

DeCovarT, a generative model for the deconvolution of heterogeneous transcriptomic samples

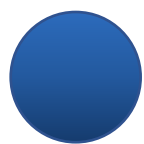


17/06/2023

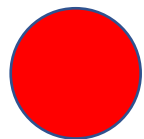
Bastien CHASSAGNOL
Gregory NUEL, Etienne BECHT, Yufei LUO

SERVIER 
moved by you

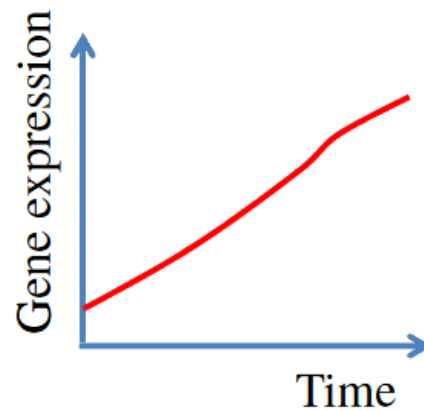
Confusing biological noise



resting cell population 1



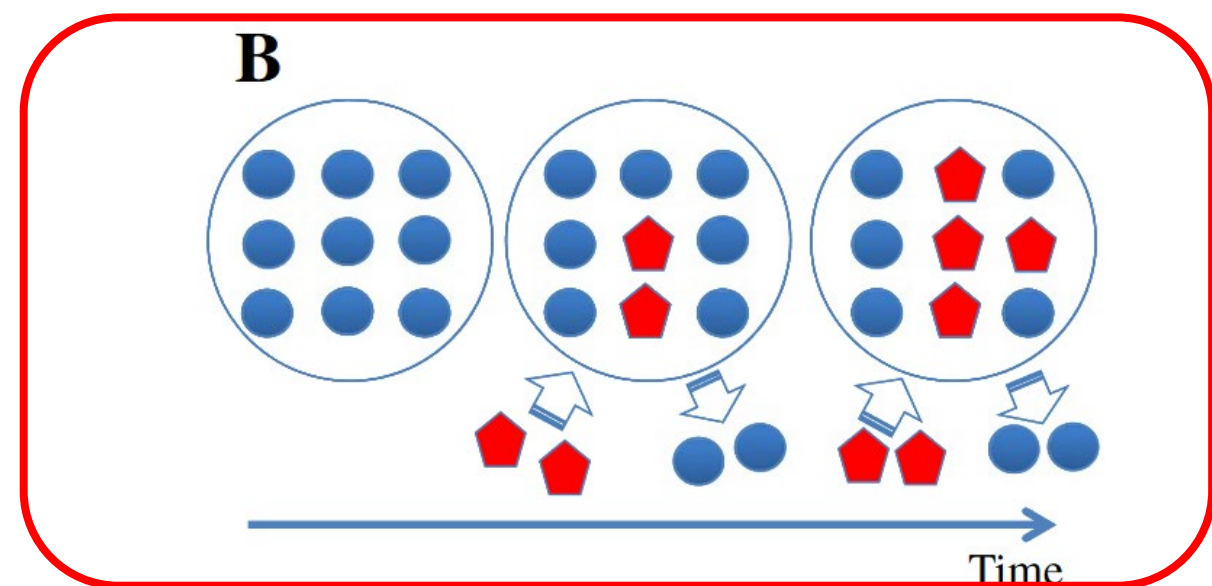
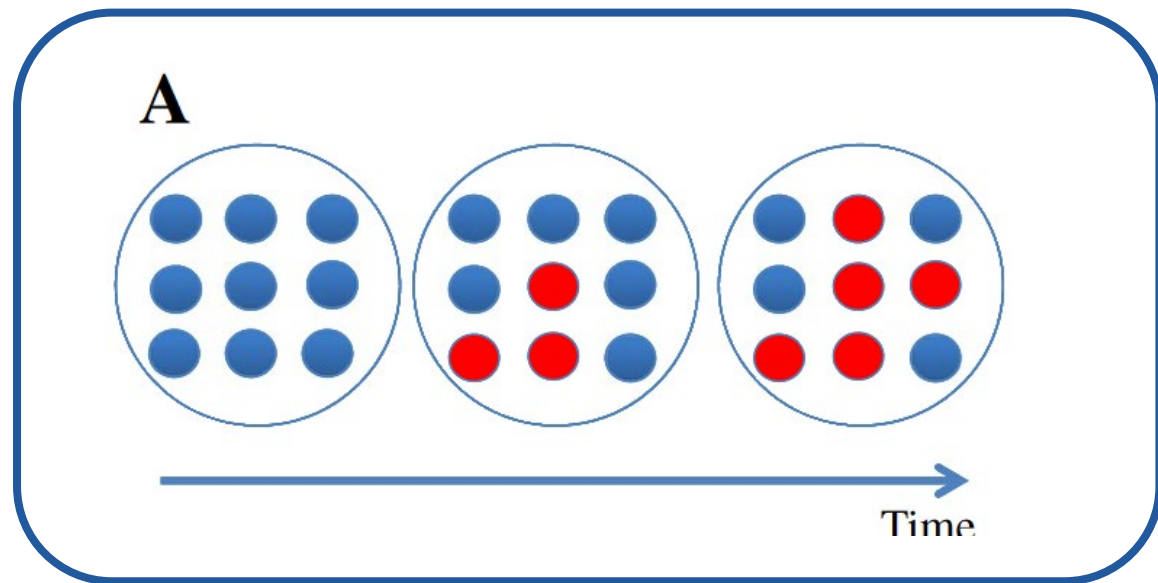
activated cell population 1



