

# The use of mixture distributions in a Bayesian linear mixed effects model

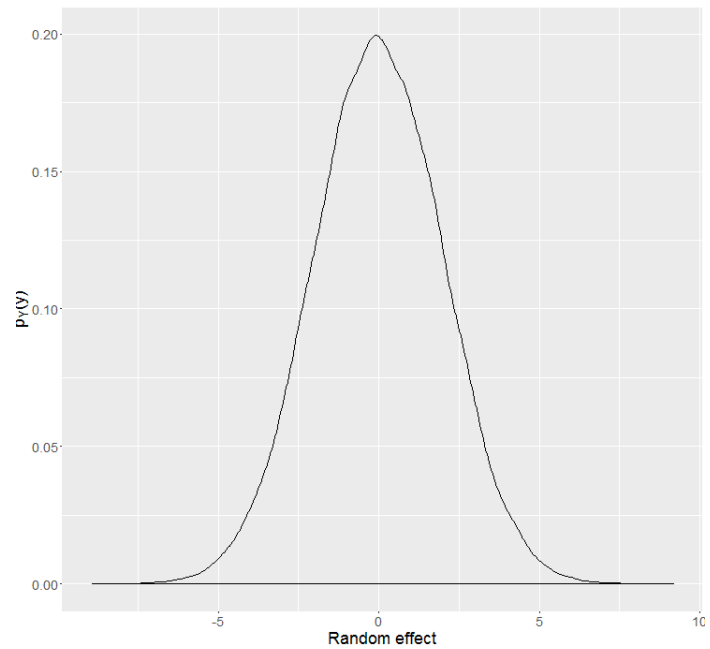
---

Anirudh TOMER

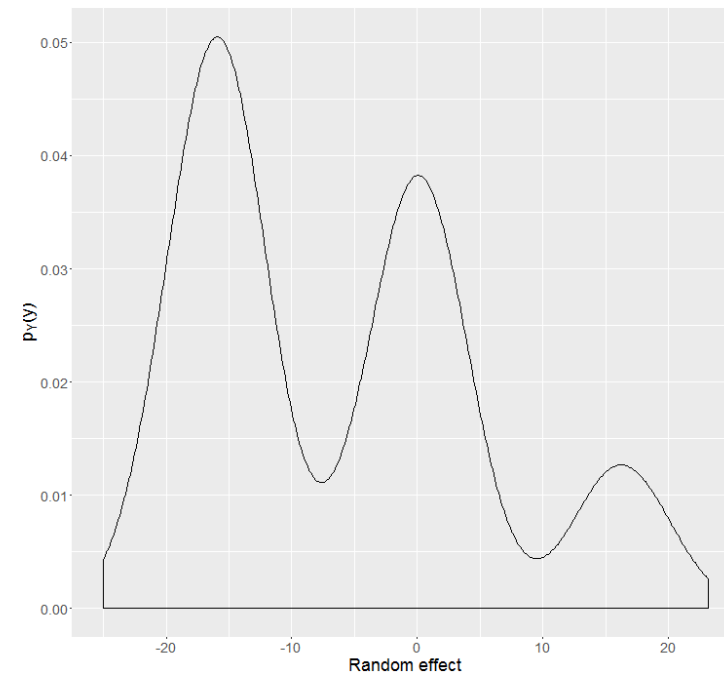
Promoter: Professor Emmanuel LESAFFRE

# Bayesian heterogeneity model

---



Vs.



# What is the problem we are facing?

---

Estimation of mixture component density parameters  $\theta_G$

Criteria for choice of number of components (DIC, Marginal likelihood, Predictive methods)

Classification of observations into groups

# DIC (deviance information criteria)

---

$$DIC = \overline{D(\boldsymbol{\theta})} + p_D$$

Where  $p_D = \overline{D(\boldsymbol{\theta})} - D(\bar{\boldsymbol{\theta}})$  is the effective number of parameters in the model.

$D(\boldsymbol{\theta}) = -2 \log p(\mathbf{y}|\boldsymbol{\theta}) + 2 \log f(\mathbf{y})$  is called the Bayesian deviance.

Various definitions of DIC proposed by Celeux et al., (2006) for missing data models.

