# Contributors

**Statistics**

- Viktorian Miok
- Wessel van Wieringen
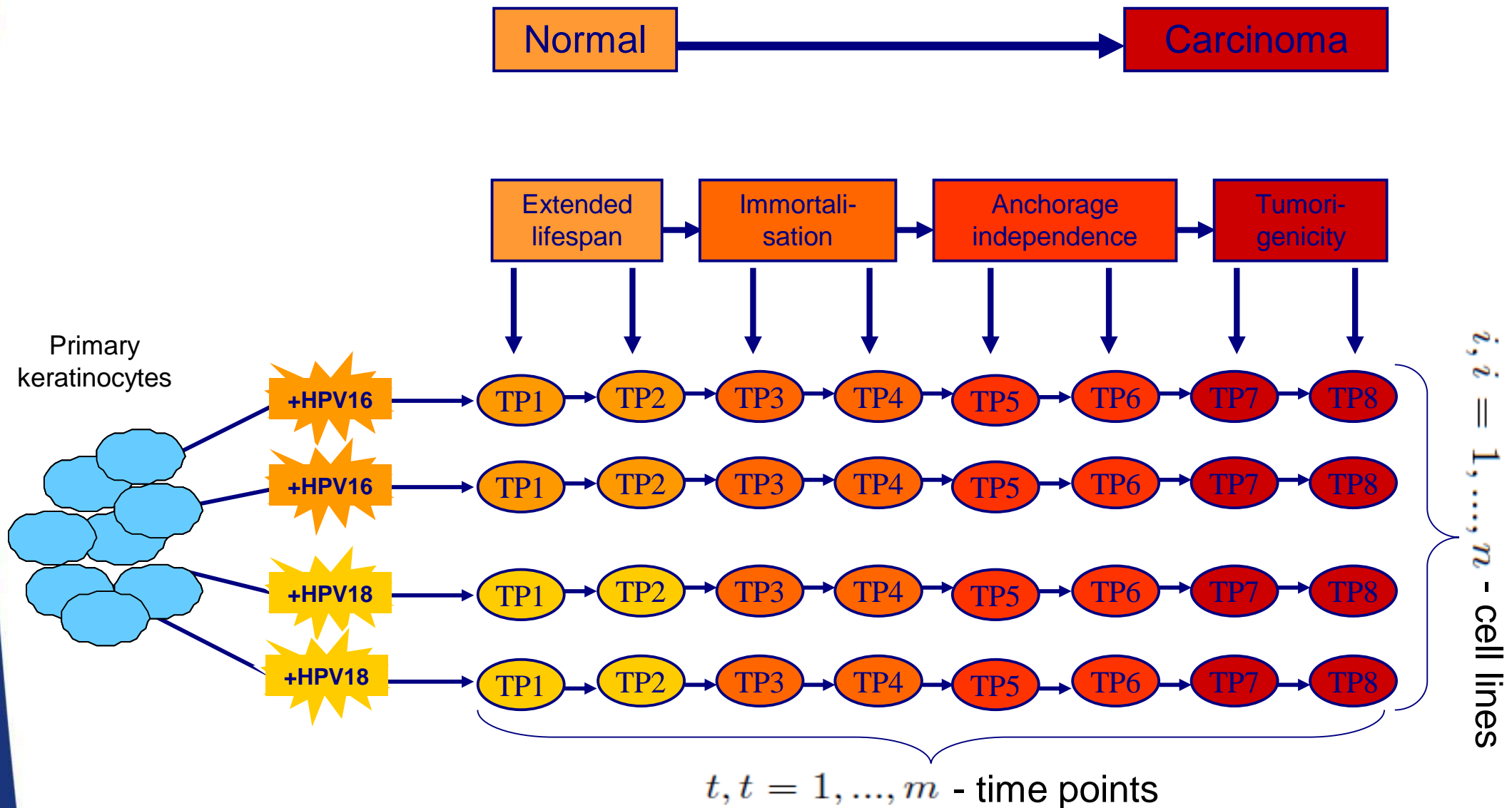- Mark van de Wiel

**Biology**

- Saskia Wilting
- Annelieke Jaspers
- Renske Steenbergen
- Peter Snijders
- Paula van Noort
- Ruud Brakenhoff

# 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

$j, j = 1, ..., p$ - genes

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

**Bayesian GLMM**: $Y_{i,j,t} \sim \mathcal{N}(\mu_{i,j,t}, \sigma_{\varepsilon,j}^2)$

Cell line effect   Time effect

$$\mu_{i,j,t} = \overbrace{f(i; \boldsymbol{\alpha}_j)} + \overbrace{h(t; \boldsymbol{\gamma}_j)}$$

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

# Fixed and random effects

Fixed effect:

Random effect:

$$f(i; \boldsymbol{\alpha}_j) = \boldsymbol{\alpha}_{i,j}$$

$$h(t; \boldsymbol{\gamma}_j) = \sum_{k=1}^{K} \gamma_{j,k} |t - \kappa_k|^3$$

Matrix notation: $\quad Y_{i,j,t} = \alpha_{i,j} + \tilde{\mathbf{Z}}_t \tilde{\boldsymbol{\gamma}}_j + \varepsilon_{i,j,t}$

$$\tilde{\mathbf{Z}}_t = \mathbf{Z}_t \mathbf{V}_\omega \mathbf{D}_\omega^{1/2}$$

$$\tilde{\boldsymbol{\gamma}}_j = \mathbf{D}_\omega^{-1/2} \mathbf{V}_\omega^{\mathrm{T}} \boldsymbol{\gamma}_j$$

Spline basis:

Spline coefficients:

$$\mathbf{Z}_t = (|t - \kappa_1|^3, \ldots, |t - \kappa_K|^3)$$

$$\boldsymbol{\gamma}_j = (\gamma_j, \ldots, \gamma_{j,K})^{\mathrm{T}}$$
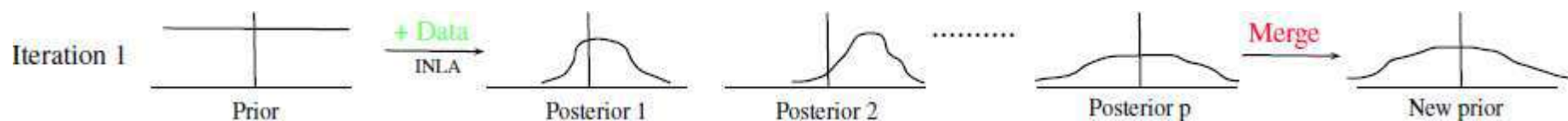
# INLA estimation

$\theta$   - model parameters

$\phi$   - hyper-parameters
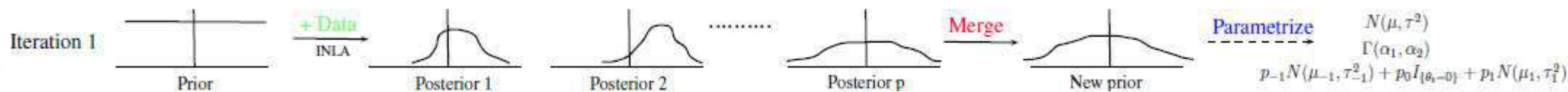
INLA (Rue et al., 2009) procedure consist in:

- Approximate full posterior of $\pi(\phi|y)$  and $\pi(\theta_l|\phi, y)$  using Laplace approximation.