Shoemaker et al. 2012c



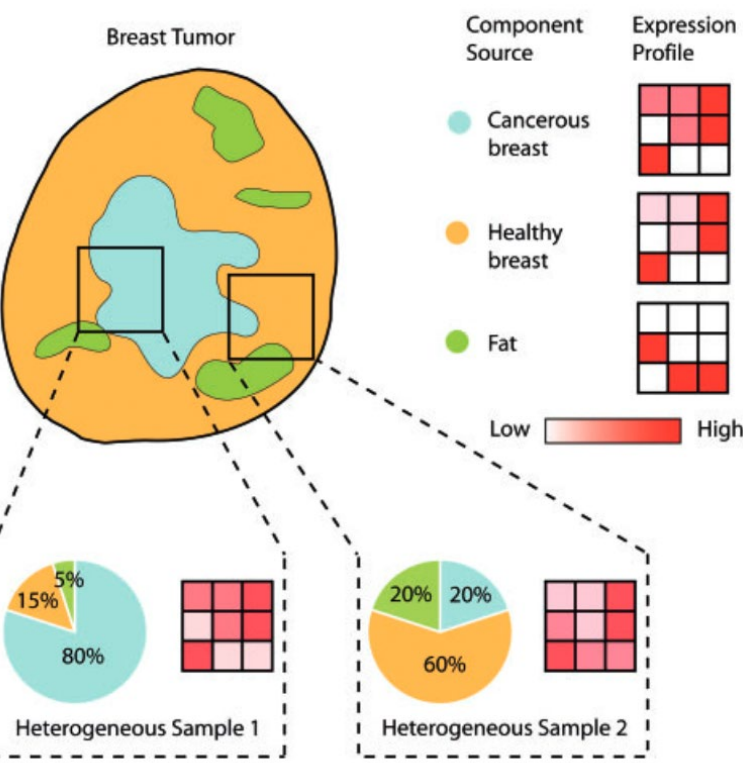
activated cell population 2



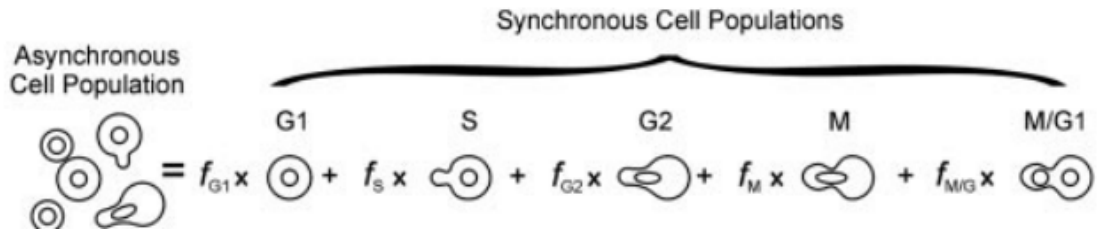
Scenario A: increase of the gene expression is generated by an **activation** of cell population 1

Scenario B: the gene expression increases due to the **arrival** of a **new** cell population 2

