 1 \\ \forall j \in \tilde{J} \quad p_j \geq 0 \end{cases} \quad (1)$$

However, in real conditions with technical and environmental variability, the strict linearity of the deconvolution does not strictly hold. Thus, an additional error term is usually added, assumed to follow a *homoscedastic* zero-centred Gaussian distribution and with pairwise independent response measures while the exogenous variables (here, the purified expression profiles) are supposed determined: this set of conditions is referred to as the Gaussian-Markow assumptions. In that configuration, the MLE (maximum likelihood estimate) that best describes this standard linear model is equal to the ordinary least squares (OLS) estimate.

In contrast to this canonical approach, in DeCovarT, we relax the *exogeneity* property by treating exogenous variables  $\mathbf{X}$  as random variables rather than determined measures, in a process close to the approach of the DSection algorithm [1]. However, to our knowledge, we are the first to weaken the independence assumption between observations by explicitly incorporating the intrinsic covariance structure of the transcriptome of each purified cell population. To do so, we conjecture that the  $G$ -dimensional vector  $\mathbf{x}_j$  characterising the transcriptomic expression of each cell population follows a multivariate Gaussian distribution:  $\mathbf{x}_j \sim \mathcal{N}_G(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ , with  $\boldsymbol{\mu}_j$  the mean purified transcriptomic expression and  $\boldsymbol{\Sigma}_j$  the covariance matrix, that we constrain to be positive-definite and of full rank and that is inferred using the output of the gLasso algorithm [2]. We display respectively the graphical models associated to the standard linear deconvolution model and our new innovative generative model used by the DeCovarT algorithm in subfigures a) and b), in Fig.1.

### 1.2 Derivation of the log-likelihood

First, we *plugged-in* the mean and covariance parameters  $\zeta_j = (\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$  inferred in the previous step. Then, by letting  $\boldsymbol{\zeta} = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\boldsymbol{\mu} = (\boldsymbol{\mu}_j)_{j \in \tilde{J}} \in \mathcal{M}_{G \times J}$ ,  $\boldsymbol{\Sigma} \in \mathcal{M}_{G \times G}$  the known parameters and  $\mathbf{p}$  the unknown cellular ratios, the conditional distribution  $\mathbf{y} | (\boldsymbol{\zeta}, \mathbf{p})$  is the convolution of pairwise independent multivariate Gaussian distributions, which is also a multivariate Gaussian distribution (Eq.(2)), deduced from the *affine invariant* property of Gaussian distributions.

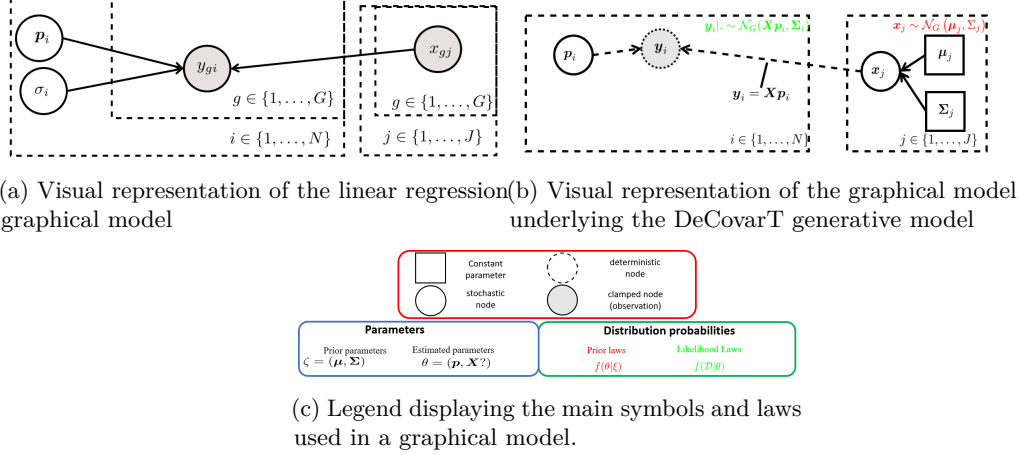


Figure 1: We use the standard graphical convention of graphical models, as depicted in RevBayes webpage. For identifiability reasons, we conjecture that all variability arises from the stochastic nature of the covariates.

$$\mathbf{y}|\zeta, \mathbf{p} \sim \mathcal{N}_G(\mu\mathbf{p}, \Sigma) \text{ with } \mu = (\mu_{.j})_{j \in \tilde{J}}, \quad \mathbf{p} = (p_1, \dots, p_J) \text{ and } \Sigma = \sum_{j=1}^J p_j^2 \Sigma_j \quad (2)$$