- Approximate marginal posterior densities of $\theta$  and $\phi$  integrating over hyper-parameters of posteriors $\pi(\phi|y)$  and $\pi(\theta_l|\phi, y)$ .
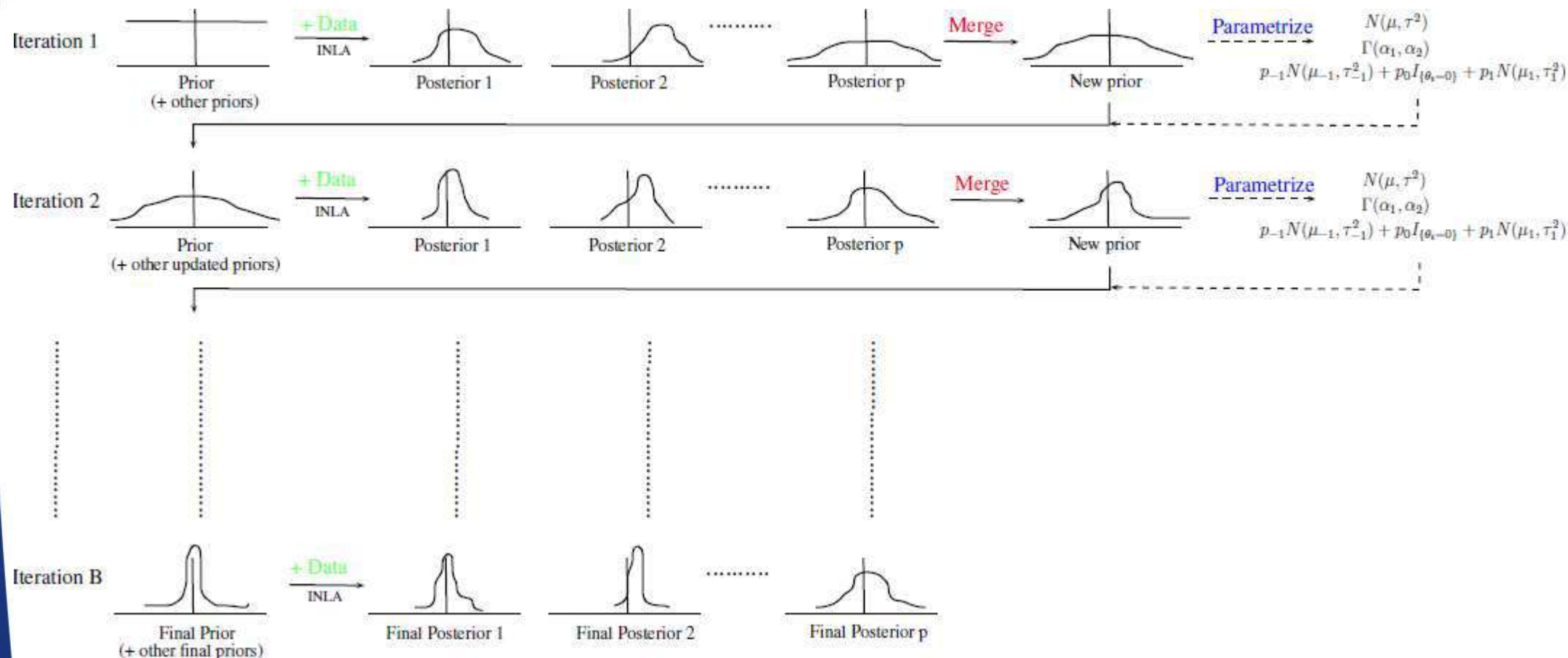
# Model parameters estimation



- Procedure start with flat prior

- Using data and INLA poster distributions are estimated

- Merge the posterior distribution – borrowing of the information

# Model parameters estimation



- Merged posterior distribution is used as prior distribution
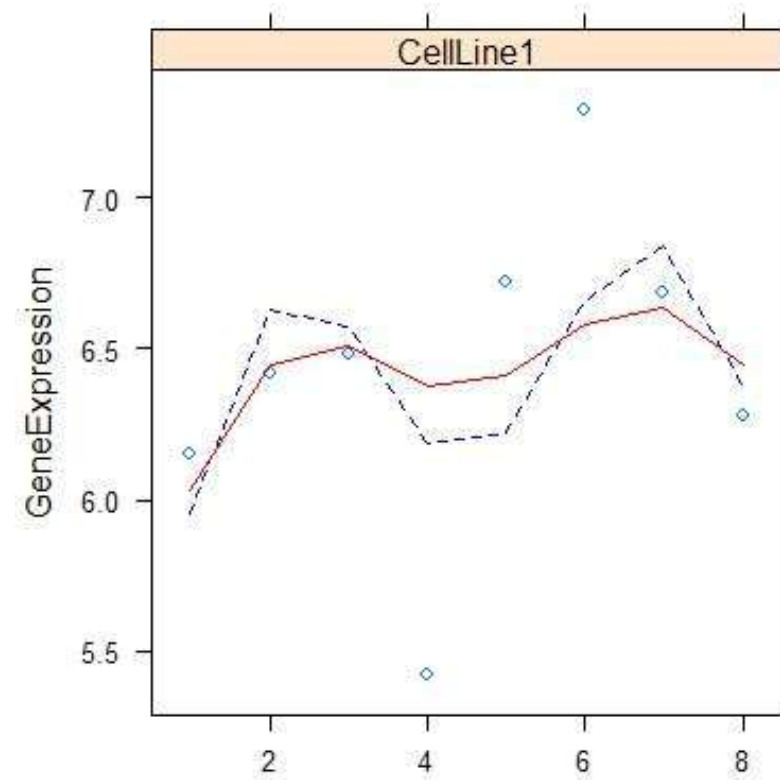
- New prior distribution is parametrized - shrinkage

# Model parameters estimation



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

# Shrinkage

- borrowing information across the genes

- better control of false positives

- leads to more stable estimates

- improvement of reproducibility

# Comparison of the methods

➢ **tigaR** – Miok et al., BMC Bioinformatics, 2014.

➢ **EDGE** – Storey et al., PNAS., 2005.

➢ **timecourse** – Tai and Speed, Annals of Statistics, 2006.

➢ **BATS** – Angelini et al., Stat. Appl. Genet. Mol. Biol., 2007.