# Marginal/Complete/Conditional likelihood

---

$$p(y|\theta) = \prod_{i=1}^n \sum_{k=1}^K f_N(y_i; X_i\beta + Z_i b_k^C, Z_i G_k Z_i^T + R_i) \eta_k$$

$$\begin{aligned} f(y^F|\theta) &= \prod_{i=1}^n f(y_i|b_i, S_i) f(b_i|S_i) f(S_i) \\ &= \prod_{i=1}^n f_N(y_i; X_i\beta + Z_i b_i, R_i) f_N(b_i; b_{S_i}^C, G_{S_i}) f(S_i; \eta) \end{aligned}$$

$X_i, Z_i$   
 $\beta, b_k^C$   
 $G_k, R_i$   
 $\eta_k, S_i, b_i$

?

$$\begin{aligned} p(y|\theta^{\text{cond}}) &= \prod_{i=1}^n f_N(y_i; X_i\beta + Z_i b_i, R_i) \\ \theta^{\text{cond}} &= (\theta, b, S). \end{aligned}$$

# Marginal DIC definitions

---

$$\text{DIC}_1 = -4\mathbb{E}_{\theta|y}(\log p(\mathbf{y}|\theta)) + 2\log p(\mathbf{y}|\bar{\theta})$$

$$\text{DIC}_2 = -4\mathbb{E}_{\theta|y}(\log p(\mathbf{y}|\theta)) + 2\log p(\mathbf{y}|\hat{\theta}_M)$$

$$\hat{\theta}_M = \arg \max_{\theta} p(\theta|\mathbf{y})$$

$$\text{DIC}_3 = -4\mathbb{E}_{\theta|y}(\log p(\mathbf{y}|\theta)) + 2\log \hat{p}(\mathbf{y})$$

$$\hat{p}(\mathbf{y}) = \prod_{i=1}^n \hat{p}(\mathbf{y}_i) = \prod_{i=1}^n \frac{1}{m} \sum_{l=1}^m \sum_{k=1}^K f_N(\mathbf{y}_i; \mathbf{X}_i \beta^{(l)} + \mathbf{Z}_i \mathbf{b}_k^{C(l)}, \mathbf{Z}_i \mathbf{G}_k^{(l)} \mathbf{Z}_i^T + \mathbf{R}_i^{(l)}) \eta_k^{(l)}$$

# Complete and conditional DIC definitions

---

$$\text{DIC}_4 = -4\mathbb{E}_{\theta,b,S|y}(\log p(\mathbf{y}^F|\theta)) + 2\mathbb{E}_{b,S|y}(\log p(\mathbf{y}^F|\bar{\theta}))$$

$$\text{DIC}_5 = -4\mathbb{E}_{\theta,b,S|y}(\log p(\mathbf{y}^F|\theta)) + 2\log p(\mathbf{y}, \hat{\mathbf{b}}_M, \hat{\mathbf{S}}_M|\hat{\boldsymbol{\theta}}_M)$$

$$\text{DIC}_6 = -4\mathbb{E}_{\boldsymbol{\theta}^{\text{cond}}}(\log p(\mathbf{y}|\boldsymbol{\theta}^{\text{cond}})) + 2\log p(\mathbf{y}|\hat{\boldsymbol{\theta}}_M^{\text{cond}})$$

$$\boldsymbol{\theta}^{\text{cond}} = (\theta, b, S).$$

$$\hat{\boldsymbol{\theta}}_M^{\text{cond}} = \arg \max_{\boldsymbol{\theta}^{\text{cond}}} p(\boldsymbol{\theta}^{\text{cond}}|\mathbf{y}), \text{ and } \mathbb{E}_{\boldsymbol{\theta}^{\text{cond}}}(\log p(\mathbf{y}|\boldsymbol{\theta}^{\text{cond}}))$$

