



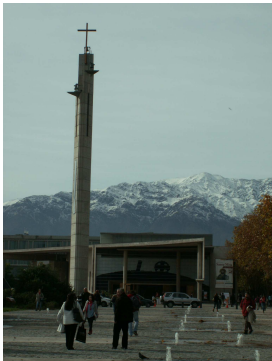
Arnošt Komárek

Katedra pravděpodobnosti a matematické statistiky

**Regrese s korelovanými
intervalově cenzorovanými daty
zatíženými nepřesnou klasifikací události**

ROBUST'18

Jetřichovice, 19. – 24. ledna 2014



Příspěvek založen na spolupráci

s **Maríou José García-Zattera**
a **Alejandro Jarou**

Pontificia Universidad Católica de Chile
Santiago de Chile



Část I

**Intervalově cenzorovaná data zatížená
nepřesnou klasifikací události**

Signal Tandmobiel® (*mobilní zub*) studie

Longitudinální zubní studie ve Flandrech, 1996 – 2001.

Stratifikovaný náhodný výběr dětí navštěvujících v roce 1996 1. třídu ZŠ.

2 315 chlapců, 2 153 dívek (7.3 % vlámské populace daného věku).

Každoroční (v 1. – 6. třídě) detailní zubní vyšetření jedním z 16 proškolených zubařů.

- Dané dítě neprohlíží v jednotlivých letech nutně stejný zubař.

V prvním roce též dotazník ohledně zubní hygieny, stravovacích návyků atd. vyplněný rodiči.

Korelované časy do události

Pro pamětníky: ROBUST 2006 (Lhota nad Rohanovem), ROBUST 2004 (Třešť)

$T_{(i,j)} \in \mathbb{R}_+$: čas do vzniku zubního kazu u j tého zuby i tého dítěte
($i = 1, \dots, N, j = 1, \dots, J$)

- obecně: čas do události u j té jednotky i tého subjektu.

$\mathbf{x}_{(i,j)}$: vektor regresorů, které mohou vysvětlit chování $T_{(i,j)}$.

Hlavní cíl: Regresní model pro závislost $T_{(i,j)}$ na $\mathbf{x}_{(i,j)}$.

Potřeba vzít v potaz:

složky vektoru $\mathbf{T}_i = (T_{(i,1)}, \dots, T_{(i,J)})^\top$ nejsou nekorelované.

Kontrola, zda nastala událost, je prováděna pouze v (předem) daných momentech.

K_i : počet kontrol u i tého subjektu.

Časy kontrol u i tého subjektu:

$$0 = v_{(i,0)} < v_{(i,1)} < v_{(i,2)} < \dots < v_{(i,K_i)} < v_{(i,K_i+1)} = \infty.$$

Čas $T_{(i,j)}$ není pozorován přesně, je pouze známo

$$T_{(i,j)} \in (v_{(i,l_{(i,j)})-1}, v_{(i,l_{(i,j)})}] \quad \text{pro } l_{(i,j)} \in \{1, \dots, K_i + 1\}$$

▮▮▮ **intervalově cenzorovaná data**

- $l_{(i,j)} \leq K_i$: $v_{(i,l_{(i,j)})}$ je čas kontroly, při které byla událost poprvé detekována;
- $l_{(i,j)} = K_i + 1$: $v_{(i,l_{(i,j)})} = \infty$ ▮▮▮ událost nenastala do času poslední kontroly.

Intervalové cenzorování zatížené nepřesnou klasifikací události

Misclassified interval-censored data

Klasifikace události (zda nastala či nikoliv) může být zatíženo chybou

- examinátor (diagnostický test, ...) s nenulovou pravděpodobností falešně pozitivního, resp. falešně negativního výsledku.

Pro (i, j) tou jednotku pozorujeme 0/1 posloupnost

$$\mathbf{Y}_{(i,j)} = (Y_{(i,j,1)}, \dots, Y_{(i,j,K_i)})^T.$$

Examinátor (test) indikoval, že událost do k té kontroly (v čase $v_{(i,k)}$)

- nastala $\Rightarrow Y_{(i,j,k)} = 1$;
- nenastala $\Rightarrow Y_{(i,j,k)} = 0$.