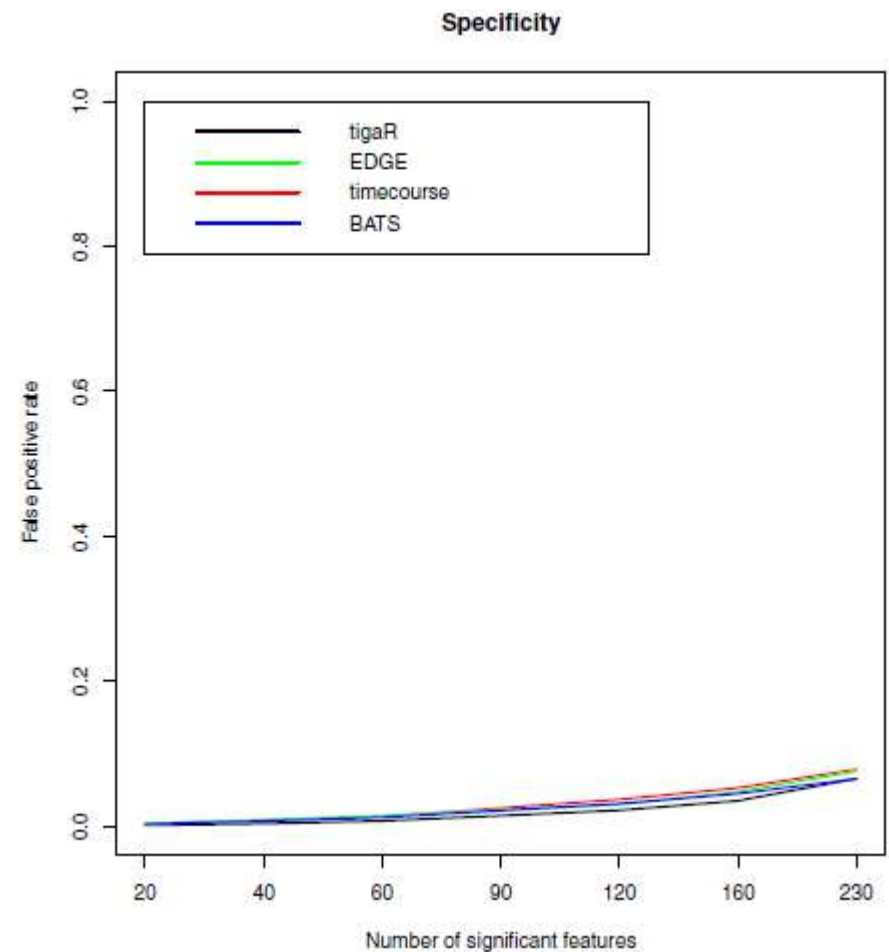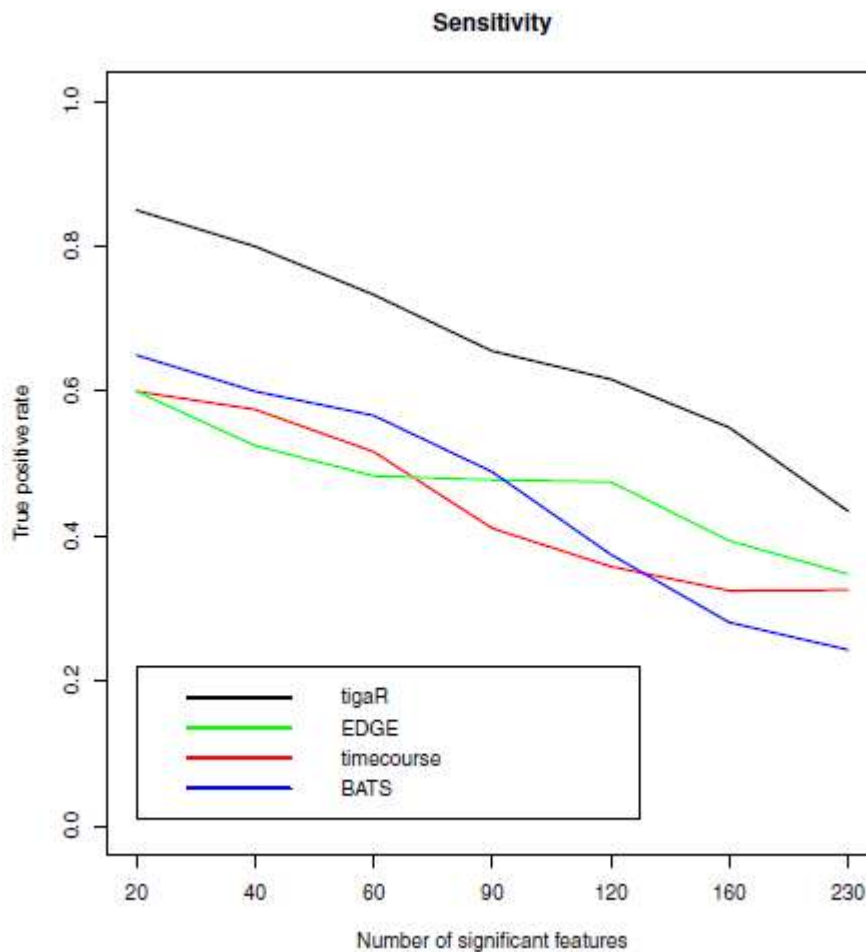# Comparison set-up

Data

> Real data from the experiment

Sensitivity and specificity

> Truth – significant genes among methods
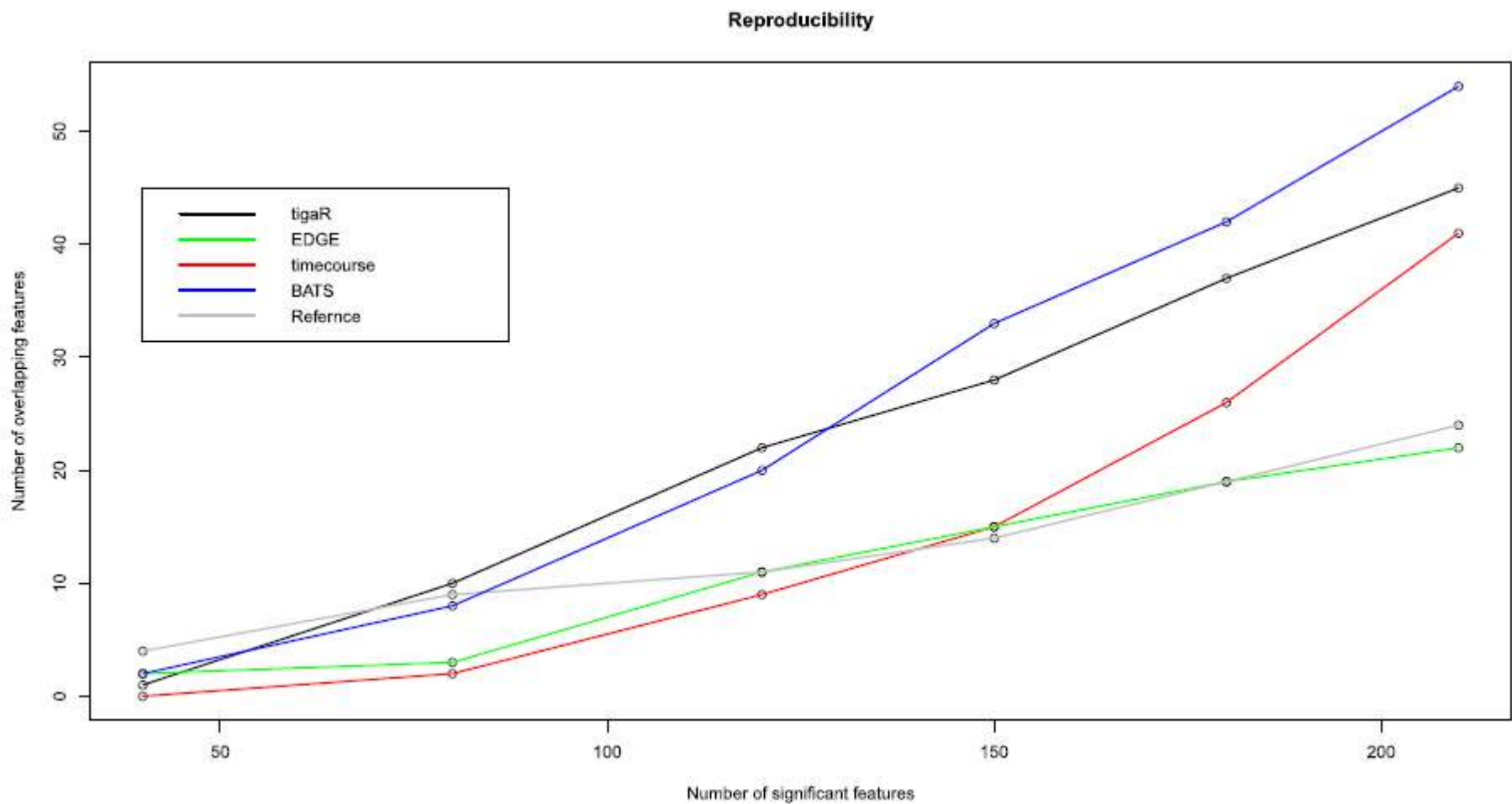
> Calculate true and false positive rate.

