

A semi-parametric empirical Bayes approach for time-course integrative genomic analysis

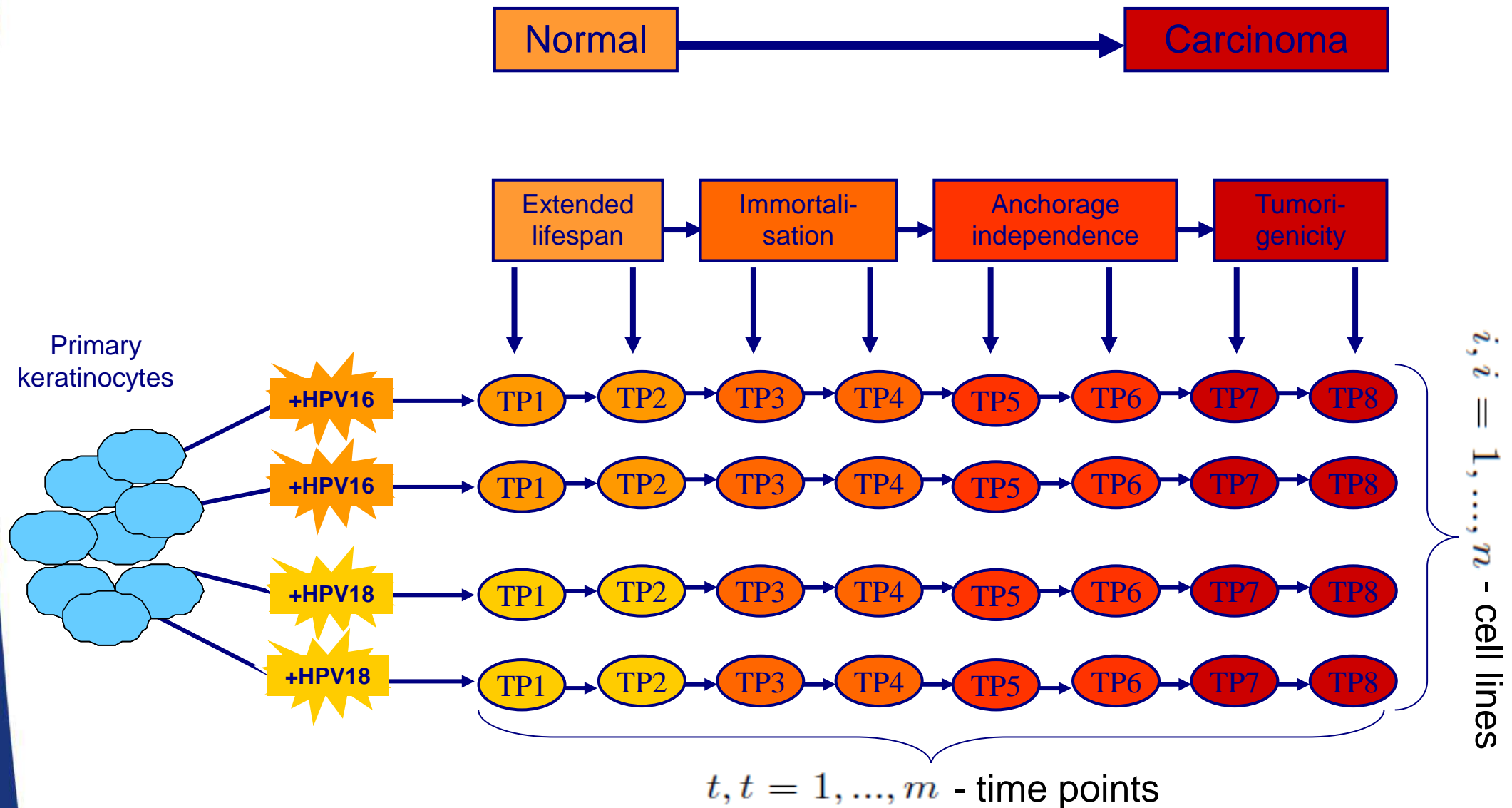
Viktorian Miok^{1,2}, Saskia Wilting², Mark van de Wiel^{1,3}, Annelieke Jaspers², Peter Snijders², Renske Steenbergen², Wessel van Wieringen^{1,3}

¹Department of Epidemiology & Biostatistics and ²Department of Pathology, VU University Medical Center, ³Department of Mathematics, VU University, Amsterdam, The Netherlands.

Cervical cancer study

- Second most common cancer in women worldwide.
- Caused by HPV virus, in 80% cases HPV16 and HPV18.
- Cell line model – in vitro model system of HPV-induced transformation.
- Integration – high-throughput multi level molecular data sets.
- Aim: identification of key genes.

Experiment



Model

$j, j = 1, \dots, p$ - genes

$\mathbf{Y}_{**t} = (\mathbf{Y}_{1*t}, \dots, \mathbf{Y}_{n*t})$ - mRNA gene expression

- **Bayesian GLMM:**

$$\mathbf{Y}_{ijt} \sim \mathcal{N}(\mu_{ijt}, \sigma_j^2)$$

Cell line effect Time effect

$$\mu_{ijt} = \overbrace{f_j(v_i; \boldsymbol{\alpha})}^{\text{Cell line effect}} + \overbrace{h_j(t; \boldsymbol{\beta})}^{\text{Time effect}} + \varepsilon_{ijt},$$

$\boldsymbol{\alpha}, \boldsymbol{\beta}$ - Gaussian distribution assumption

Fixed and random effects

Fixed effect:

$$f_{ij}(v_i; \boldsymbol{\alpha}) = \alpha_{ij} v_i,$$

Random effect:

$$h_j(t; \boldsymbol{\beta}) = \sum_{k=1}^K \beta_{jk} |t - \kappa_k|^3$$

Matrix notation: $\mathbf{Y}_{it} = \tilde{\mathbf{V}} \tilde{\boldsymbol{\alpha}} + \tilde{\mathbf{Z}} \tilde{\boldsymbol{\beta}} + \varepsilon_{it}$

$$\tilde{\mathbf{Z}} = \mathbf{Z}_K \boldsymbol{\Omega}_{K \times K}^{-1/2}$$

$$\tilde{\boldsymbol{\beta}} = \boldsymbol{\Omega}_{K \times K}^{1/2} \boldsymbol{\beta}$$