\Rightarrow intervalové cenzorování zatížené nepřesnou klasifikací události
(*misclassified interval-censored data*)

Intervalové cenzorování zatížené nepřesnou klasifikací události

Misclassified interval-censored data

S perfektním examínátorem/testem (nulová pravděpodobnost falešně pozitivních i negativních výsledků) by posloupnost $Y_{(i,j)}$ byla (samozřejmě) monotónní.

Chování examínátora/testu v čase $v_{i,k}$ nezávislé na výsledcích získaných v minulosti ($Y_{(i,j,l)}, l < k$):

- posloupnost $Y_{(i,j)}$ **není nutně monotónní**.

Regrese s intervalové cenzorovanými daty zatíženými nepřesnou klasifikací události

Hlavní cíl:

i nadále regresní model závislosti času do události $T_{(i,j)}$ na prediktorech $\mathbf{x}_{(i,j)}$.

Čas $T_{(i,j)}$ však pozorujeme pouze skrze $\mathbf{Y}_{(i,j)}$
– sadu 0/1 indikátorů události zatížených klasifikační chybou.

Pro itý subjekt (dítě), $i = 1, \dots, N$

$\mathbf{T}_i = (T_{(i,1)}, \dots, T_{(i,J)})^\top$: časy do události pro J jednotek (zuby)

▮▮▮ „skrytá“ data.

$\mathbf{Y}_i = (\mathbf{Y}_{(i,1)}^\top, \dots, \mathbf{Y}_{(i,J)}^\top)^\top$: posloupnosti 0/1 indikátorů události zatíženými klasifikační chybou

▮▮▮ pozorovaná data.

$\mathbf{x}_i = (\mathbf{x}_{(i,1)}^\top, \dots, \mathbf{x}_{(i,J)}^\top)^\top$: regresory pro vysvětlení chování časů \mathbf{T}_i .

$\mathbf{v}_i = (v_{(i,1)}, \dots, v_{(i,K_i)})^\top$: časy kontrol stavu události.

Model pro pozorovaná data (**věrohodnost**) specifikujeme hierarchicky.

Sdružená věrohodnost skrytých a pozorovaných dat *i*tého subjektu

$$p(\mathbf{Y}_i, \mathbf{T}_i) = p(\mathbf{Y}_i | \mathbf{T}_i) p(\mathbf{T}_i).$$

- $p(\mathbf{Y}_i | \mathbf{T}_i)$: model pro klasifikační proces
 - časy \mathbf{v}_i kontrol stavu události zde budou vystupovat v roli regresorů.
- $p(\mathbf{T}_i)$: (regresní) model pro (korelované) časy do události
 - \mathbf{x}_i zde bude vystupovat klasicky v roli regresorů.

Věrohodnost pozorovaných dat *i*tého subjektu

$$p(\mathbf{Y}_i) = \int_{\mathbb{R}_+^J} p(\mathbf{Y}_i, \mathbf{T}_i) d\mathbf{T}_i.$$

Budeme předpokládat nezávislost mezi subjekty, tj. pro věrohodnosti odpovídající všem subjektům:

$$p(\mathbf{Y}_1, \dots, \mathbf{Y}_N, \mathbf{T}_1, \dots, \mathbf{T}_N) = \prod_{i=1}^N p(\mathbf{Y}_i, \mathbf{T}_i) = \prod_{i=1}^N \left\{ p(\mathbf{Y}_i | \mathbf{T}_i) p(\mathbf{T}_i) \right\},$$

$$p(\mathbf{Y}_1, \dots, \mathbf{Y}_N) = \prod_{i=1}^N p(\mathbf{Y}_i) = \prod_{i=1}^N \int_{\mathbb{R}_+^J} p(\mathbf{Y}_i, \mathbf{T}_i) d\mathbf{T}_i.$$

Část II

Model pro klasifikační proces

Nyní se budeme zabývat $p(\mathbf{Y}_i | \mathbf{T}_i)$ částí sdružené věrohodnosti skrytých a pozorovaných dat.

Budeme předpokládat, že klasifikace události pro danou jednotku v daném čase ($Y_{(i,j,k)}$) je (podmíněně) **nezávislá** na

- (a) klasifikaci události u jiných jednotek (jiné j) v libovolném čase (libovolné k);
- (b) klasifikaci události u stejné jednotky (stejně j) v jiném čase (jiné k);
- (c) času události u jiných jednotek (jiné j).