Reproducibility

> Equally divided data set in two groups

> Methods applied on the groups and calculate number of overlaps
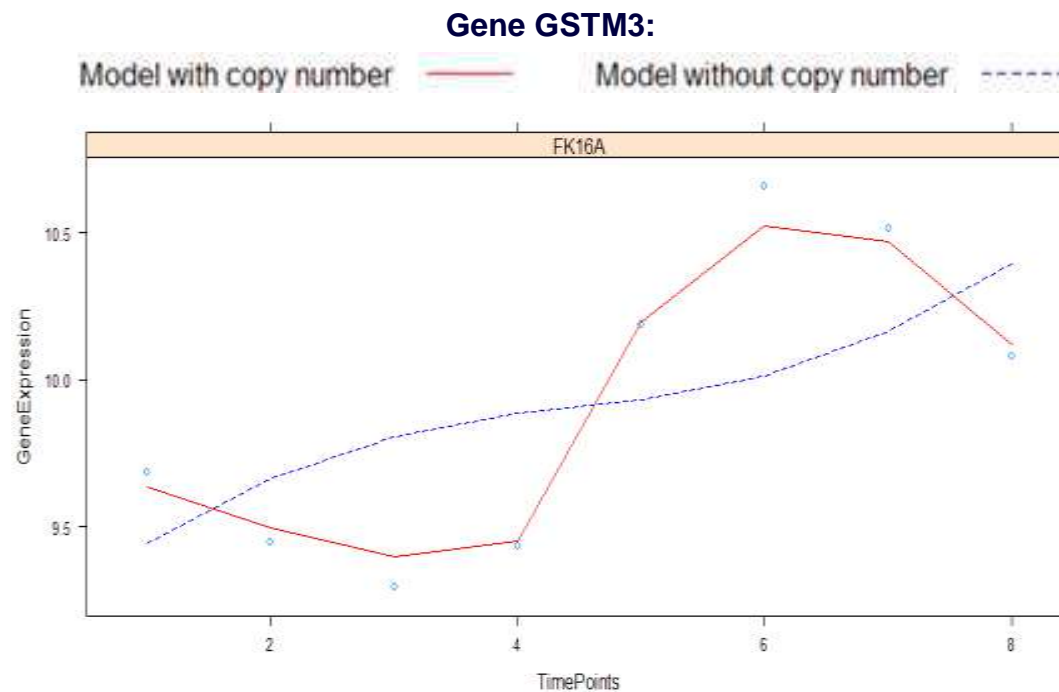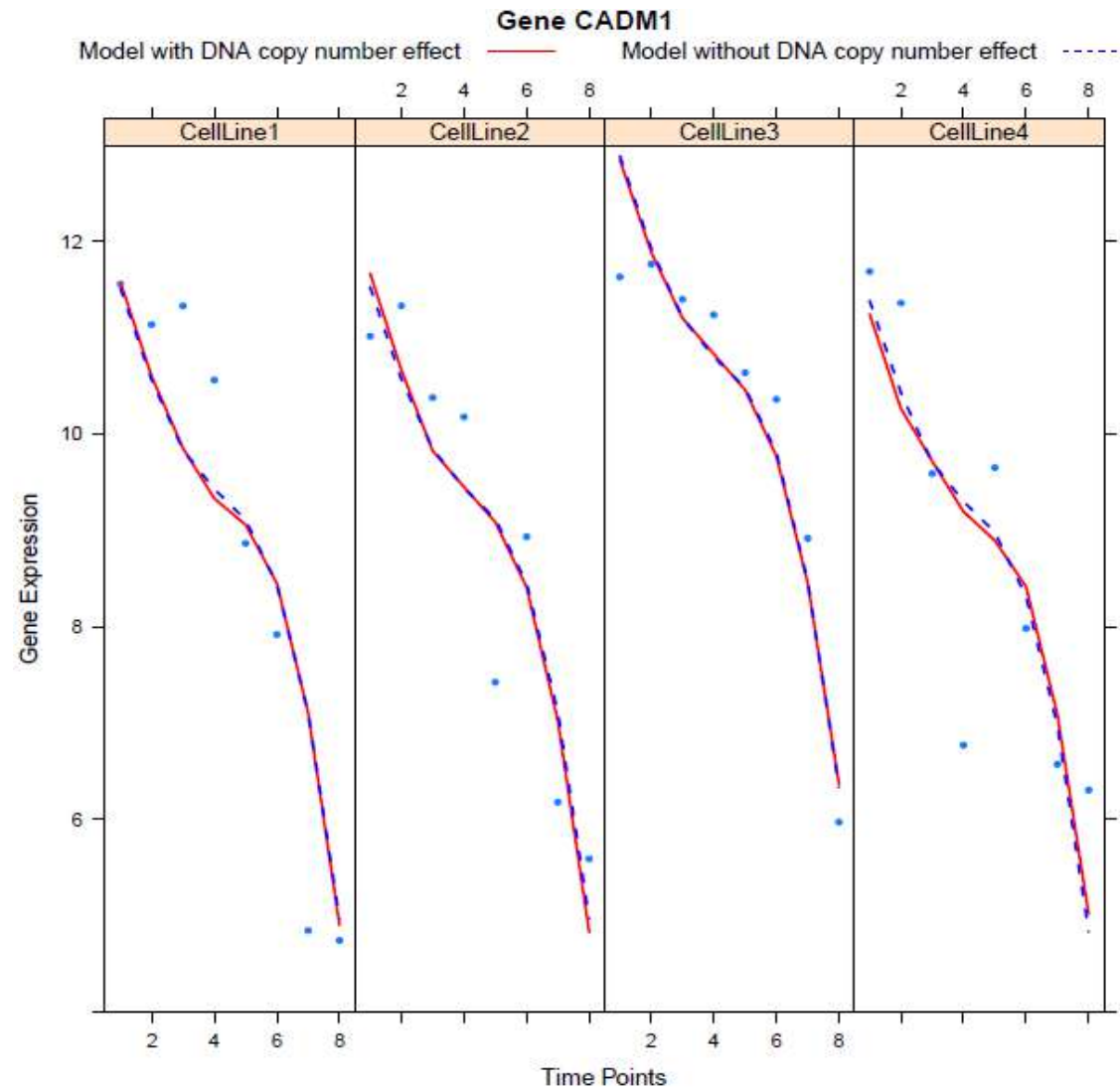
# Sensitivity and specificity

# Reproducibility



Reproducibility

# DNA copy number (CN)

$$\mathbf{X}_{*,*,t} \;=\; \left(\mathbf{X}_{1,*,t}, ..., \mathbf{X}_{n,*,t}\right)$$ - CN observations

Cell line · · · CN · · · · Time · · · Error

$$Y_{i,j,t} \;=\; \alpha_{i,j} + \beta_j\, x_{i,j,t} + \tilde{\mathbf{Z}}_t \tilde{\boldsymbol{\gamma}}_j + \varepsilon_{i,j,t}$$

**Gene GSTM3:**



16

# CADM1- gene without CN effect



Gene CADM1

Model with DNA copy number effect ——— Model without DNA copy number effect - - - - - -