Spline basis: $\mathbf{Z}_K = \left\{ |t - \kappa_1|^3, \dots, |t - \kappa_K|^3 \right\}$ **Penalty matrix:** $\boldsymbol{\Omega}_{K \times K}$

Model parameters estimation

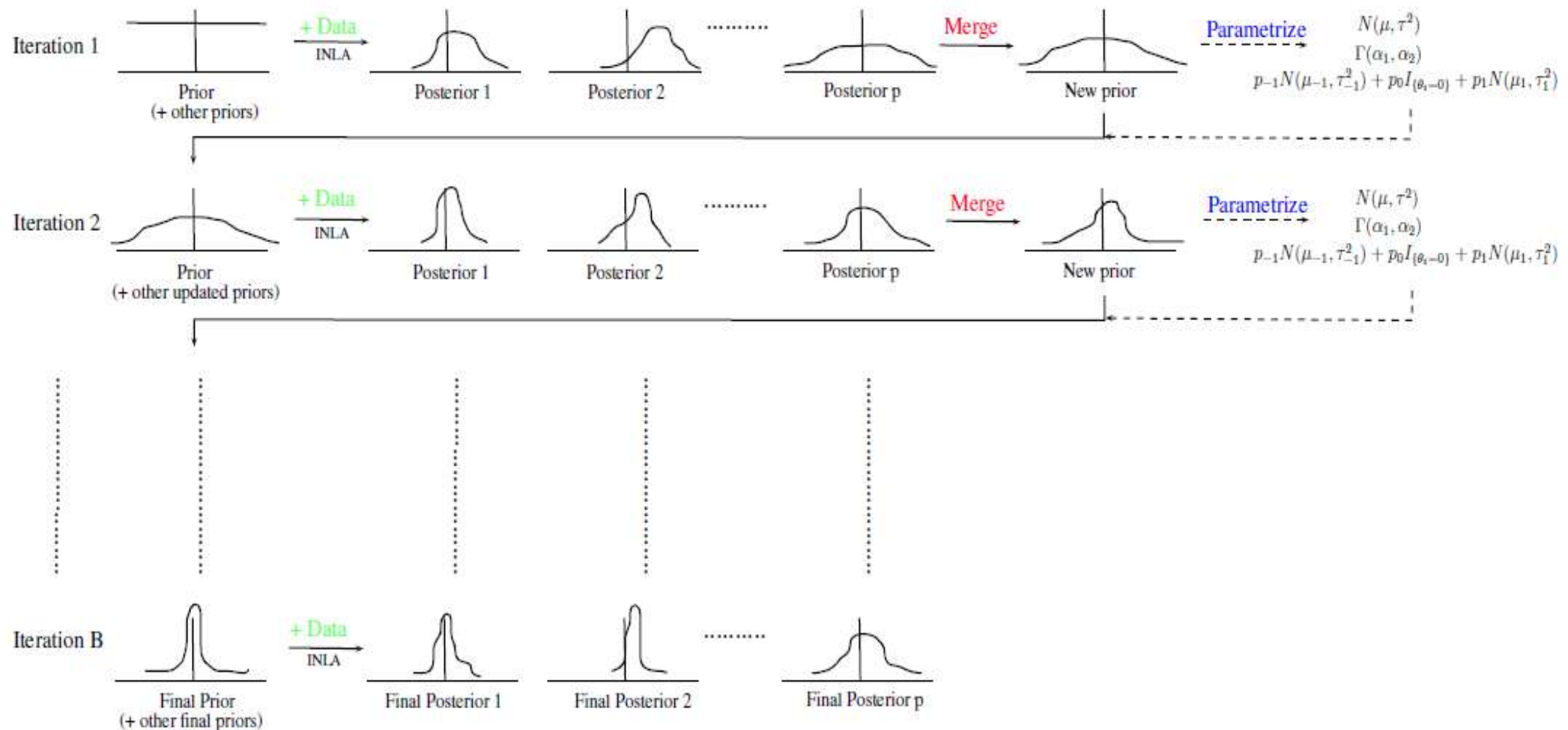
θ - model parameters

ϕ - hyper-parameters

INLA (Rue et al., 2009) procedure consist in:

- Approximate full posterior of $\pi(\phi|\mathbf{y})$ and $\pi(\theta_l|\phi, \mathbf{y})$ using Laplace approximation.
- Approximate marginal posterior densities of θ and ϕ integrating over hyper-parameters of posteriors $\pi(\phi|\mathbf{y})$ and $\pi(\theta_l|\phi, \mathbf{y})$.