From Eq.(2), we readily compute the associated conditional log-likelihood (Eq.(3)):

$$\ell_{\mathbf{y}|\zeta}(\mathbf{p}) = C + \log \left( \text{Det} \left( \sum_{j=1}^J p_j^2 \Sigma_j \right)^{-1} \right) - \frac{1}{2} (\mathbf{y} - \mathbf{p}\mu)^\top \left( \sum_{j=1}^J p_j^2 \Sigma_j \right)^{-1} (\mathbf{y} - \mathbf{p}\mu) \quad (3)$$

### 1.3 First and second-order derivation of the unconstrained DeCovarT log-likelihood function

The stationary points of a function and notably maxima, are given by the roots (the values at which the function crosses the  $x$ -axis) of its gradient, in our context, the vector:  $\nabla \ell : \mathbb{R}^J \rightarrow \mathbb{R}^J$  evaluated at point  $\nabla \ell(\mathbf{p}) : ]0, 1[^J \rightarrow \mathbb{R}^J$ . Since the computation is the same for any cell ratio  $p_j$ , we give an explicit formula for only one of them (Eq.(4)):

$$\begin{aligned} \frac{\partial \ell_{\mathbf{y}|\zeta}(\mathbf{p})}{\partial p_j} &= \frac{\partial \log(\text{Det}(\Theta))}{\partial p_j} - \frac{1}{2} \left[ \frac{\partial (\mathbf{y} - \mu\mathbf{p})^\top}{\partial p_j} \Theta (\mathbf{y} - \mu\mathbf{p}) + (\mathbf{y} - \mu\mathbf{p})^\top \frac{\partial \Theta}{\partial p_j} (\mathbf{y} - \mu\mathbf{p}) + (\mathbf{y} - \mu\mathbf{p})^\top \Theta \frac{\partial (\mathbf{y} - \mu\mathbf{p})}{\partial p_j} \right] \\ &= -\text{Tr} \left( \Theta \frac{\partial \Sigma}{\partial p_j} \right) - \frac{1}{2} \left[ -\mu_{.j}^\top \Theta (\mathbf{y} - \mu\mathbf{p}) - (\mathbf{y} - \mu\mathbf{p})^\top \Theta \frac{\partial \Sigma}{\partial p_j} \Theta (\mathbf{y} - \mu\mathbf{p}) - (\mathbf{y} - \mu\mathbf{p})^\top \Theta \mu_{.j} \right] \\ &= -2p_j \text{Tr}(\Theta \Sigma_j) + (\mathbf{y} - \mu\mathbf{p})^\top \Theta \mu_{.j} + p_j (\mathbf{y} - \mu\mathbf{p})^\top \Theta \Sigma_j \Theta (\mathbf{y} - \mu\mathbf{p}) \end{aligned} \quad (4)$$

Since the solution to  $\nabla(\ell_{\mathbf{y}|\zeta}(\mathbf{p})) = 0$  is not closed, we had to approximate the MLE using iterated numerical optimisation methods. Some of them, such as the Levenberg–Marquardt algorithm, require a second-order approximation of the function, which needs the computation of the Hessian matrix. Deriving once more Eq.(5) yields the Hessian matrix,  $\mathbf{H} \in \mathcal{M}_{J \times J}$  is given by:

$$\begin{aligned}
\mathbf{H}_{i,i} &= \frac{\partial^2 \ell}{\partial^2 p_i} = -2 \text{Tr}(\Theta \Sigma_i) + 4p_i^2 \text{Tr}\left((\Theta \Sigma_i)^2\right) - 2p_i(\mathbf{y} - \mu \mathbf{p})^\top \Theta \Sigma_i \Theta \mu_{.i} - \mu_{.i}^\top \Theta \mu_{.i} - \\
&\quad 2p_i(\mathbf{y} - \mu \mathbf{p})^\top \Theta \Sigma_i \Theta \mu_{.i} - (\mathbf{y} - \mu \mathbf{p})^\top \Theta (4p_i^2 \Sigma_i \Theta \Sigma_i - \Sigma_i) \Theta (\mathbf{y} - \mu \mathbf{p}), \quad i \in \tilde{J} \\
\mathbf{H}_{i,j} &= \frac{\partial^2 \ell}{\partial p_i \partial p_j} = 4p_j p_i \text{Tr}(\Theta \Sigma_j \Theta \Sigma_i) - 2p_i(\mathbf{y} - \mu \mathbf{p})^\top \Theta \Sigma_i \Theta \mu_{.j} - \mu_{.i}^\top \Theta \mu_{.j} - \\
&\quad 2p_j(\mathbf{y} - \mu \mathbf{p})^\top \Theta \Sigma_j \Theta \mu_{.i} - 4p_i p_j (\mathbf{y} - \mu \mathbf{p})^\top \Theta \Sigma_i \Theta \Sigma_j \Theta (\mathbf{y} - \mu \mathbf{p}), \quad (i, j) \in \tilde{J}^2, i \neq j
\end{aligned} \tag{5}$$