**Overmeer et al., J Pathol., 2008.**

# SLC25A36 – gene with CN effect



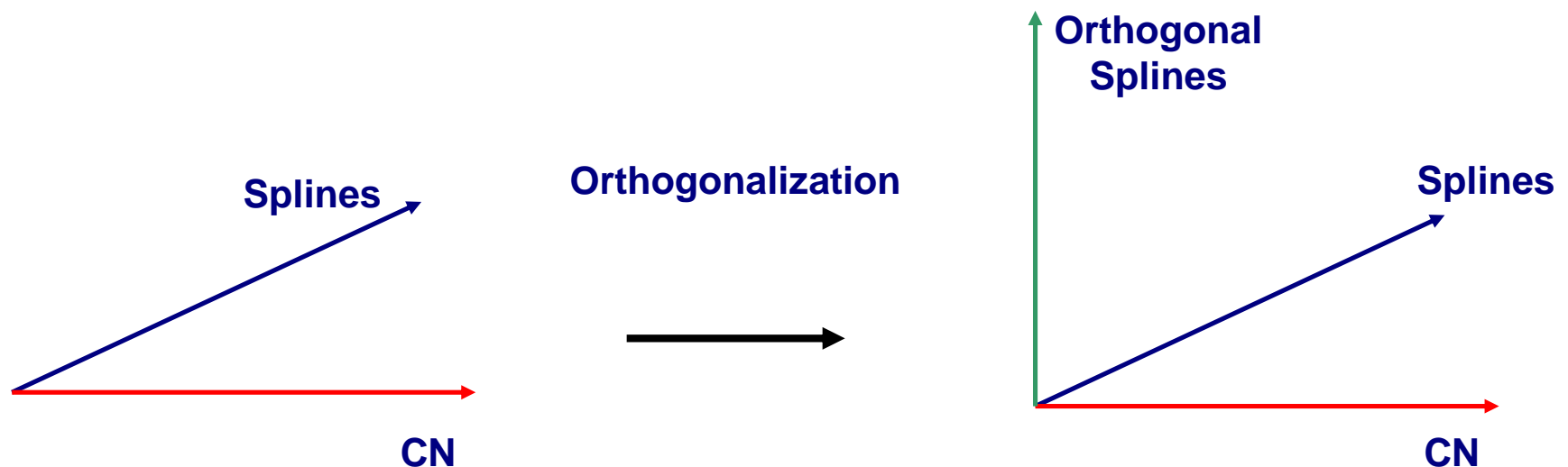**Wilting et al., Genes, Chromosomes and Cancer, 2008.**
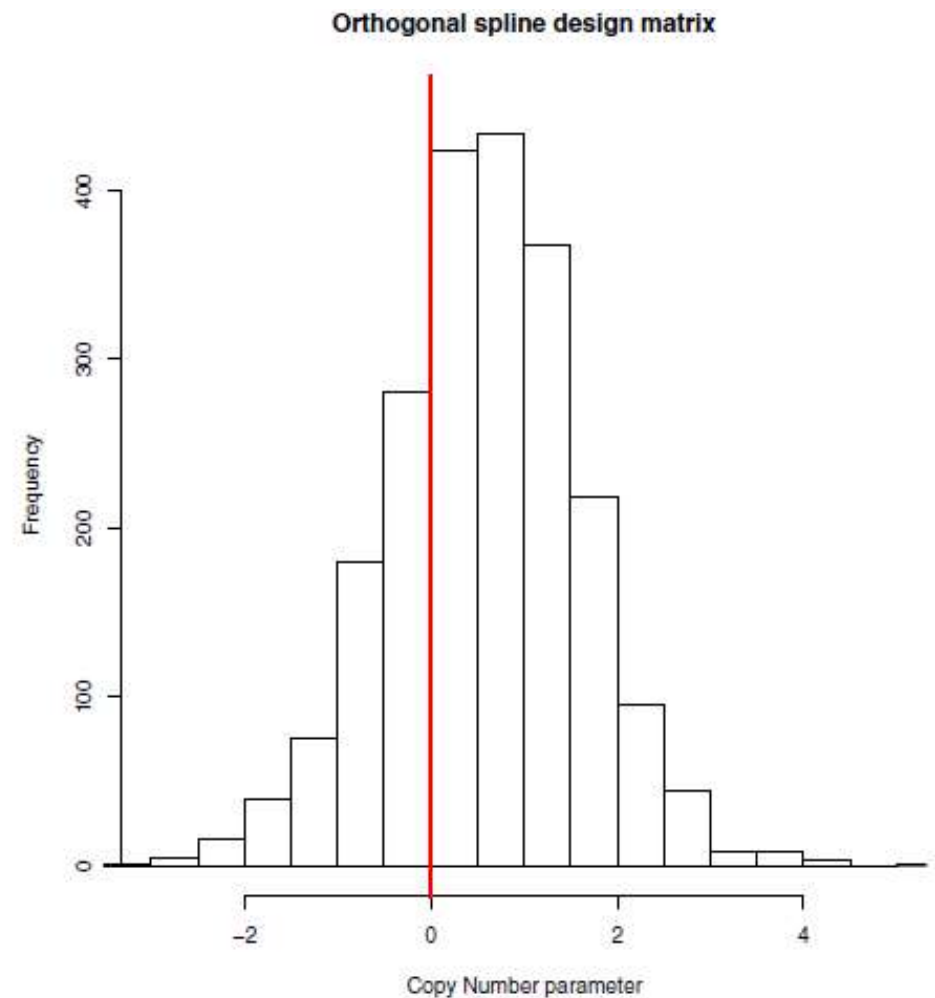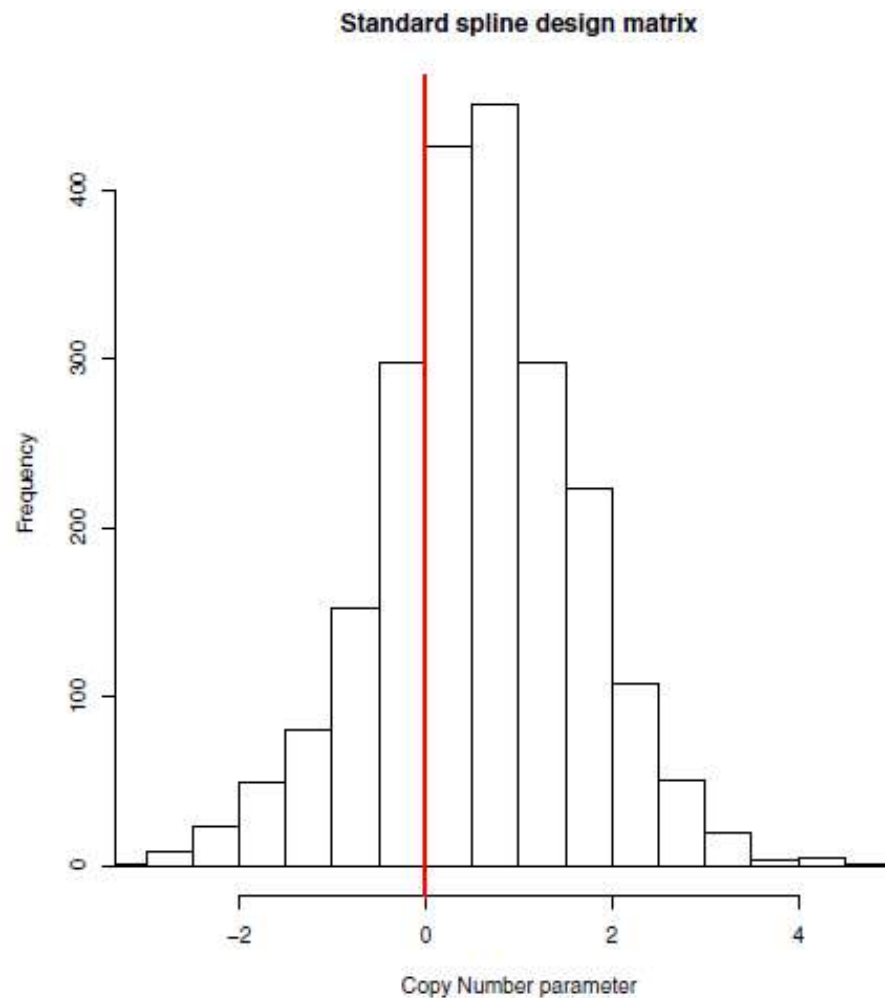
# Time vs. CN effect

Probelem:

➢ Flexibility of the splines(time) consumes effect of CN
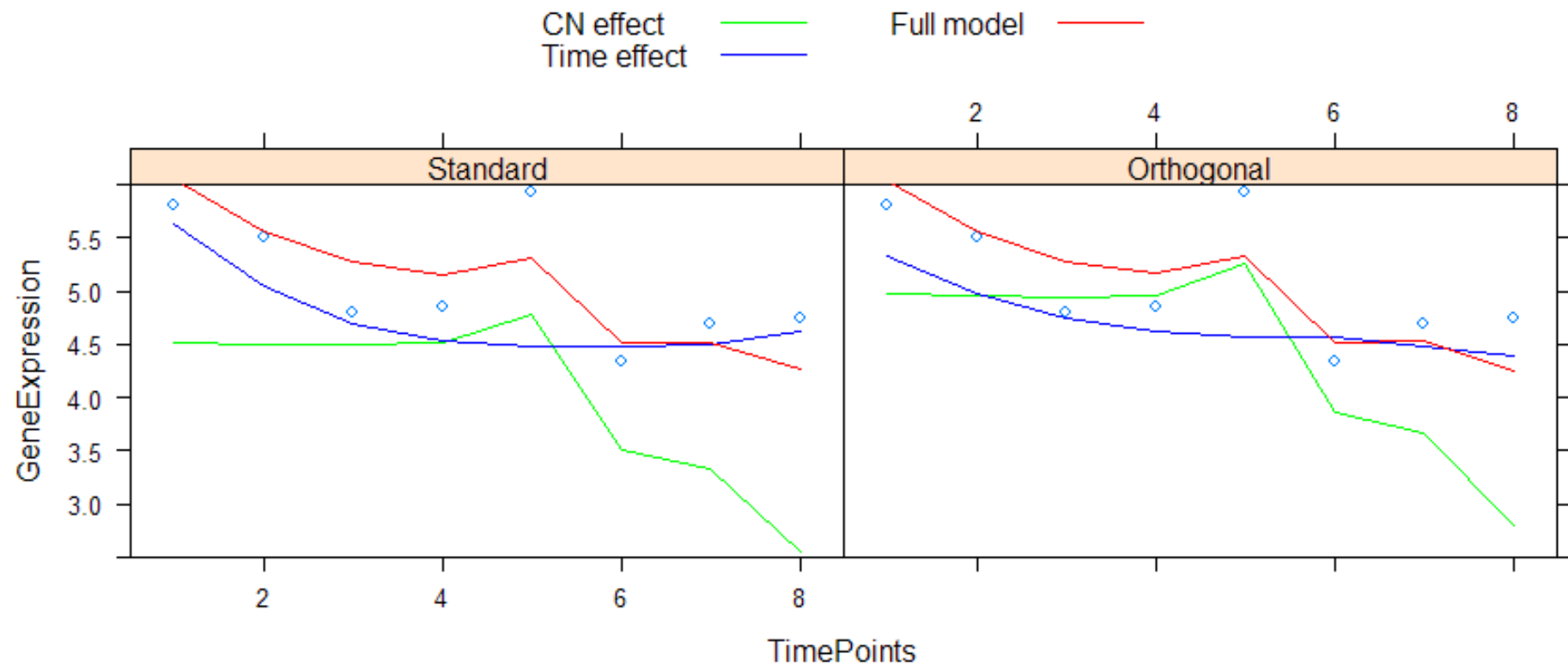
Potential solution

➢ Orthgonalization of the splines onto CN design matrix

# CN parameter

# Fit of the model

# Spatial multivariate prior for CN

$\beta_j$ follow the first-order autoregressive process along the genome:
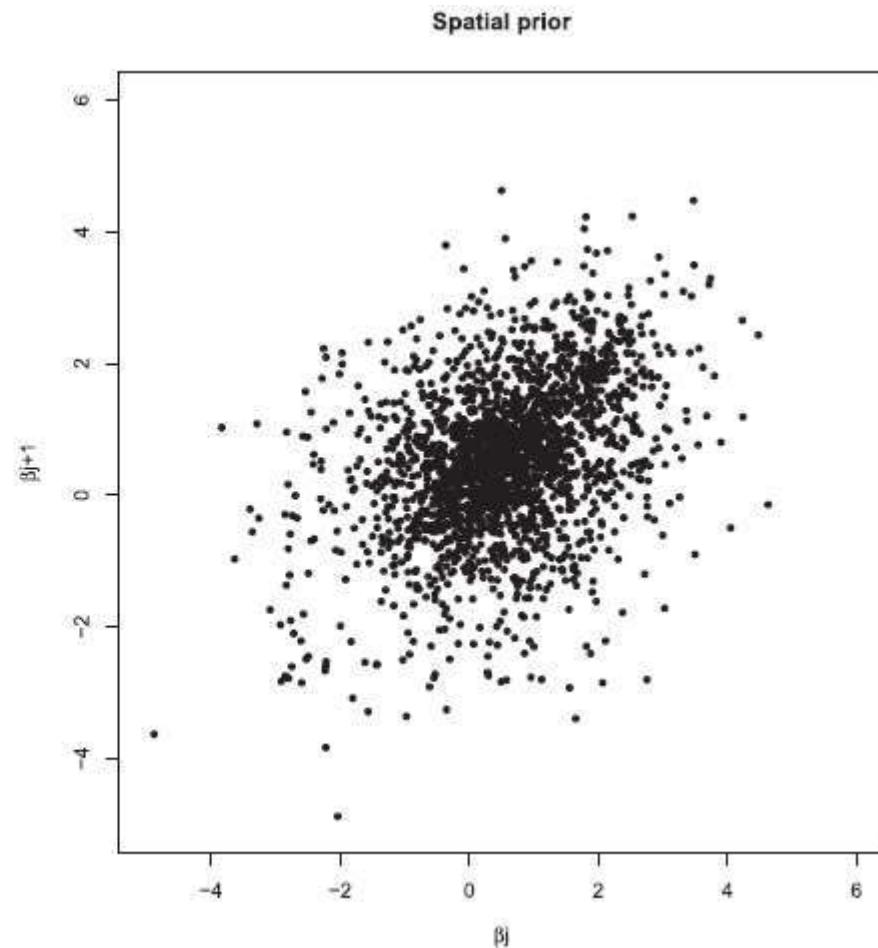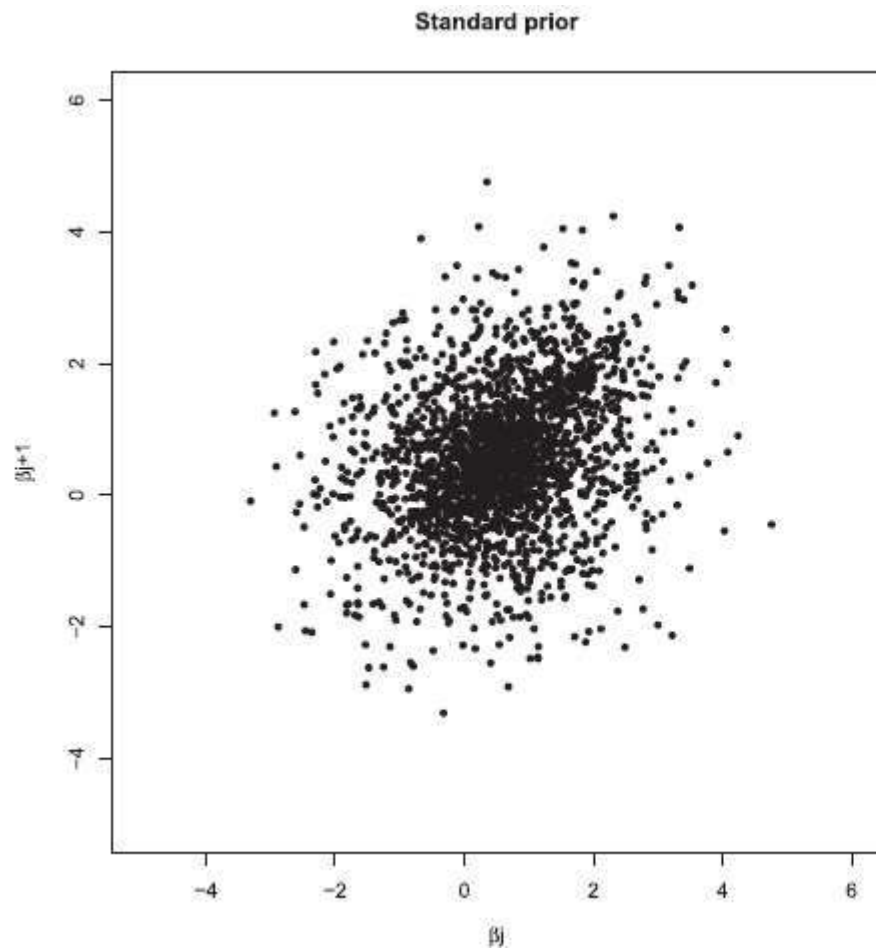
$$\beta_j = \rho\beta_{j-1} + \varepsilon_j$$