# Marginal likelihood using Chib's approximation

$$m(\mathbf{y}) = p(\mathbf{y}|M) = \frac{L(\boldsymbol{\theta}|\mathbf{y}, M)p(\boldsymbol{\theta}|M)}{p(\boldsymbol{\theta}|\mathbf{y}, M)}$$

$$\log \hat{m}(\mathbf{y}) = \log L(\boldsymbol{\theta}^*|\mathbf{y}) + \log p(\boldsymbol{\theta}^*) - \log p(\boldsymbol{\theta}^*|\mathbf{y})$$

$$\begin{aligned} \log p(\boldsymbol{\theta}^*|\mathbf{y}) = & \sum_{k=1}^K \log p(G_k^*|\mathbf{y}) + \sum_{k=1}^K \log p(\mathbf{b}_k^{C*}|G_k^*, \mathbf{y}) + \log p(\sigma^{2*}|G_k^*, \mathbf{b}_k^{C*}, \mathbf{y}) \\ & + \log p(\beta^*|G_k^*, \mathbf{b}_k^{C*}, \sigma^{2*}, \mathbf{y}) + \log p(\eta^*|G_k^*, \mathbf{b}_k^{C*}, \sigma^{2*}, \beta^*, \mathbf{y}) \end{aligned}$$

$$\begin{aligned} \prod_{k=1}^K p(G_k^*|\mathbf{y}) &= \int \prod_{k=1}^K p(G_k^*|\mathbf{y}, \mathbf{b}, \mathbf{S}, \mathbf{b}_k^C) p(\mathbf{b}_1^C, \mathbf{b}_2^C, \dots, \mathbf{b}_K^C, \mathbf{b}, \mathbf{S}|\mathbf{y}) d\mathbf{b}_1^C d\mathbf{b}_2^C \dots d\mathbf{b}_K^C d\mathbf{b} d\mathbf{S} \\ &\approx \frac{1}{m} \sum_{l=1}^m \prod_{k=1}^K p(G_k^*|\mathbf{y}, \mathbf{b}^{(l)}, \mathbf{S}^{(l)}, \mathbf{b}_k^{C(l)}) \\ &\approx \frac{1}{m} \sum_{l=1}^m \prod_{k=1}^K f_{\mathcal{W}^{-1}}(G_k^*; n_k^{(l)} + n_0, \Psi + \sum_{i=1}^{n_k^{(l)}} (\mathbf{b}_i^{(l)} - \mathbf{b}_k^{C(l)})(\mathbf{b}_i^{(l)} - \mathbf{b}_k^{C(l)})^T) \end{aligned}$$



# Posterior predictive checks

---

$$T(\mathbf{r}) = \frac{1}{\sum_{i=1}^n m_i} \sum_{i=1}^n \sum_{j=1}^{m_i} (r_{ij} - \bar{r}_{i.})^2$$

$$r_{ij} = z_{ij} \tilde{\mathbf{b}}_i + \varepsilon_{ij} = y_{ij} - x_{ij} \beta,$$

$$\bar{r}_{i.} = \frac{1}{m_i} \sum_{j=1}^{m_i} r_{ij}.$$

Motivation: Testing for overfitting.

Idea: Sample big values from empty components to obtain inflated test statistic.

# Data sets

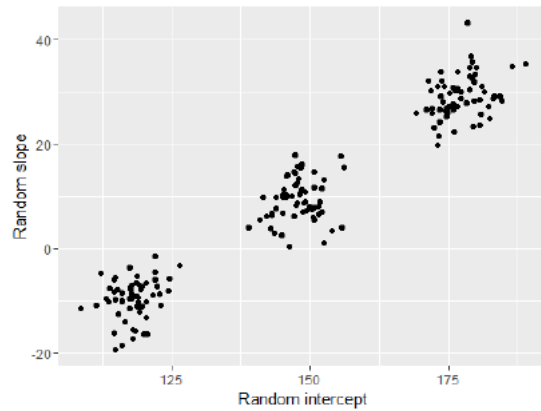
---

## Zebu cow weights

Predictors were gender, birth year, age, time of measurement.

Various versions of the dataset which differed in number of mixture components for random effects, number of subjects, separation of mixture components and number of subjects per component.