Prior parameter estimation

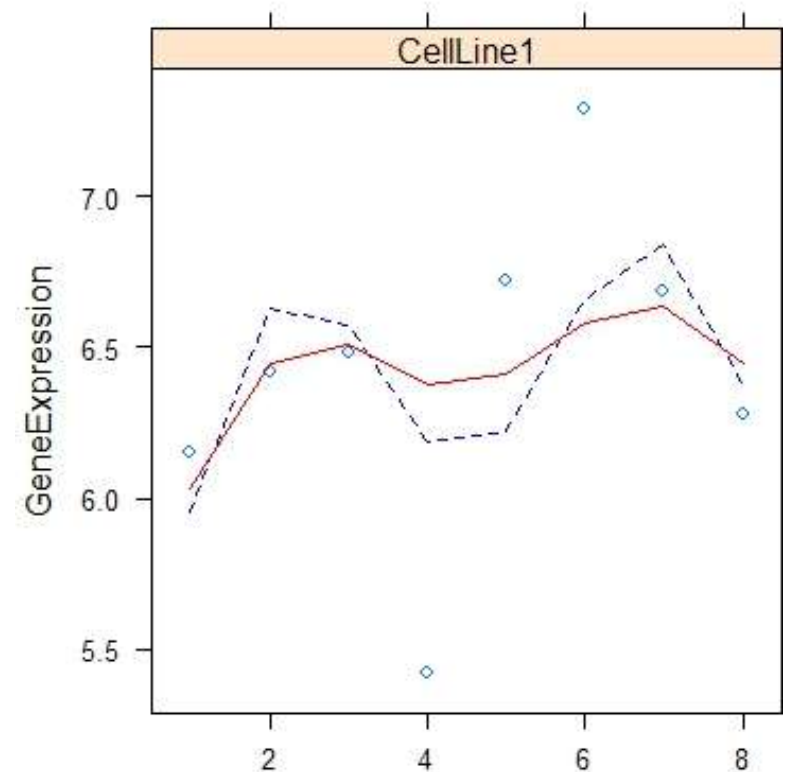


van de Wiel et al. (2013), Biostatistics.

Shrinkage

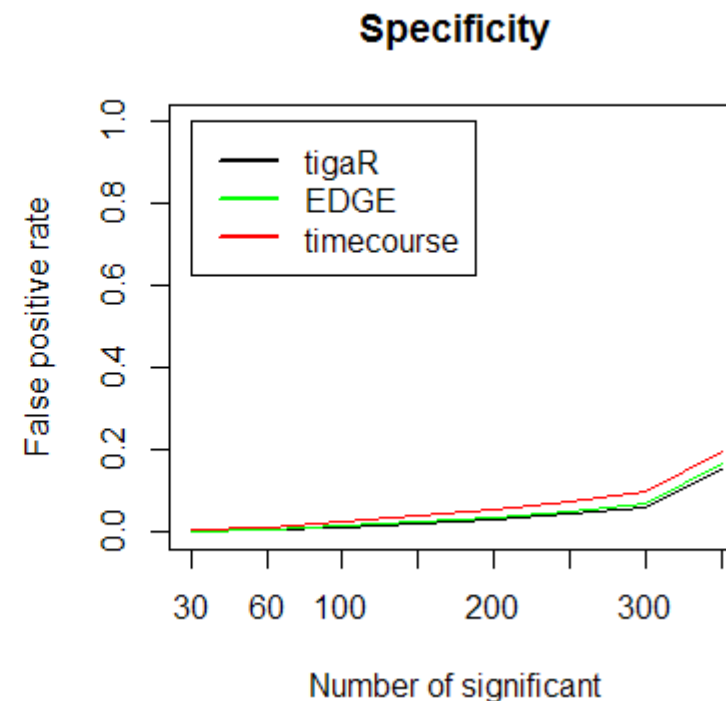
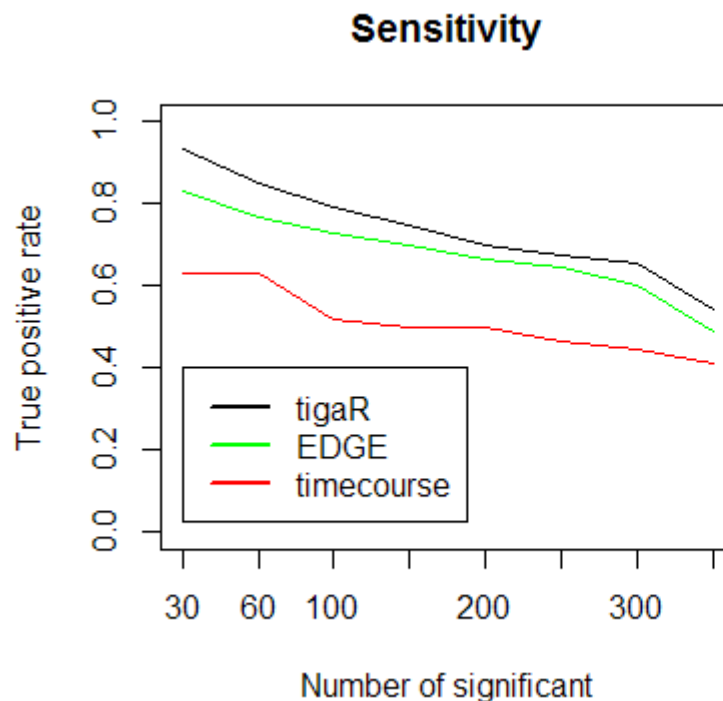
$$\sum_{i=1}^n \sum_{t=1}^m [y_{ijt} - f_{ij}(v_i; \alpha) - h_j(t; \beta)]^2 + \lambda \beta^T \mathbf{D} \beta, \text{ where } \lambda = \frac{\sigma_{\epsilon}^2}{\sigma_{\beta}^2}$$

- borrowing information across the genes
- better control of false positives
- leads to more stable estimates
- improvement of reproducibility