To jest, předpokládáme, že lze psát

$$p(\mathbf{Y}_i | \mathbf{T}_i) = \prod_{j=1}^J \prod_{k=1}^{K_j} p(Y_{(i,j,k)} | T_{(i,j)}).$$

V dalším už se budeme zabývat jenom tím, jak může vypadat $p(Y_{(i,j,k)} | T_{(i,j)})$.

Jeden examinator/test

α : **sensitivita** examinatora/testu.

η : **specificita** examinatora/testu.

To jest,

$$\alpha = P(Y_{(i,j,k)} = 1 \mid T_{(i,j)} \leq v_{(i,k)}),$$

$$\eta = P(Y_{(i,j,k)} = 0 \mid T_{(i,j)} > v_{(i,k)}).$$

α, η : neznámé parametry modelu.

$Y_{(i,j,k)} \mid T_{(i,j)}$ část modelu:

$$p(Y_{(i,j,k)} \mid T_{(i,j)}) = p(Y_{(i,j,k)} \mid T_{(i,j)}; \alpha, \eta, v_{i,k})$$

$$= \begin{cases} \alpha^{Y_{(i,j,k)}} (1 - \alpha)^{1 - Y_{(i,j,k)}}, & \text{pokud } T_{(i,j)} \leq v_{(i,k)} \\ & \text{(správné } Y_{(i,j,k)} \text{ je rovno 1),} \\ (1 - \eta)^{Y_{(i,j,k)}} \eta^{1 - Y_{(i,j,k)}}, & \text{pokud } T_{(i,j)} > v_{(i,k)} \\ & \text{(správné } Y_{(i,j,k)} \text{ je rovno 0).} \end{cases}$$

Q **examinátorů/testů, různé sensitivity/specificity pro různé jednotky (různá j)**

Bylo potřeba v kontextu studie mobilního zubu.

Klasifikaci kazu provádělo $Q = 16$ různých zubařů:

- každý má jinou schopnost detekovat zubní kaz;
- potřeba uvážit závislost sensitivit/specificit na examinatorovi.

Různá $j \equiv$ různé zuby (šestka vlevo dole, čtyřka vpravo nahoře, ...):

- detekce kazu je na různých zubech různě obtížná;
- potřeba uvážit závislost sensitivit/specificit na j (jednotce).

Q **examinátorů/testů, různé sensitivity/specificity pro různé jednotky (různá j)**

Další „regresor“ v modelu: $\xi_{(i,k)} \in \{1, \dots, Q\}$

- index (číslo) examinatora, který prováděl klasifikaci události u (všech jednotek) i tého subjektu během jeho $kté$ kontroly v čase $v_{(i,k)}$.

Trochu více neznámých parametrů modelu ($q = 1, \dots, Q$):

$$\alpha_q = (\alpha_{(q,1)}, \dots, \alpha_{(q,J)})^\top,$$

$$\eta_q = (\eta_{(q,1)}, \dots, \eta_{(q,J)})^\top.$$

$\alpha_{(q,j)}, \eta_{(q,j)}$: sensitivita a specificita klasifikace, jestliže ji provádí examinator q na jednotce j , to jest,

$$\alpha_{(q,j)} = P(Y_{(i,j,k)} = 1 \mid T_{(i,j)} \leq v_{(i,k)}; \xi_{(i,k)} = q),$$

$$\eta_{(q,j)} = P(Y_{(i,j,k)} = 0 \mid T_{(i,j)} > v_{(i,k)}; \xi_{(i,k)} = q).$$

Q **examinátorů/testů, různé sensitivity/specificity pro různé jednotky (různá j)**

Vektory všech (neznámých) sensitivit a specificit:

$$\alpha = (\alpha_1^\top, \dots, \alpha_Q^\top)^\top,$$

$$\eta = (\eta_1^\top, \dots, \eta_Q^\top)^\top.$$

$Y_{(i,j,k)} \mid T_{(i,j)}$ část modelu, jenom nějaké indexy navíc ☺

$$p(Y_{(i,j,k)} \mid T_{(i,j)}) = p(Y_{(i,j,k)} \mid T_{(i,j)}; \alpha, \eta, v_{i,k}, \xi_{(i,k)})$$

$$= \begin{cases} \alpha_{(\xi_{(i,k)}, j)}^{Y_{(i,j,k)}} (1 - \alpha_{(\xi_{(i,k)}, j)})^{1 - Y_{(i,j,k)}}, & \text{pokud } T_{(i,j)} \leq v_{(i,k)} \\ & \text{(správné } Y_{(i,j,k)} \text{ je rovno 1),} \\ (1 - \eta_{(\xi_{(i,k)}, j)})^{Y_{(i,j,k)}} \eta_{(\xi_{(i,k)}, j)}^{1 - Y_{(i,j,k)}}, & \text{pokud } T_{(i,j)} > v_{(i,k)} \\ & \text{(správné } Y_{(i,j,k)} \text{ je rovno 0).} \end{cases}$$