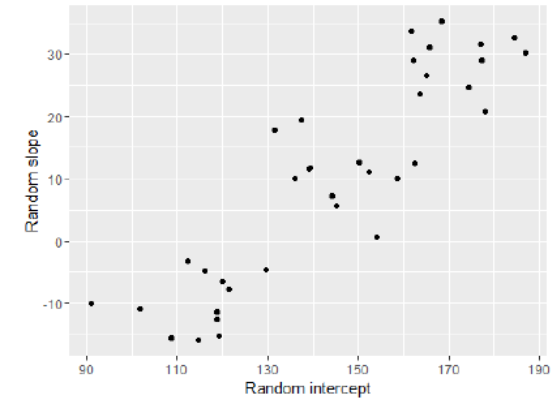
# DIC results



(a) Data set 2



# Components Fitted	DIC <sub>1</sub>	DIC <sub>2</sub>	DIC <sub>3</sub>	DIC <sub>4</sub>	DIC <sub>5</sub>	DIC <sub>6</sub>
1	9966	9959	9965	12921	10531	7855
2	9865	9849	9864	12498	10458	7860
3	9664	9665	9663	11847	10244	7870



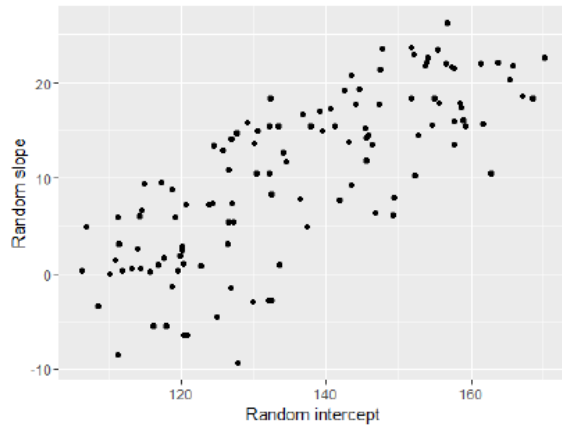
(b) Data set 3



# Components Fitted	DIC <sub>1</sub>	DIC <sub>2</sub>	DIC <sub>3</sub>	DIC <sub>4</sub>	DIC <sub>5</sub>	DIC <sub>6</sub>
1	2013	2012	2012	2611	2118	1570
2	1989	1949	1987	2497	2013	1562
3	1942	1942	1940	2339	2039	1571

DIC4 is most discerning.

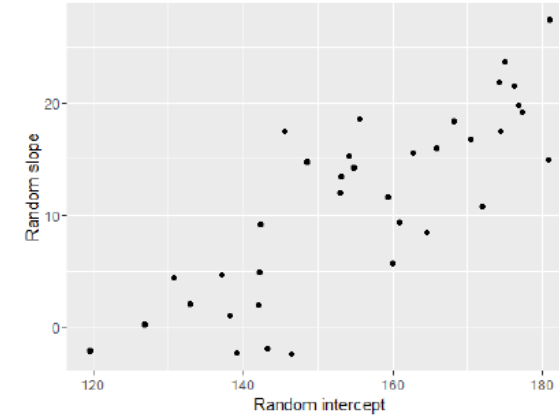
# DIC results



(a) Data set 4



# Components Fitted	DIC <sub>1</sub>	DIC <sub>2</sub>	DIC <sub>3</sub>	DIC <sub>4</sub>	DIC <sub>5</sub>	DIC <sub>6</sub>
1	6568	6566	6566	8454	6899	5197
2	6531	6523	6530	8263	6946	5253
3	6497	6492	6497	8017	6898	5263



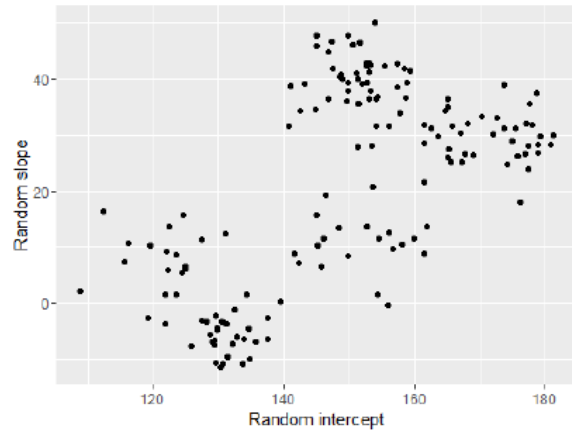
(b) Data set 5



# Components Fitted	DIC <sub>1</sub>	DIC <sub>2</sub>	DIC <sub>3</sub>	DIC <sub>4</sub>	DIC <sub>5</sub>	DIC <sub>6</sub>
1	1944	1943	1943	2500	1879	1364
2	1934	1936	1936	2438	2042	1526
3	1921	1920	1922	2350	2023	1537