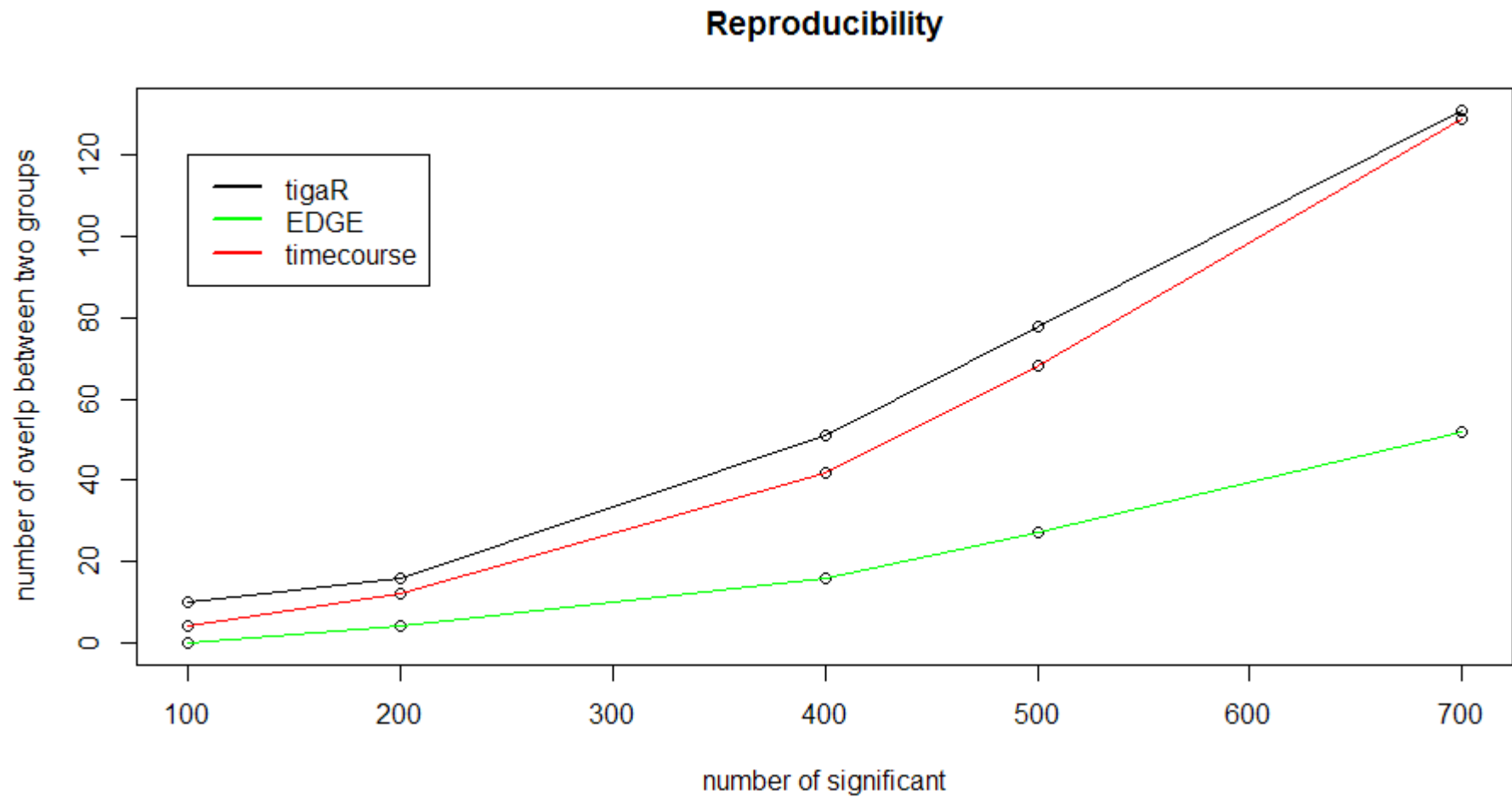


Sensitivity and specificity

- Comparison of methods
 - **tigaR** – temporal integrative genomic analysis in R
 - **EDGE** – Storey et al., PNAS., 2005.
 - **timecourse** – Tai and Speed, Annals of Statistics, 2006.



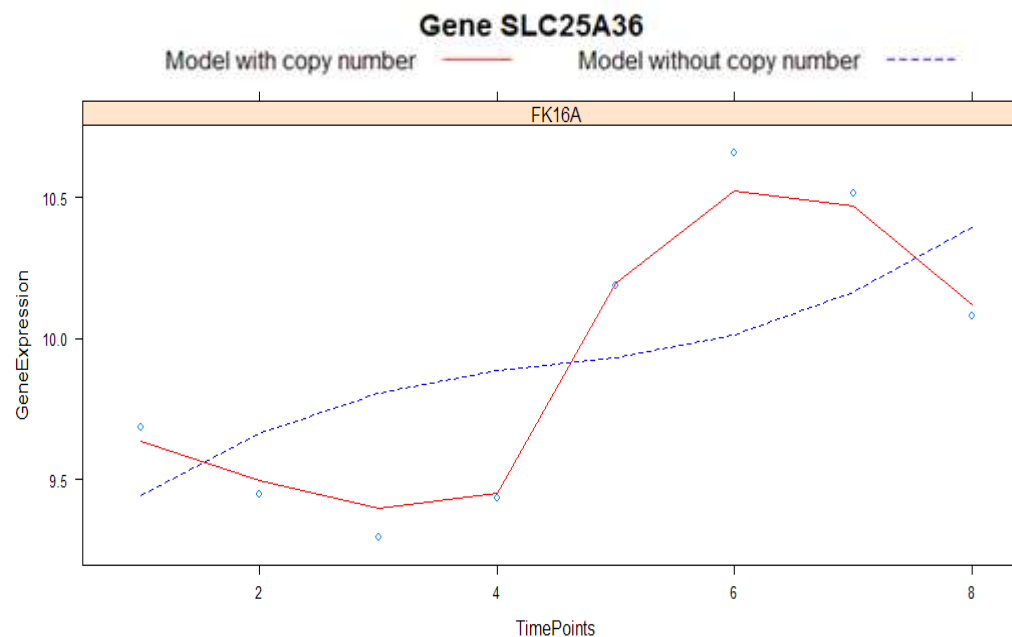
Reproducibility



DNA copy number (CN)

$\mathbf{X}_{**t} = (\mathbf{X}_{1*t}, \dots, \mathbf{X}_{n*t})$ - CN observations

Fixed effect: $f_{ij}(v_i, x_{ijt}; \alpha, \gamma) = \underbrace{\alpha_{ij}v_i}_{\text{Cell line effect}} + \underbrace{\gamma_j x_{ijt}}_{\text{CN effect}}$