For each triplet trivariate normal prior is assumed:
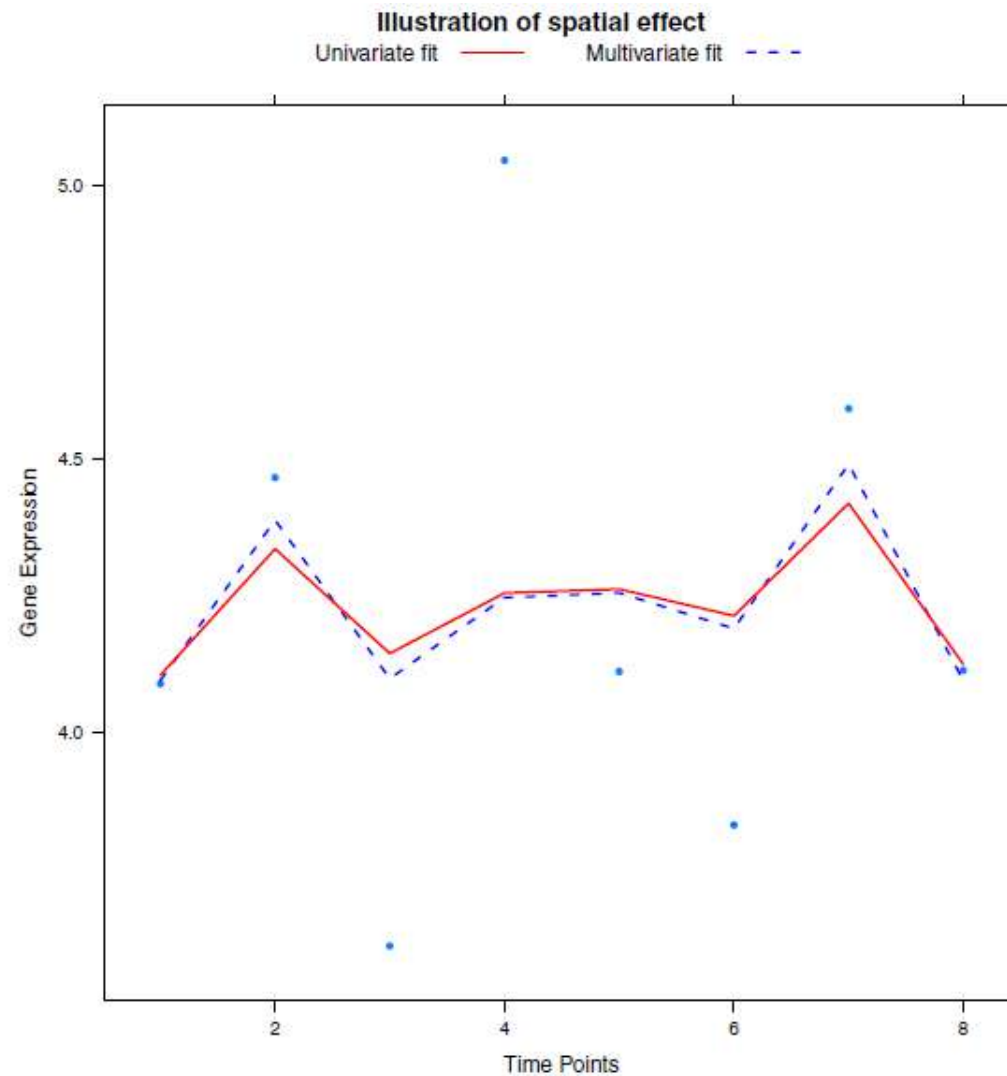
$$\begin{pmatrix} \beta_{j-1} \\ \beta_j \\ \beta_{j+1} \end{pmatrix} \sim \mathcal{N}\left( \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_{j-1}^2 & \sigma_{j-1}\sigma_j\rho & \sigma_{j-1}\sigma_{j+1}\rho^2 \\ \sigma_{j-1}\sigma_j\rho & \sigma_j^2 & \sigma_j\sigma_{j+1}\rho \\ \sigma_{j-1}\sigma_{j+1}\rho^2 & \sigma_j\sigma_{j+1}\rho & \sigma_{j+1}^2 \end{pmatrix} \right)$$
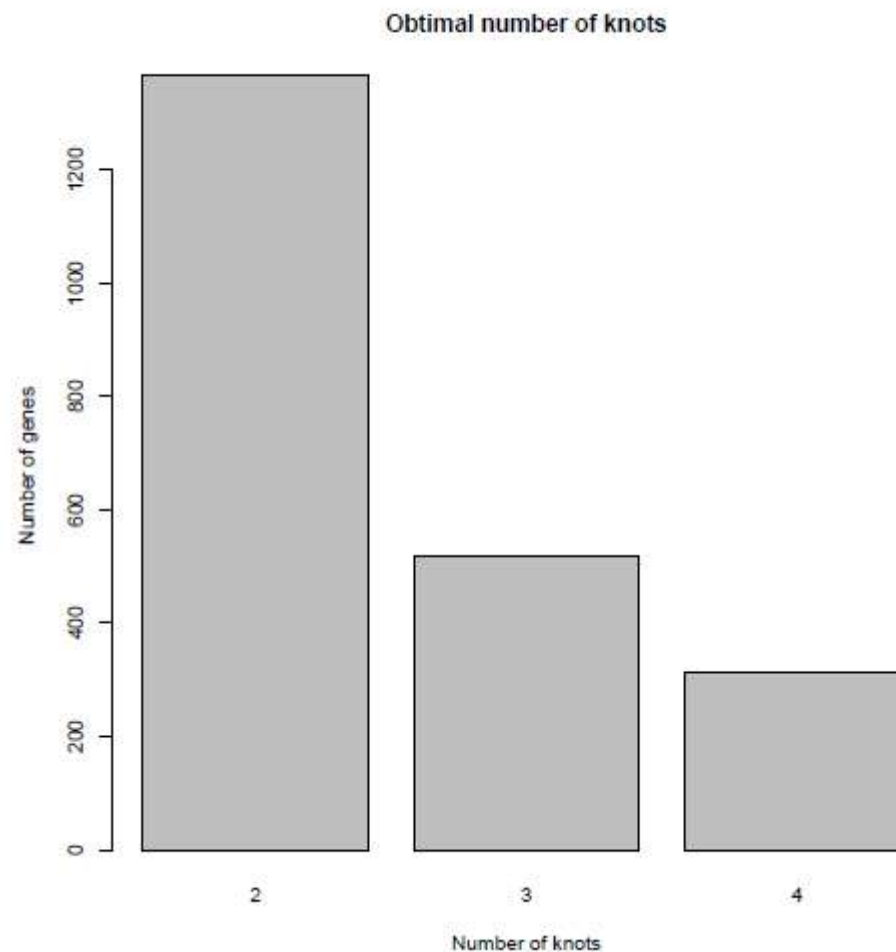
# CN parameters

Partial correlation of CN parameter:



Standard prior

Spatial prior

# Fit of the model

# Optimal number of knots for splines
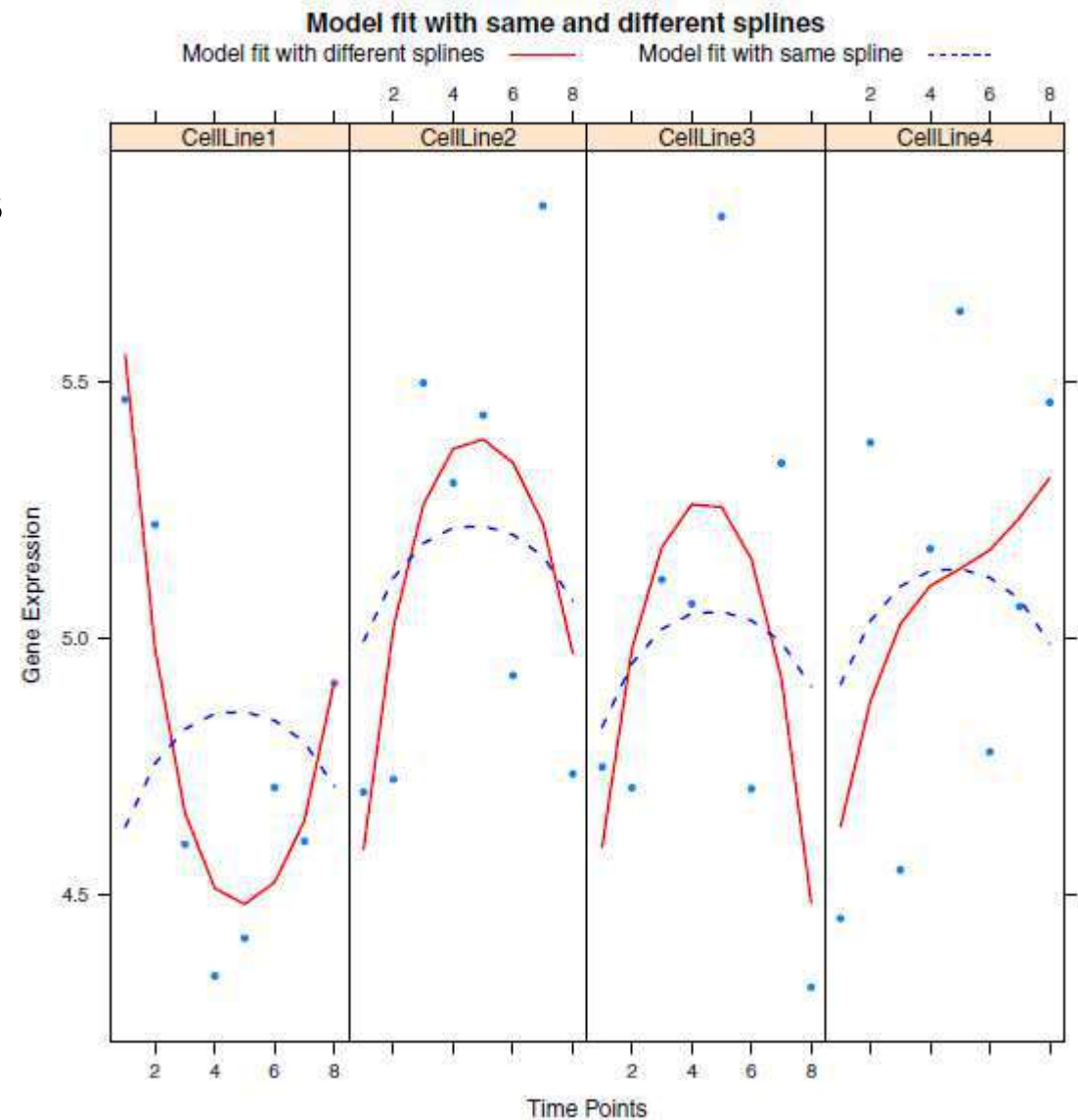


Obtimal number of knots

# Splines flexibility

Same spline – up/down regulated genes

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

Different spline – allow more flexibility

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



Model fit with same and different splines

# Application

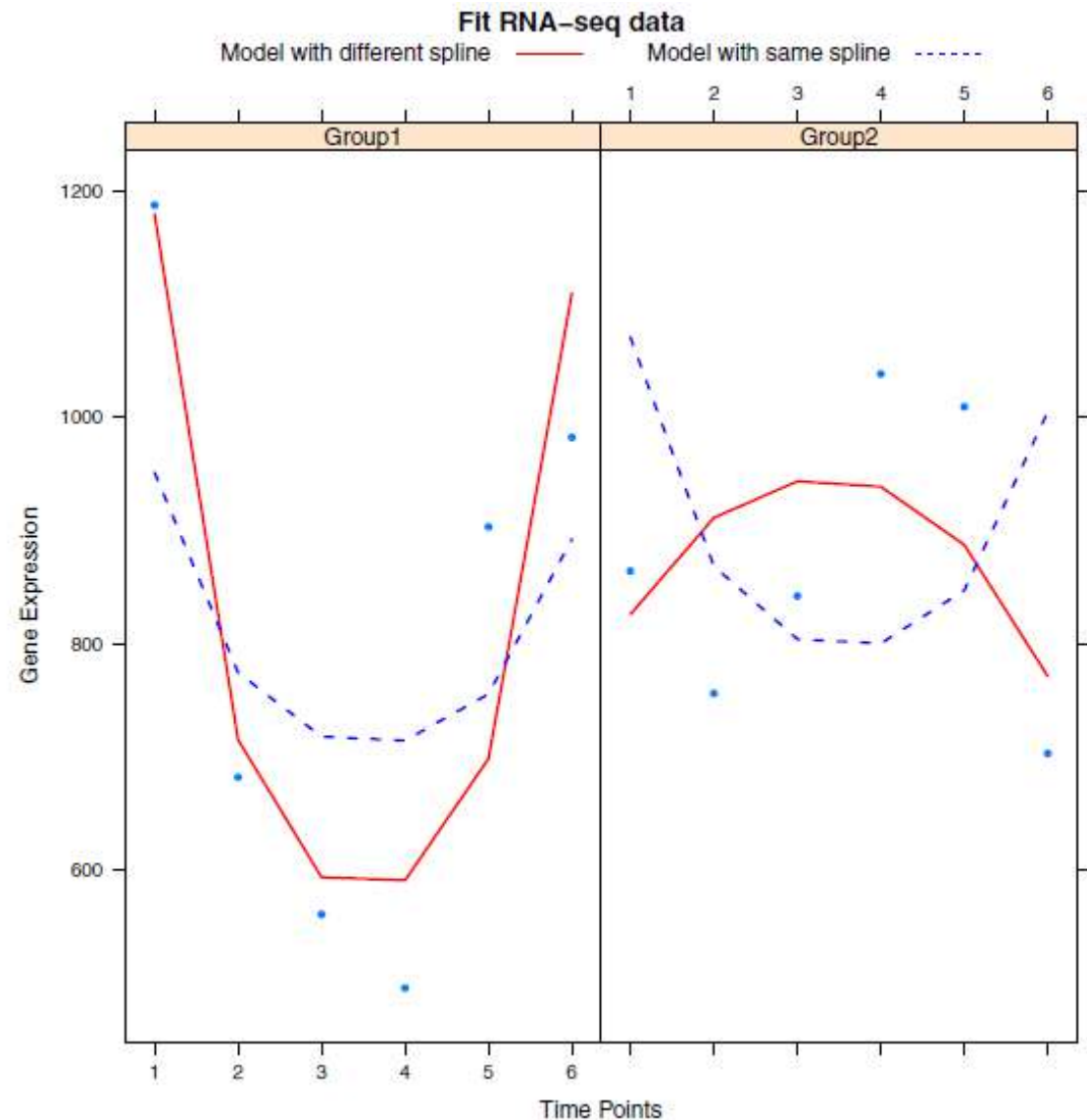| Effect | Model | Same spline | | Different spline | |
|---|---|---|---|---|---|
| | | Standard | Orthogonal | Standard | Orthogonal |
| Time | Splines | 417 | | 583 | |
| | CN+Splines | 204 | 203 | 421 | 421 |
| CN | CN+Splines | 402 | 403 | 380 | 380 |
| | Multivariate | 398 | 399 | 377 | 380 |

Analysis is performed only on 2202 features, which represent one chromosome.

Method identify genes with time and CN effect allowing for:
- ➢ flexibility in modeling of time effect
- ➢ additional stability of CN parameters

# RNA-seq data

- Changing link function method can deal with count data.

- Two group time-course RNA-seq data.



**Fit RNA–seq data**

Model with different spline ——    Model with same spline ------

# 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 orhogonalization and imposing spatial multivariate prior.

- Identification of significant up or down regulated genes.

- As a proof of principle biologically relevant genes **SLC25A36** and **CADM1** are identified.

# Thank you for your attention!