Less sample size, fused components → DIC is not much discerning for anything other than DIC4

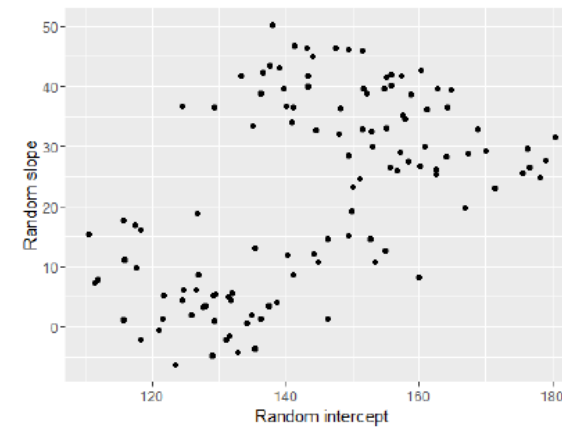
# DIC results



(a) Data set 6



# Components Fitted	DIC <sub>1</sub>	DIC <sub>2</sub>	DIC <sub>3</sub>	DIC <sub>4</sub>	DIC <sub>5</sub>	DIC <sub>6</sub>
1	8982	8981	8980	11847	9251	6655
2	8829	8827	8827	11327	9293	6838
3	8745	8742	8744	11036	9251	6895
4	8669	8672	8677	10737	9208	6925
5	8649	8643	8648	10601	9165	6909



(b) Data set 7



# Components Fitted	DIC <sub>1</sub>	DIC <sub>2</sub>	DIC <sub>3</sub>	DIC <sub>4</sub>	DIC <sub>5</sub>	DIC <sub>6</sub>
1	6708	6707	6706	8819	6977	5071
2	6606	6605	6604	8443	6946	5135
3	6539	6538	6537	8178	6944	5204
4	6506	6514	6521	8078	6915	5196
5	6505	6500	6508	7984	6896	5202

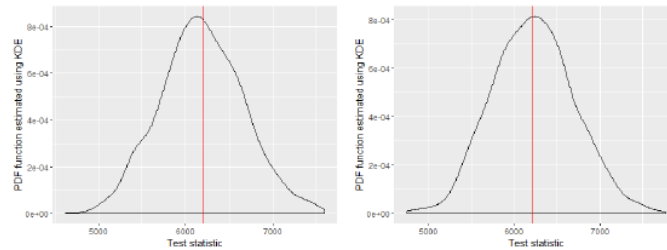
# Marginal likelihood results ( $\log \hat{m}(y)$ )

---

	Fitted	1 Comp	2 Comp	3 Comp	4 Comp	5 Comp
	Data set 1	-2120				
✓	Data set 2	-5019	-4989	-4937		
✗	Data set 3	-1038	-1044	-1042		
✗	Data set 4	-3317	-3318	-3322		
✗	Data set 5	-1001	-986	-993		
✗	Data set 6	-4545	-4492	-4477	-4467	-4473
✓	Data set 7	-3397	-3379	-3373	-3380	-2749

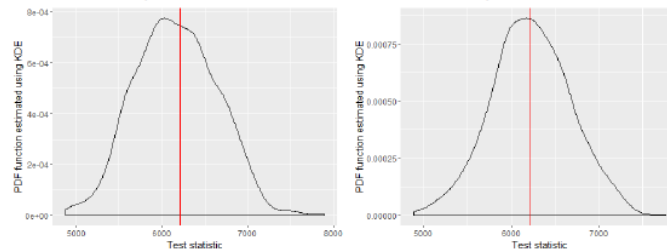
# PPC results (data set 6)