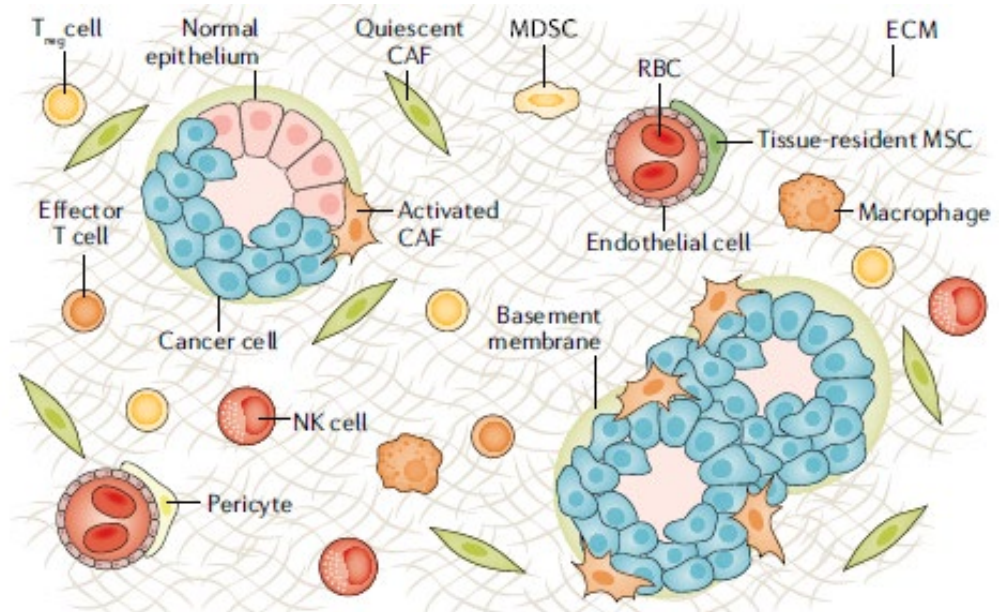
Heterogeneity of tissues



Mixture of tissues
Quon and Morris, 2009



Mixture of cell phases
Lu et al, 2003



Mixture of cell populations
Finotello and Trajanoski 2018

Deconvolution inputs

$$\begin{pmatrix} x_{1,1} & \dots & x_{1,J} \\ \vdots & \ddots & \vdots \\ x_{G,1} & \dots & x_{G,J} \end{pmatrix} \times \begin{pmatrix} p_{1,1} & \dots & p_{1,N} \\ \vdots & \ddots & \vdots \\ p_{J,1} & \dots & p_{J,N} \end{pmatrix} = \begin{pmatrix} y_{1,1} & \dots & y_{1,N} \\ \vdots & \ddots & \vdots \\ y_{G,1} & \dots & y_{G,N} \end{pmatrix}$$

X stores purified cellular expression profiles

p is the individual vector of cell ratios

Y stores the resulting bulk expression values

$$\begin{pmatrix} x_{G_1,1} & \dots & 0 \\ 0 & x_{G_2,2} & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & x_{G_k,k} \end{pmatrix}$$

Marker-based

Linear box constraints

$$X p_i = y_i$$

$$\begin{cases} \sum_{j=1}^J p_{ji} = 1 \\ \forall j \in \{1, \dots, J\}, \quad p_{ji} \geq 0 \end{cases}$$