Sensitivity/specificity lze dále modelovat pomocí charakteristik jednotlivých examinátorů/testů, resp. charakteristik jednotek, jsou-li nějaké charakteristiky k dispozici.

Lze použít např. logistickou regresi

- další hierarchická úroveň celkového modelu;
- detaily až po 22. hodině.

Tím máme z krku specifikaci $p(\mathbf{Y}_i | \mathbf{T}_i)$ části modelu.

Zbývá ještě $p(\mathbf{T}_i)$ část (model pro časy do události).

Část III

Model pro čas do události

Připomínka

$$\mathbf{T}_i = (T_{(i,1)}, \dots, T_{(i,J)})^\top$$

≡ ne nutně nekorelované časy událostí u J jednotek i tého subjektu.

$$\mathbf{x}_i = (\mathbf{x}_{(i,1)}^\top, \dots, \mathbf{x}_{(i,J)}^\top)^\top$$

≡ možné regresory vysvětlující chování časů do události.

Tvar $p(\mathbf{T}_i)$ lze principiálně odvodit z libovolného regresního modelu pro korelovaná „přeživací“ data (o kterém jsme přesvědčeni, že se podle něho časy do události skutečně řídí):

- *frailty* Coxův model;
- *accelerated failure time (AFT)* model s náhodným absolutním členem;
- \vdots

AFT model s náhodným absolutním členem

Pro pamětníky: ROBUST 2006 (Lhota nad Rohanovem), ROBUST 2004 (Třešť)

$$\log(T_{(i,j)}) = \mathbf{x}_{(i,j)}^T \boldsymbol{\beta} + b_i + \varepsilon_{(i,j)} \quad i = 1, \dots, N, j = 1, \dots, J,$$

- $\boldsymbol{\beta}$: neznámé regresní koeficienty;
- $\varepsilon_{(1,1)}, \dots, \varepsilon_{(N,J)}$: i.i.d. náhodné veličiny s hustotou $g_\varepsilon(\cdot)$ s nulovou střední hodnotou;
- b_1, \dots, b_N : i.i.d. náhodné veličiny s hustotou $g_b(\cdot)$
 - b_i (společné pro všechna j) indukuje jistou formu závislosti mezi $T_{(i,1)}, \dots, T_{(i,J)}$;
- $\varepsilon_{(1,1)}, \dots, \varepsilon_{(N,J)}, b_1, \dots, b_N$ nezávislé.

Příslušná transformace hustot g_ε a g_b určuje rozdělení jednotlivých časů do události a naše $p(\mathbf{T}_i)$.

Zde:

$$g_\varepsilon(\cdot) \sim \mathcal{N}(0, \sigma_\varepsilon^2),$$

$$g_b(\cdot) \sim \mu + \tau \underbrace{\sum_{l=-M}^M w_l \mathcal{N}(\kappa_l, \zeta^2)}_{\text{penalizovaná normální směs}}$$

- Neznámé parametry: σ_ε^2 , $\mathbf{w} = (w_{-M}, \dots, w_M)^\top$, μ , τ .