## Using Wishart prior for covariance matrix



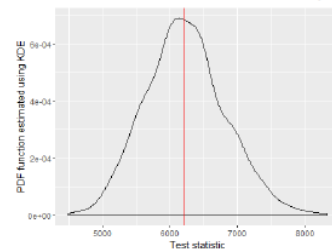
(a) # Components fitted = 5

(b) # Components fitted = 4

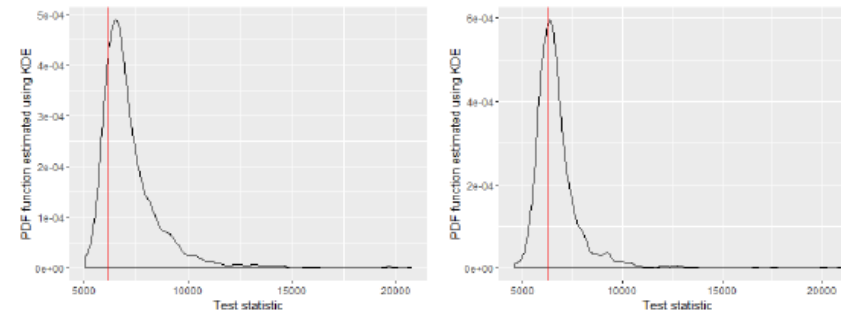


(c) # Components fitted = 3

(d) # Components fitted = 2

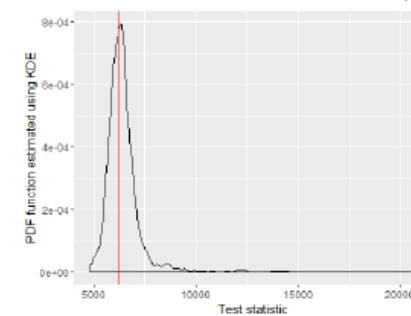


(e) # Components fitted = 1



(a) # Components fitted = 8

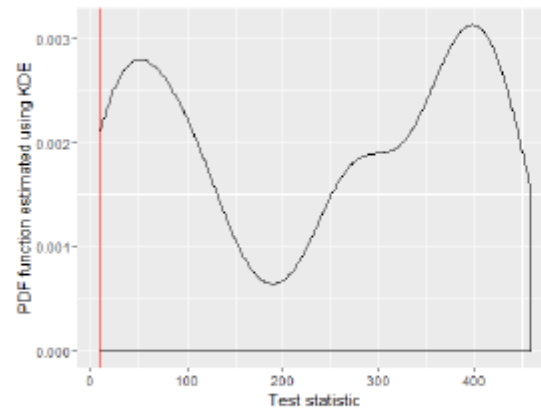
(b) # Components fitted = 7



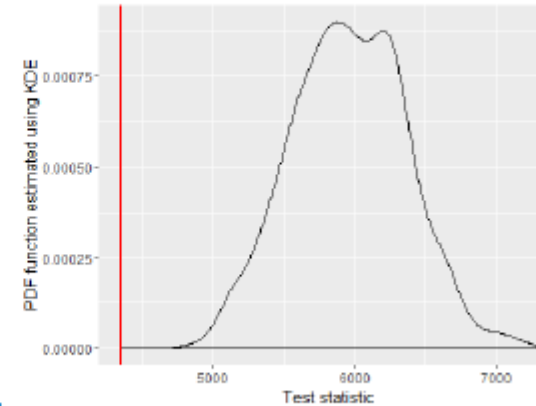
(c) # Components fitted = 6

# PPC results

Using  $U(-1,1)$  prior on  $\rho$  Gamma  $(10^{-4}, 10^{-4})$  on precision



(a) 4 components fitted for data set 2. log scale is used.