Deconvolution ecosystem

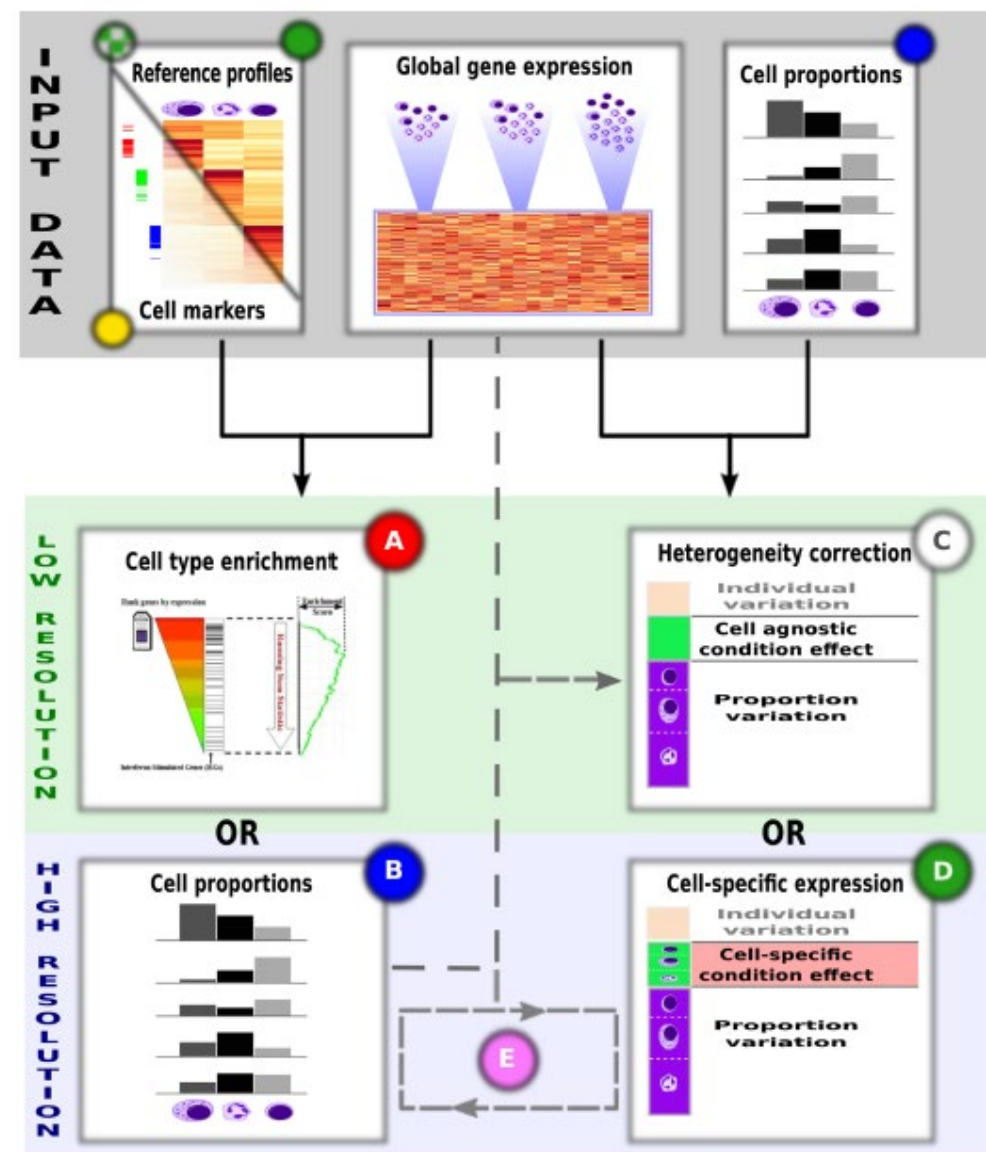
Partial
deconvolution

Estimate the ratios p for all individuals with the purified cell signature X and bulk mixture y .

- Try to infer cell specific expression profiles X based on p and y .

Complete
deconvolution

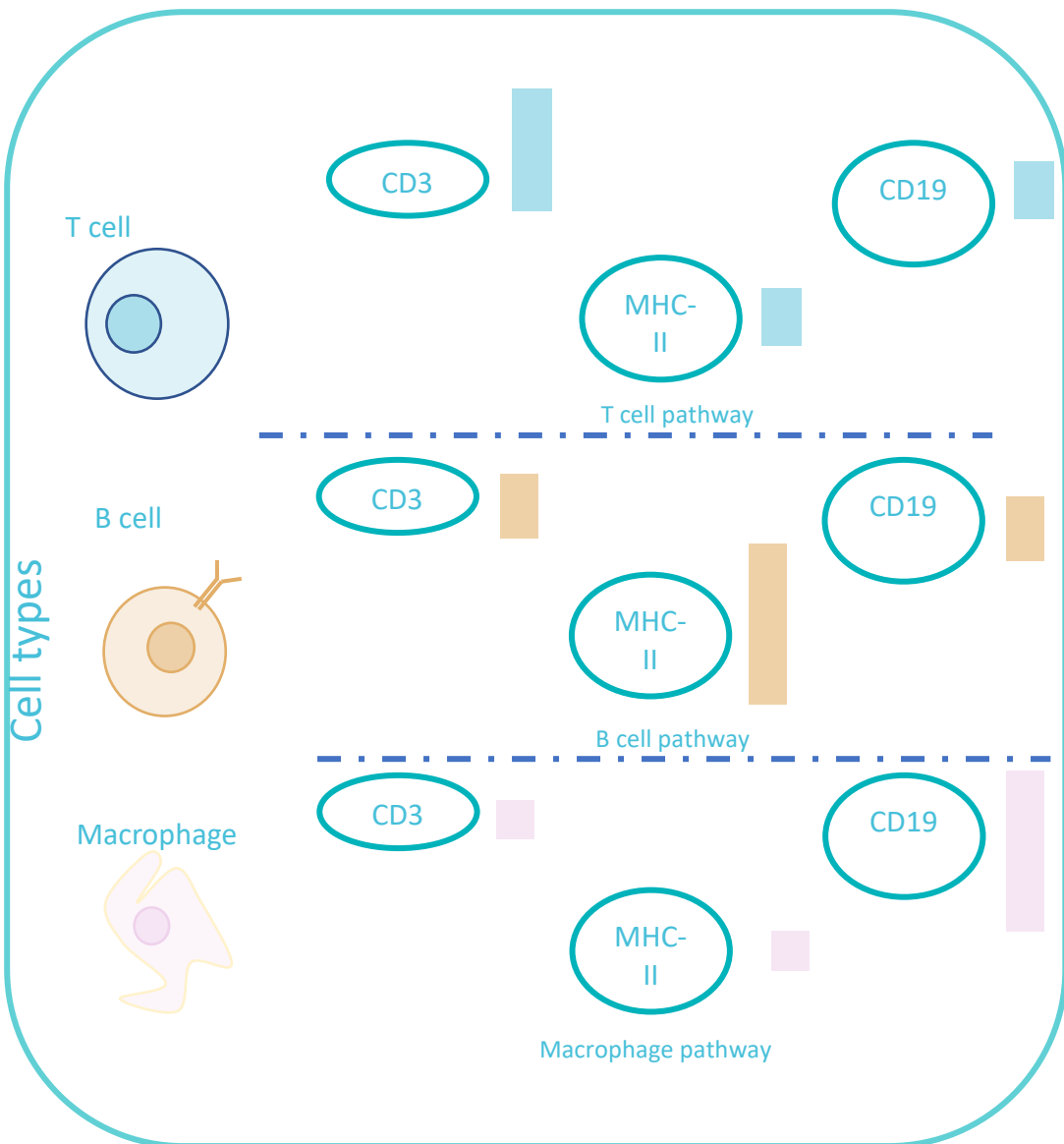
- Try to infer alternatively both p and X (unsupervised, reference-free methods).
Undetermined problem without prior.