in which the coloured sections pair one by one with the corresponding coloured sections of the gradient, given in Eq.(4). Matrix calculus can largely ease the derivation of complex algebraic expressions, thus we remind in Appendix Matrix calculus relevant matrix properties and derivations <sup>1</sup>.

However, the explicit formulas for the gradient and the Hessian matrix of the log-likelihood function, given in Eq.(4) and Eq.(5) respectively, do not take into account the simplex constraint assigned to the ratios. While some optimisation methods use heuristic methods to solve this problem, we consider alternatively a reparametrised version of the problem, detailed comprehensively in Appendix Reparametrised log-likelihood.

## 1.4 Iterated optimisation

The MLE is traditionally retrieved from the roots of the gradient of the log-likelihood. However, in our generative framework, cancelling the gradient of Equation (3) reveals a non-closed form. Instead, iterated numerical optimisation algorithms can be used to proxy the roots, most of them considering first or second-order approximations of the function to optimise.

The *Levenberg-Marquardt algorithm* bridges the gap between between the steepest descent method (first-order) and the Newton-Raphson method (second-order) by inflating the diagonal terms of the Hessian matrix. Away from the endpoint, a second-order descent is favoured for its faster convergence pace, while the steepest approach is privileged close to the extremum, as it allows careful refinement of the step size. We use function **marqLevAlg**, since it notably introduces a stringent convergence criteria, the relative distance to the maximum (RDM), which sets apart extrema from spurious saddle points [3].

We provide additional theoretical results, such as analytical formulas for the Gradient and the Hessian in their constrained and unconstrained versions as well as simulation outputs in the vignette of the DeCovarT Github webpage.

## References

- [1] T. Erkkilä, S. Lehmusvaara, P. Ruusuvuori, T. Visakorpi, I. Shmulevich, and H. Lähdesmäki, “Probabilistic analysis of gene expression measurements from heterogeneous tissues,” *Bioinformatics*, vol. 26, no. 20, pp. 2571–2577, Oct. 2010, doi: 10.1093/bioinformatics/btq406.
- [2] R. Mazumder and T. Hastie, “The Graphical Lasso: New Insights and Alternatives,” *Electronic Journal of Statistics*, vol. 6, Nov. 2011, doi: 10.1214/12-EJS740.
- [3] M. Prague, D. Commenges, J. Guedj, J. Drylewicz, and R. Thiébaud, “NIMROD: A program for inference via a normal approximation of the posterior in models with random effects based on ordinary differential equations,” *Computer Methods and Programs in Biomedicine*, vol. 111, no. 2, pp. 447–458, Aug. 2013, doi: 10.1016/j.cmpb.2013.04.014.

---

<sup>1</sup>The numerical consistency of these derivatives was asserted with the **numDeriv numDeriv** package, using the more stable Richardson’s extrapolation

## A Theoretical details

### First and second-order derivation of the constrained DeCovarT log-likelihood function

To reparametrise the log-likelihood function (Eq.(3)) in order to explicitly handling the unit simplex constraint (Eq.(1)), we consider the following mapping function:  $\boldsymbol{\psi} : \boldsymbol{\theta} \rightarrow \mathbf{p} \mid \boldsymbol{\theta} \in \mathbb{R}^{J-1}, \mathbf{p} \in ]0, 1[^J$  (Eq.(6)):

$$\mathbf{p} = \boldsymbol{\psi}(\boldsymbol{\theta}) = \begin{cases} p_j = \frac{e^{\theta_j}}{\sum_{k < J} e^{\theta_k} + 1}, & j < J \\ p_J = \frac{1}{\sum_{k < J} e^{\theta_k} + 1} \end{cases} \quad \boldsymbol{\theta} = \boldsymbol{\psi}^{-1}(\mathbf{p}) = \left( \ln \left( \frac{p_j}{p_J} \right) \right)_{j \in \{1, \dots, J-1\}} \quad (6)$$

that is a  $C^2$ -diffeomorphism, since  $\boldsymbol{\psi}$  is a bijection between  $\mathbf{p}$  and  $\boldsymbol{\theta}$  twice differentiable.