Orthogonalization of splines

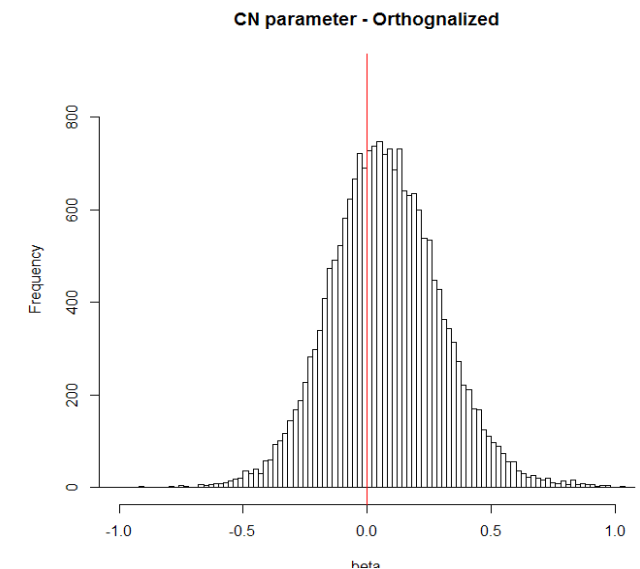
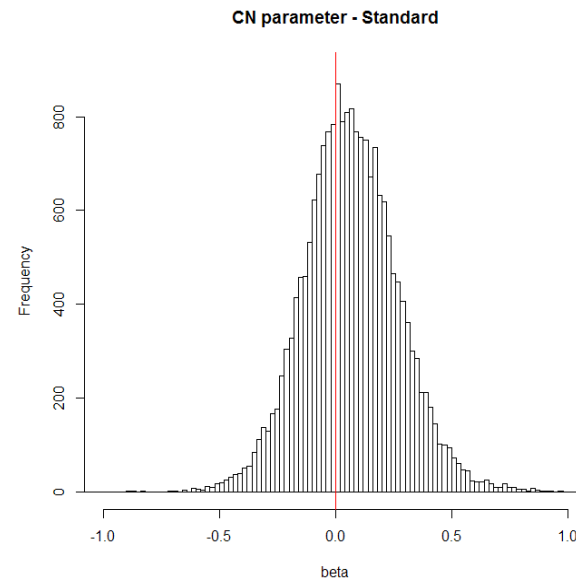
- Time effect consumes CN effect.
- Attribute effect to biological factor rather than to dynamic factor.
- Orthogonalize splines design matrix to CN.

$$\tilde{\mathbf{Z}}^{\perp} = \left(\mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right) \tilde{\mathbf{Z}}$$



Effect of orthogonalization

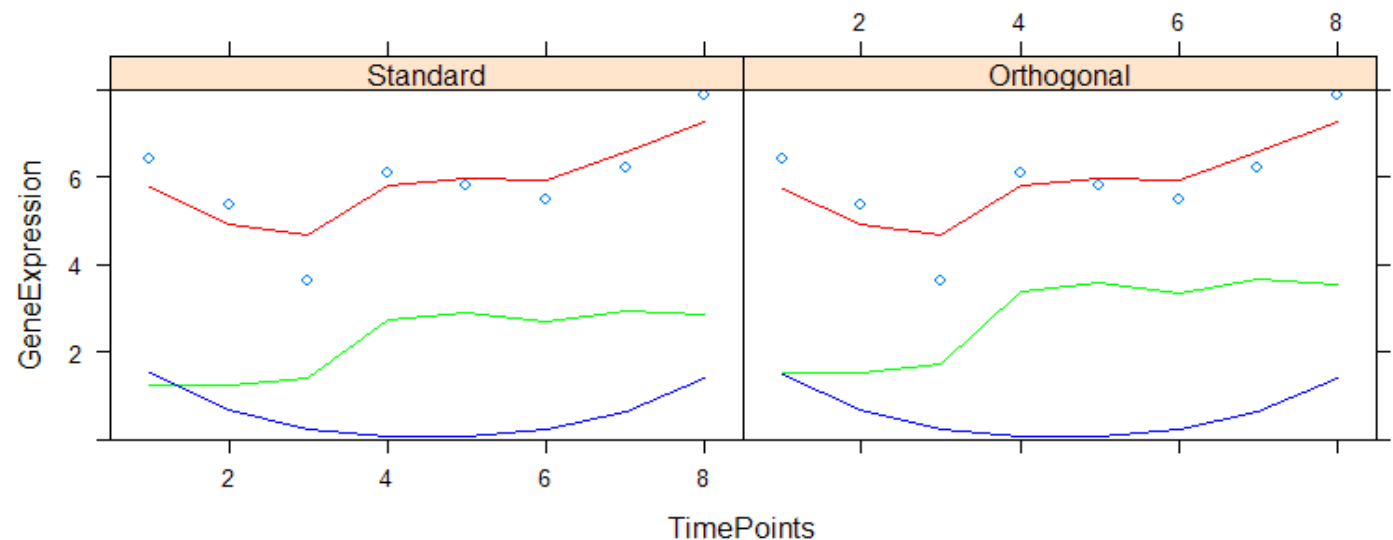
Parameter of CN:



CN effect — (green line)
Time effect — (blue line)

Full model — (red line)

Standard vs.
orthogonal:

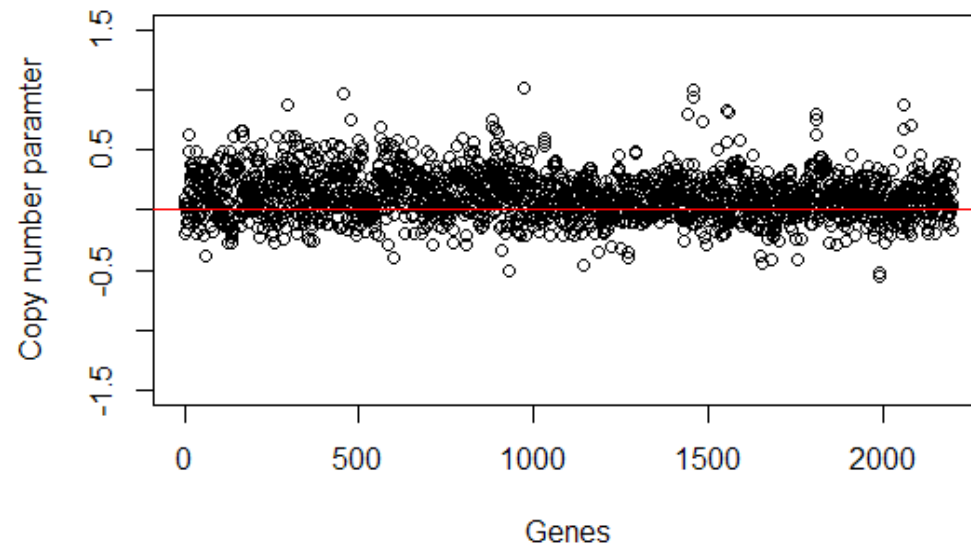


Spatial multivariate prior for CN

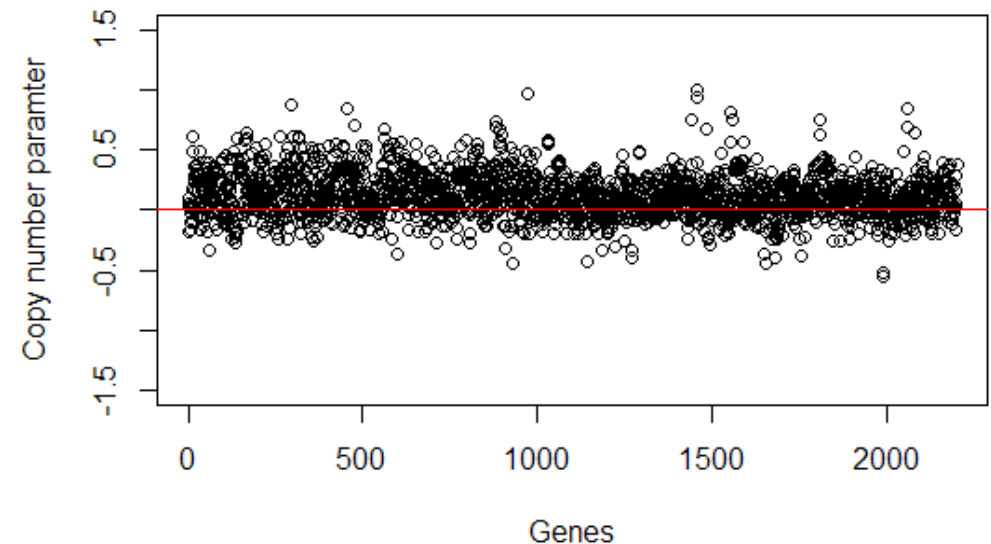
Multivariate prior: $\gamma | \sigma_\gamma^2 \sim \mathcal{N}_3(\mathbf{0}, \Sigma \sigma_\gamma^2)$ $\Sigma = \begin{bmatrix} 1 & \rho & \rho^2 \\ \rho & 1 & \rho \\ \rho^2 & \rho & 1 \end{bmatrix}$

- ϕ - hyper-parameters estimated univariate
- ρ - correlation estimated trivariate
- θ - model parameters estimated multivariate per triplets

Univariate parameter estimation



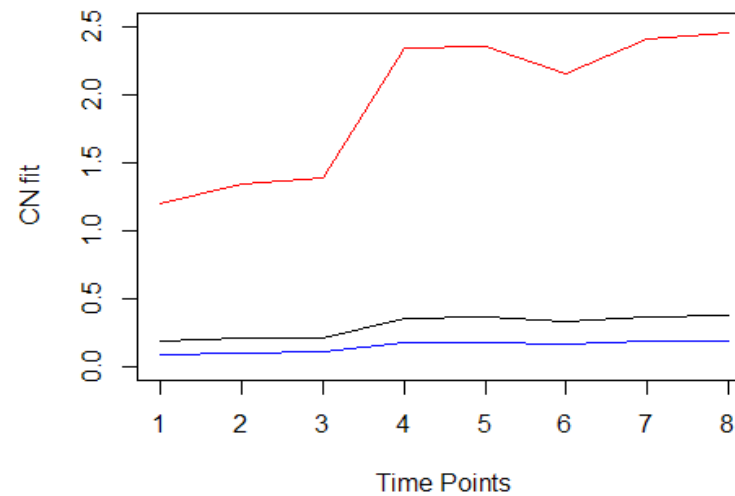
Multivariate spatial prior



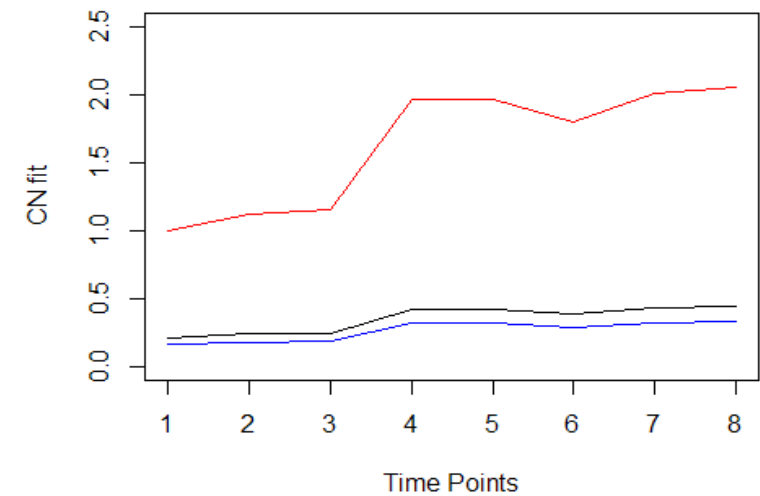
Effect of multivariate spatial prior

Triplet:

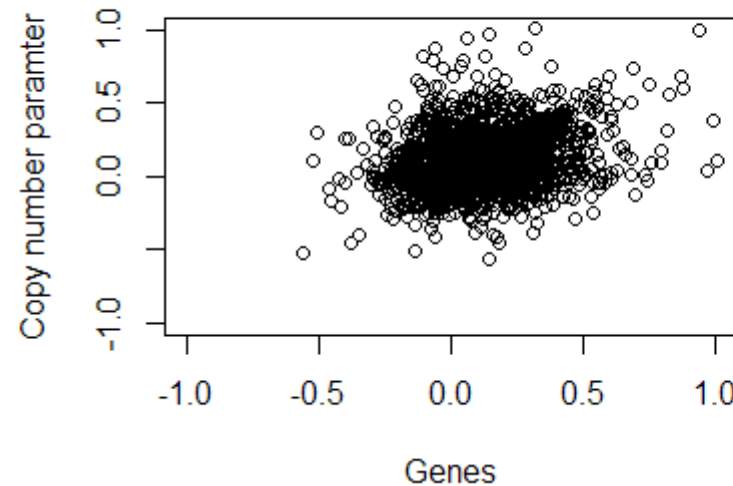
Univariate parameter estimation



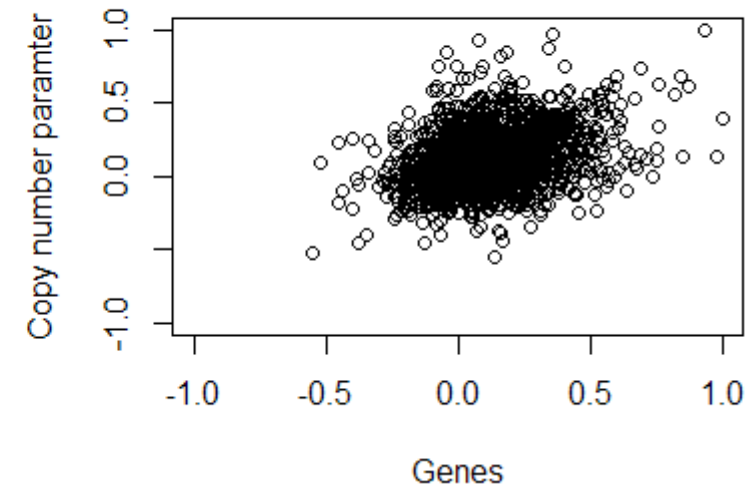
Multivariate spatial prior



Univariate parameter estimation



Multivariate spatial prior



Partial correlation:

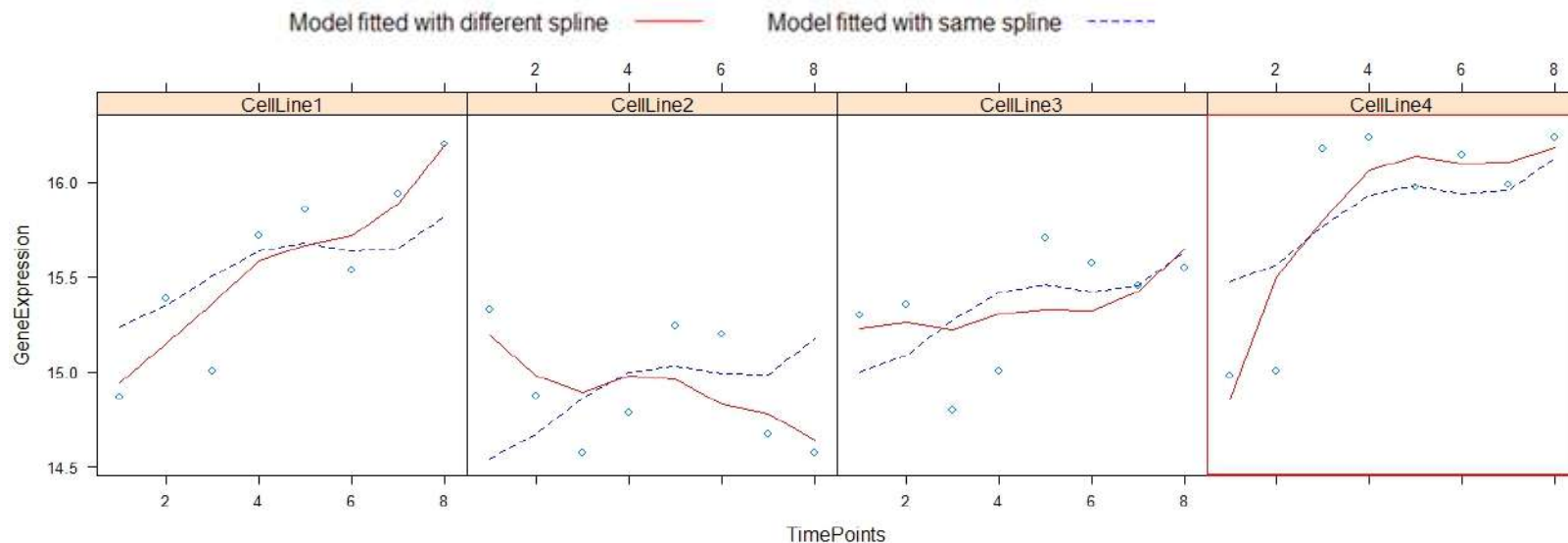
Same and different splines

Same spline – identify up or down regulated genes

$$\tilde{\mathbf{Z}} = \tilde{\mathbf{Z}} \otimes \mathbf{1}_{n \times n}$$

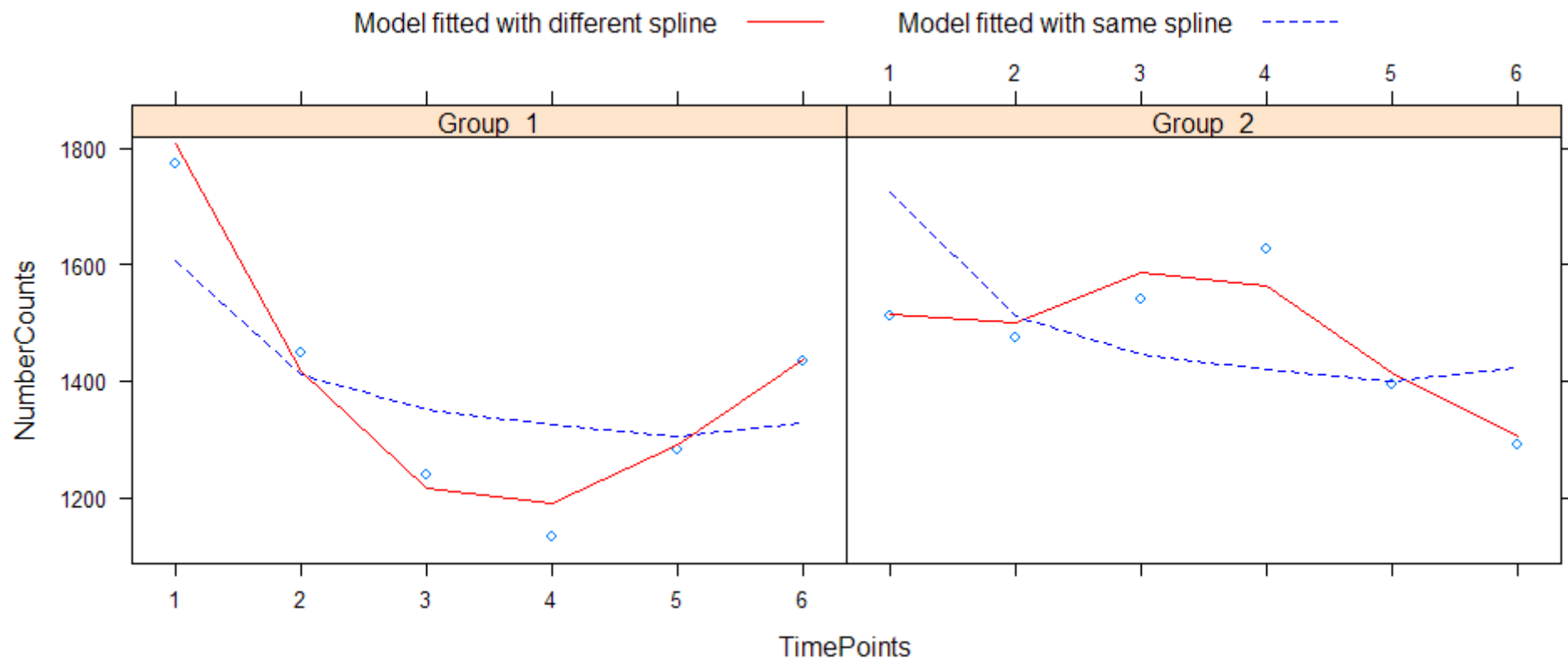
Different spline – allow more flexibility

$$\tilde{\mathbf{Z}} = \tilde{\mathbf{Z}} \otimes \mathbf{I}_{n \times n}$$



RNA-seq data

- Changing link function method can deal with count data.
- Two group time-course RNA-seq data.



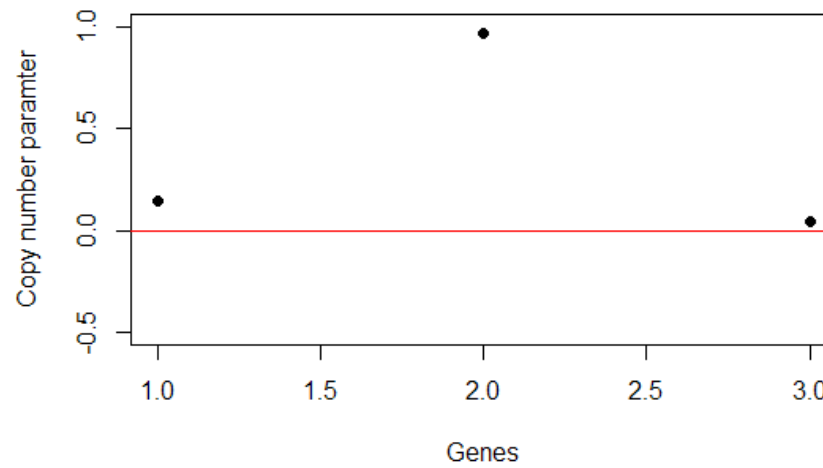
Summary

- Improved identification of temporal differential gene expression (TDGE) using penalized splines and empirical Bayes shrinkage.
- Identification of TDGE induced by CN.
- Identification of TDGE in count RNA-seq data.
- Improvement of CN estimates, with orthogonalization and imposing spatial multivariate prior.
- Identification of significant up or down regulated genes.
- As a proof of principle gene **SLC25A36** and **CADM1** are identified.

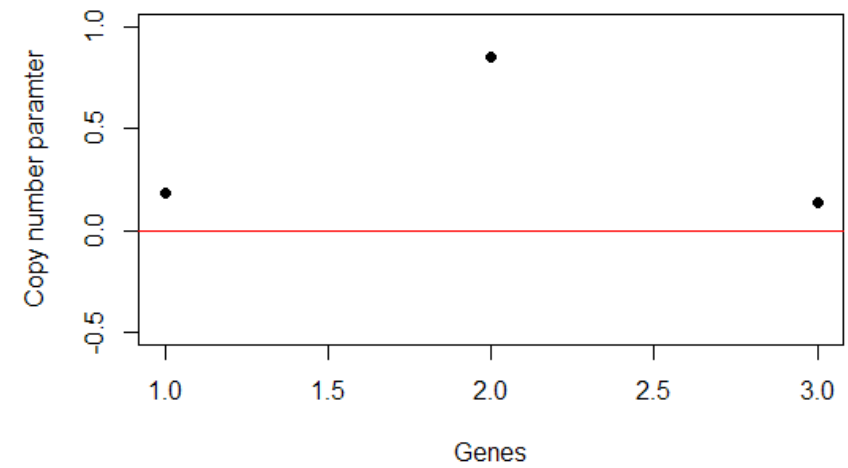
**Thank you for your
attention!**

Triplet:

Univariate parameter estimation



Multivariate spatial prior



Multivariate fit of triplets