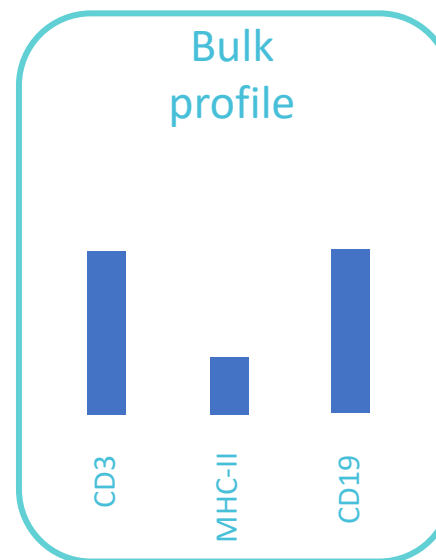
Shen-Orr et al, 2013

Standard deconvolution framework

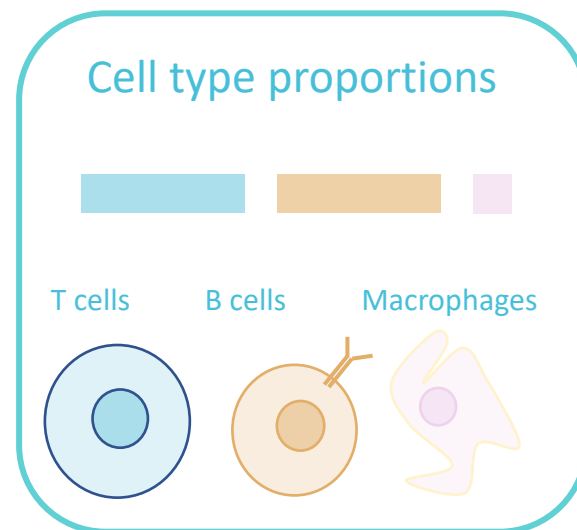
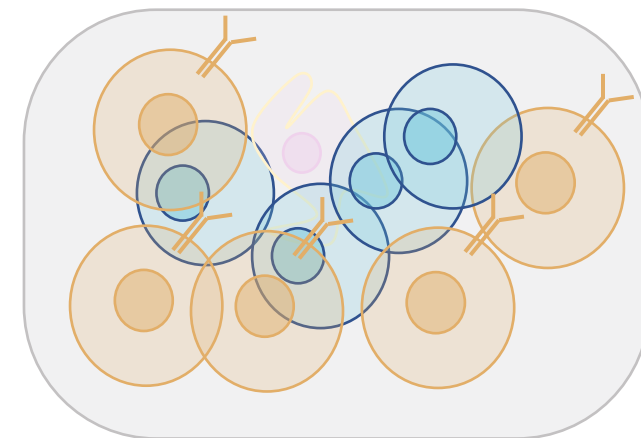
Cell types