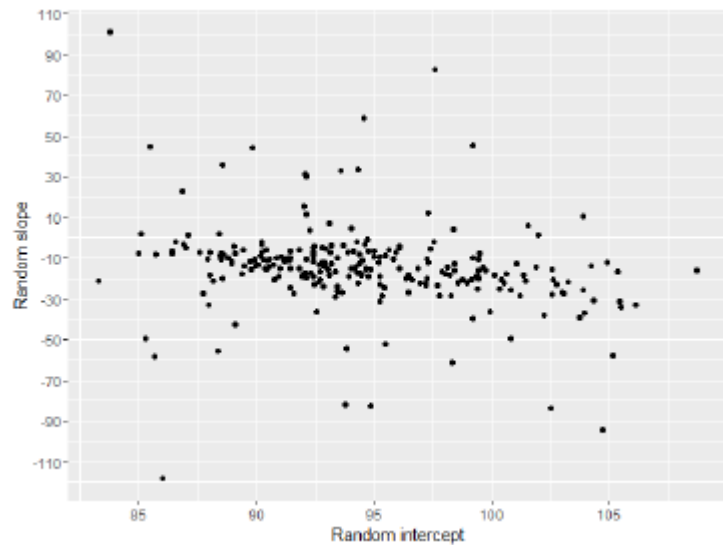


(b) 3 components fitted for data set 2

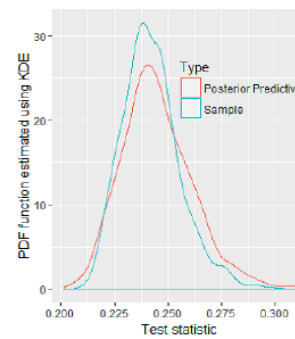
- a) Overfitting has a big penalty
- b) posterior variances of random effects were underestimating the sample data's variance covariance of the random effects. Bad model fit is detected with the test statistic



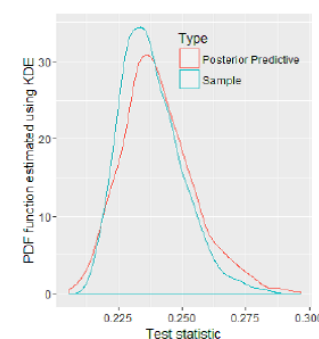
# Blood donor data set analysis



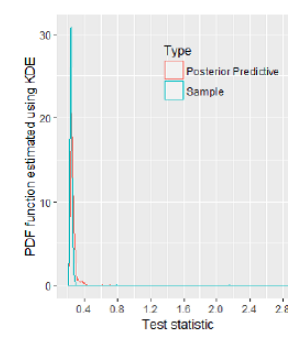
# Components Fitted	DIC <sub>1</sub>	DIC <sub>2</sub>	DIC <sub>3</sub>	DIC <sub>4</sub>	DIC <sub>5</sub>	DIC <sub>6</sub>
1 comp	4817	4816	4818	7077	7532	4353
2 comp	4808	4805	4811	6956	7376	4306



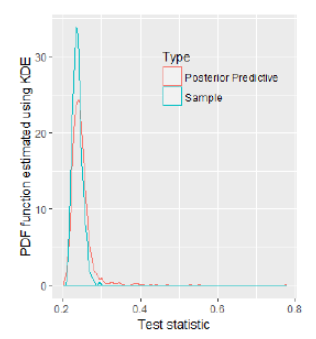
(a) #Components fitted = 1



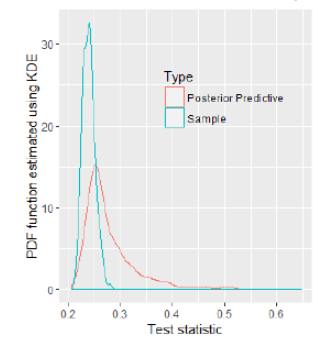
(b) #Components fitted = 2



(a) #Components fitted = 3



(b) #Components fitted = 4



(c) #Components fitted = 5

# Blood donor data set analysis

---

## Component 1

- Baseline random component of Hb level is low
- Random slope is less negative relative to group 2.
- Correlation between random slope and random intercept is most likely 0.

## Component 2

- Baseline random component of Hb level is higher.
- Average subject from group 2 has a more rapid decrease in Hb level levels with frequent donations.
- Higher the baseline Hb level level, more likely to have a rapid decrease in Hb level levels with donations

# Other interesting results

---

Dirichlet prior with large parameter values. Incrementally increasing hyperparameter?

Chib's approximation and chain length. Chib's approximation when posterior's do not belong to well known parametric families.

Wishart prior: underestimated posterior component variances if actual within subject variability > between subject variability.