Its Jacobian,  $\mathbf{J}_{\boldsymbol{\psi}} \in \mathcal{M}_{J \times (J-1)}$  is given by Eq.(7):

$$\mathbf{J}_{i,j} = \frac{\partial p_i}{\partial \theta_j} = \begin{cases} \frac{e^{\theta_i} B_i}{A^2}, & i = j, i < J \\ \frac{-e^{\theta_j} e^{\theta_i}}{A^2}, & i \neq j, i < J \\ \frac{-e^{\theta_j}}{A^2}, & i = J \end{cases} \quad (7)$$

with  $i$  indexing vector-valued  $\mathbf{p}$  and  $j$  indexing the first-order partial derivatives of the mapping function,  $A = \sum_{j' < J} e^{\theta_{j'}} + 1$  the sum over exponential (denominator of the mapping function) and  $B = A - e^{\theta_i}$  the sum over ratios minus the exponential indexed with the currently considered index  $i$ .

The Hessian (which fortunately is symmetric for each component  $j$ , as expected according to the Schwarz's theorem) of the vectorial mapping function  $\boldsymbol{\psi}(\boldsymbol{\theta})$  is a third-order tensor of rank  $(J-1)(J-1)J$ , given by Eq.(8):

$$\frac{\partial^2 p_i}{\partial k \partial j} = \begin{cases} \frac{e^{\theta_i} e^{\theta_l} (-B_i + e^{\theta_i})}{A^3}, & (i < J) \wedge ((i \neq j) \oplus (i \neq k)) \quad (a) \\ \frac{2e^{\theta_i} e^{\theta_j} e^{\theta_k}}{A^3}, & (i < J) \wedge (i \neq j \neq k) \quad (b) \\ \frac{e^{\theta_i} e^{\theta_j} (-A + 2e^{\theta_j})}{A^3}, & (i < J) \wedge (j = k \neq i) \quad (c) \\ \frac{B_i e^{\theta_i} (B_i - e^{\theta_i})}{A^3}, & (i < J) \wedge (j = k = i) \quad (d) \\ \frac{e^{\theta_j} (-A + 2e^{\theta_j})}{A^3}, & (i = J) \wedge (j = k) \quad (e) \\ \frac{2e^{\theta_j} e^{\theta_k}}{A^3}, & (i = J) \wedge (j \neq k) \quad (f) \end{cases} \quad (8)$$

with  $i$  indexing  $\mathbf{p}$ ,  $j$  and  $k$  respectively indexing the first-order and second-order partial derivatives of the mapping function with respect to  $\boldsymbol{\theta}$ . In line (a),  $\oplus$  refers to the Boolean XOR operator,  $\wedge$  to the AND operator and  $l = \{j, k\} \setminus i$ .

To derive the log-likelihood function in Eq.(4), we reparametrise  $\mathbf{p}$  to  $\boldsymbol{\theta}$ , using a standard \textit{chain rule} formula). Considering the original log-likelihood function, Eq.(3), and the mapping function, Eq.(6), the differential at the first order and at the second order is given by Eq.(9) and Eq.(10), respectively defined in  $\mathbb{R}^{J-1}$  and  $\mathcal{M}_{(J-1) \times (J-1)}$ :

$$\left[ \frac{\partial \ell_{\mathbf{y}|\boldsymbol{\zeta}}}{\partial \theta_j} \right]_{j < J} = \sum_{i=1}^J \frac{\partial \ell_{\mathbf{y}|\boldsymbol{\zeta}}}{\partial p_i} \frac{\partial p_i}{\partial \theta_j} \quad (9)$$

$$\left[ \frac{\partial^2 \ell_{\mathbf{y}|\boldsymbol{\zeta}}}{\partial \theta_k \partial \theta_j} \right]_{j < J, k < J} = \sum_{i=1}^J \sum_{l=1}^J \left( \frac{\partial p_i}{\partial \theta_j} \frac{\partial^2 \ell_{\mathbf{y}|\boldsymbol{\zeta}}}{\partial p_i \partial p_l} \frac{\partial p_l}{\partial \theta_k} \right) + \sum_{i=1}^J \left( \frac{\partial \ell_{\mathbf{y}|\boldsymbol{\zeta}}}{\partial p_i} \frac{\partial^2 p_i}{\partial \theta_k \partial \theta_j} \right) \quad (d) \quad (10)$$