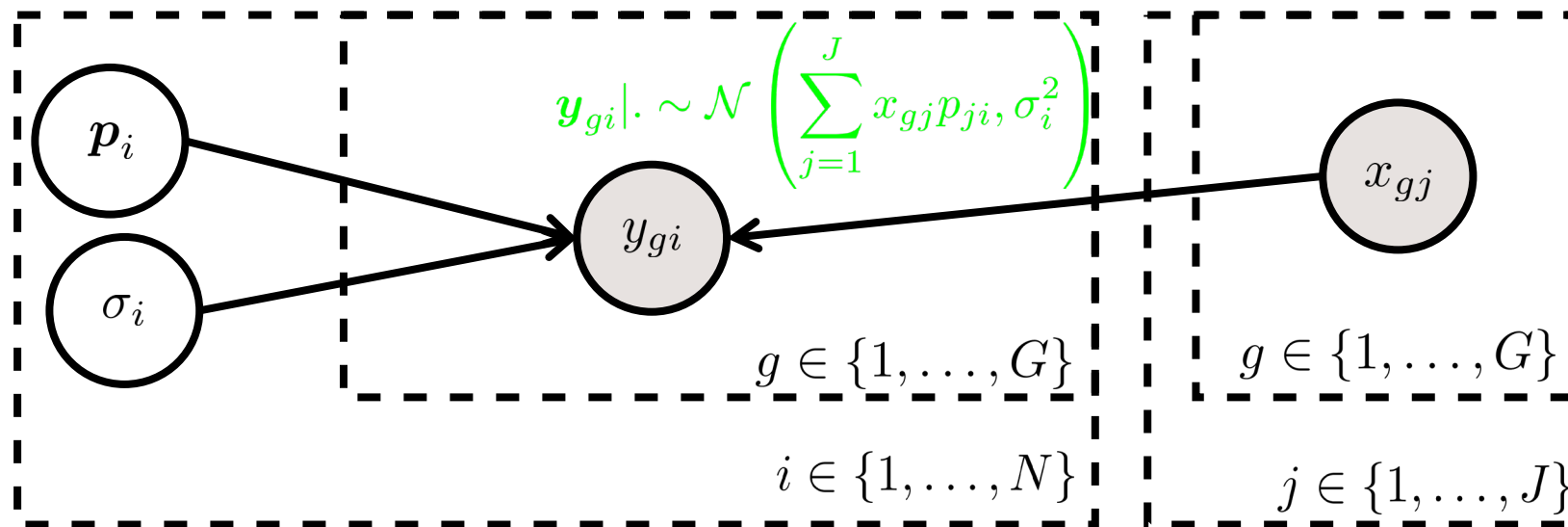
Purified cellular profiles



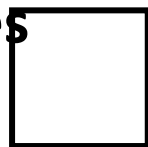
Standard deconvolution



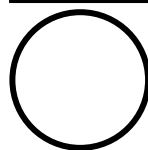
Graphical model of linear regression



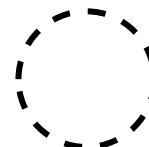
Nodes



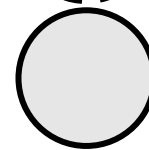
Constant



Stochastic
variable



Deterministic
variable



Observations

Parameters

Estimated parameters

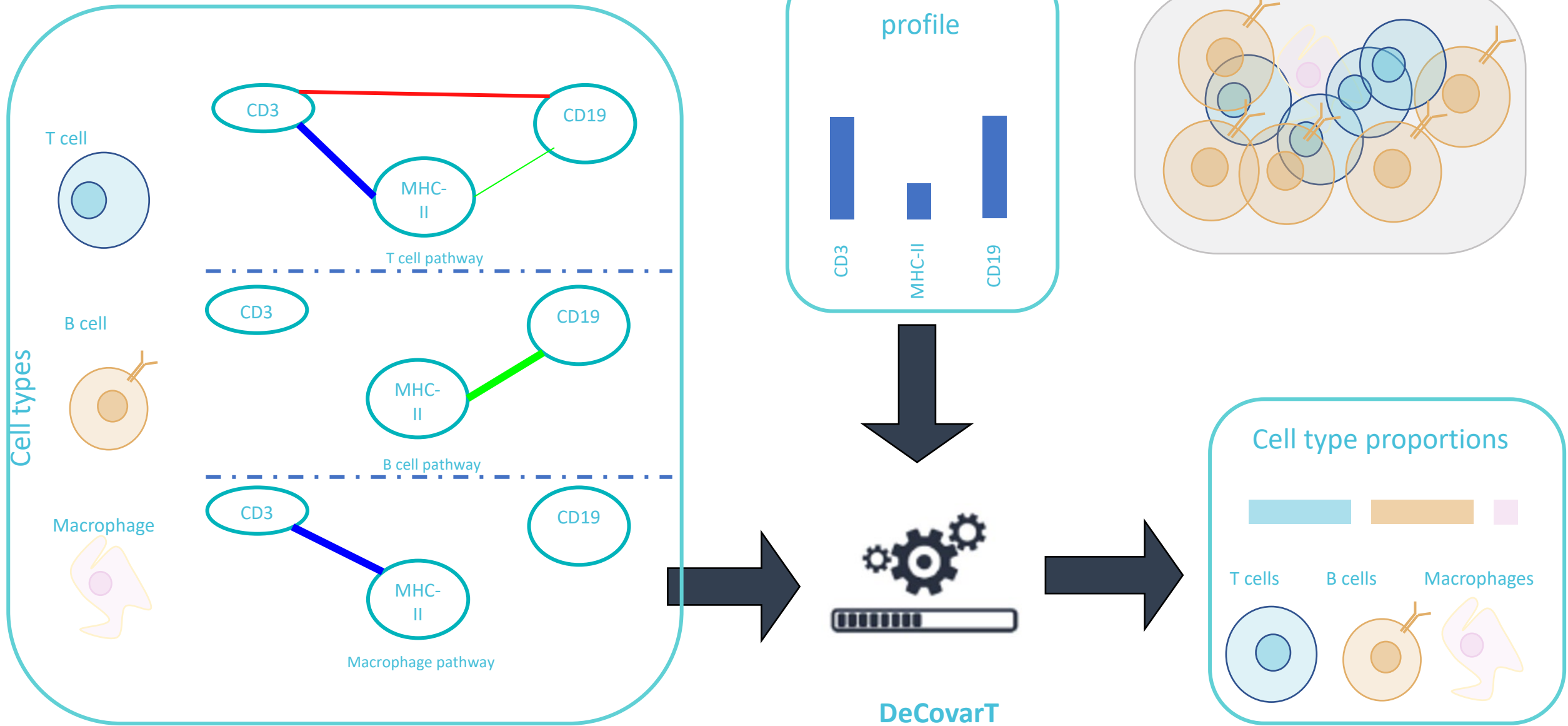
$$\theta = (\mathbf{p}, \sigma)$$

Distribution probabilities

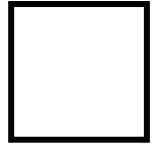
Likelihood Laws

$$f(\mathcal{D} | \theta)$$

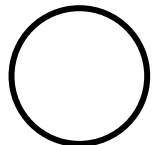
DeCovarT framework



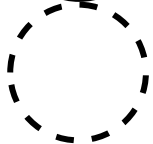
Nodes



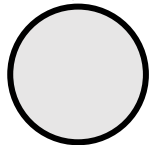
Constant



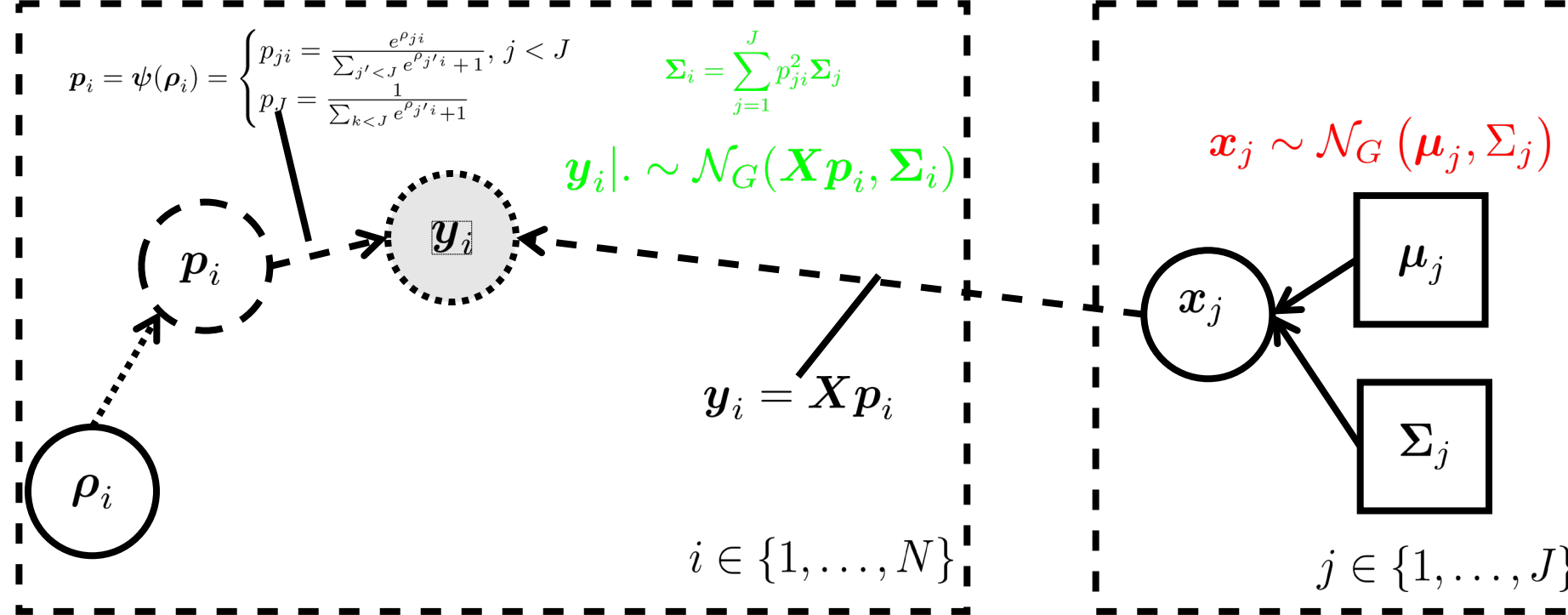
Stochastic variable



Deterministic variable



Observations



Graphical model of DeCoVarT

Parameters

Prior parameters

$$\zeta = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Estimated parameters

$$\theta = (\boldsymbol{p}, \boldsymbol{X})$$

Distribution probabilities

Prior laws

$$f(\theta|\xi)$$

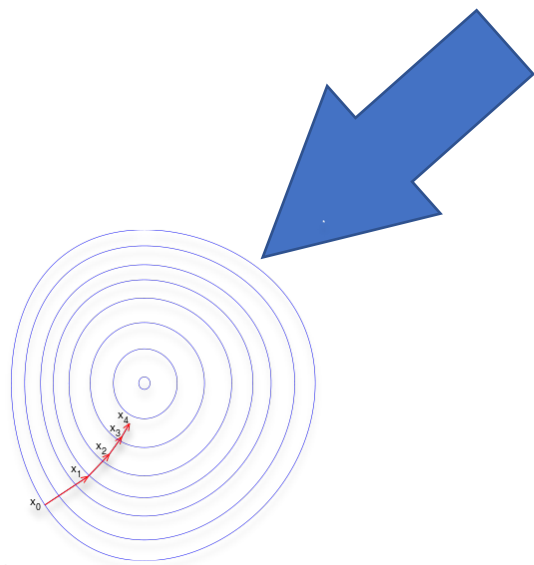
Likelihood Laws

$$f(\mathcal{D}|\theta)$$

Optimisation algorithms in R

$$\ell_{\mathbf{y}|\mathbf{X},\Sigma}(\mathbf{p}) = C + \log \left(\det \left(\sum_{j=1}^J p_j^2 \Sigma_j \right)^{-1} \right) - \frac{1}{2} (\mathbf{y} - \mathbf{Xp})^\top \left(\sum_{j=1}^J p_j^2 \Sigma_j \right)^{-1} (\mathbf{y} - \mathbf{Xp})$$

Log-likelihood associated to the model = log-likelihood of a multivariate Gaussian distribution



Steepest descent method

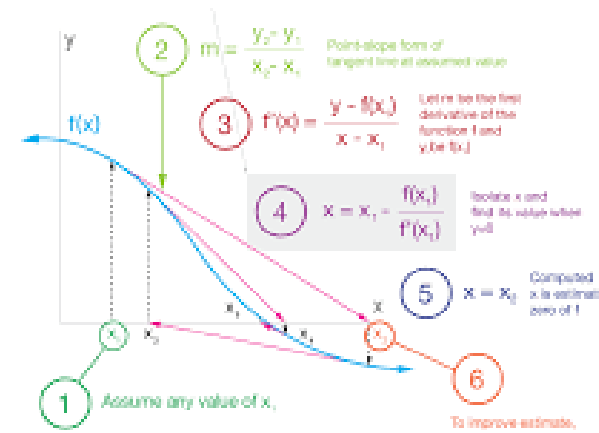
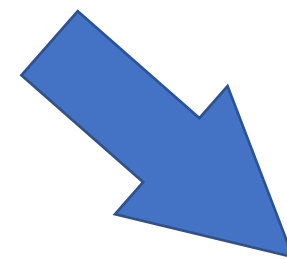
$$\theta^{(q+1)} = \theta^{(q)} + \gamma \nabla \ell(\theta^{(q)})$$



Updated **Marquardt-Levenberg**,
in package [*marqLevAlg*](#)

BP Hejblum, 2021, R Journal

$$\theta^{(q+1)} = \theta^{(q)} + \gamma_q \frac{\nabla \ell(\theta^{(q)})}{\tilde{\mathbf{H}}_\ell(\theta^{(q)})}$$



Newton Raphson method

$$\theta^{(q+1)} = \theta^{(q)} + \frac{\nabla \ell(\theta^{(q)})}{\mathbf{H}_\ell(\theta^{(q)})}$$