- Penalizovaná normální směs:

$$M \approx 15, \zeta \approx 0.2,$$

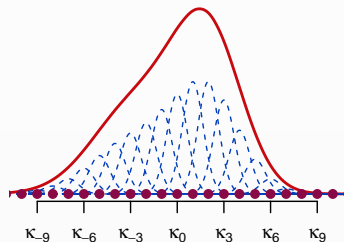
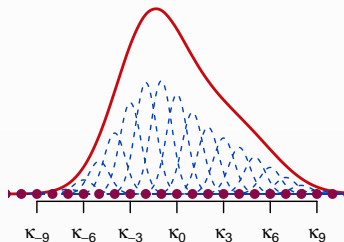
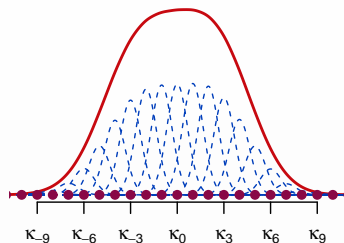
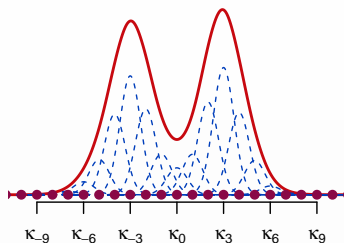
$\kappa_{-M}, \dots, \kappa_M$: ekvidistantní uzly na intervalu přibližně $[-4.5, 4.5]$;

▮ flexibilní model pro rozdělení s přibližně nulovou střední hodnotou a jednotkovým rozptylem.

- Regularizace pomocí penalizace diferencí (transformovaných) vah w_{-M}, \dots, w_M .

Penalizovaná normální směs

Pro pamětníky: ROBUST 2006 (Lhota nad Rohanovem), ROBUST 2004 (Třešť)



Tvar $p(\mathbf{T}_i)$ odvozen z následujícího

- $\log(T_{(i,j)}) = \mathbf{x}_{(i,j)}^\top \boldsymbol{\beta} + b_i + \varepsilon_{(i,j)}, \quad i = 1, \dots, N, j = 1, \dots, J;$
- $\varepsilon_{(1,1)}, \dots, \varepsilon_{(N,J)}, b_1, \dots, b_N$ nezávislé;
- $\varepsilon_{(i,j)} \sim \mathcal{N}(0, \sigma_\varepsilon^2);$
- $b_i \sim \mu + \tau \sum_{l=-M}^M w_l \mathcal{N}(\kappa_l, \zeta^2).$

Neznámé parametry modelu

$$\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \mathbf{w}^\top, \mu, \tau^2, \sigma_\varepsilon^2)^\top.$$

Rozdělení času do události $T_{(i,j)}$ je (po log-transformaci) konvolucí plně parametrického rozdělení (\mathcal{N}) a penalizované normální směsi.

- Penalizovaná normální směs v rozdělení náhodného absolutního členu AFT modelu zajišťuje kromě jiného flexibilní model pro rozdělení času do události $T_{(i,j)}$.

Část IV

Inference

Věrohodnost

$$\begin{aligned} p(\mathbf{Y}_1, \dots, \mathbf{Y}_N) &= \prod_{i=1}^N p(\mathbf{Y}_i) = \prod_{i=1}^N \int_{\mathbb{R}_+^J} p(\mathbf{Y}_i, \mathbf{T}_i) d\mathbf{T}_i \\ &= \prod_{i=1}^N \int_{\mathbb{R}_+^J} p(\mathbf{Y}_i | \mathbf{T}_i) p(\mathbf{T}_i) d\mathbf{T}_i. \end{aligned}$$

$p(\mathbf{Y}_i | \mathbf{T}_i)$: model pro klasifikační proces

- neznámé parametry: $\boldsymbol{\alpha} = (\alpha_{(1,1)}, \dots, \alpha_{(Q,J)})^\top$, $\boldsymbol{\eta} = (\eta_{(1,1)}, \dots, \eta_{(Q,J)})^\top$: sensitivity a specificity pro jednotlivé examinátory a jednotky.

$p(\mathbf{T}_i)$: model pro časy do události

- AFT model s náhodným absolutním členem a penalizovanou normální směsí v jeho rozdělení;
- neznámé parametry: $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \mathbf{w}^\top, \mu, \tau^2, \sigma_\varepsilon^2)^\top$.

Všechny parametry modelu jsou identifikovatelné, omezíme-li parametrický prostor pro sensitivity a specificity podmínkou

$$\alpha_{(q,j)} + \eta_{(q,j)} > 1, \quad q = 1, \dots, Q, j = 1, \dots, J.$$

Výpočet věrohodnosti, respektive její maximalizace (pokud bychom chtěli maximálně věrohodné odhady) je mírně (spíše více) zkomplikována integrálem ve vyjádření $p(\mathbf{Y}_1, \dots, \mathbf{Y}_N)$.

Bayesovská inference se slabě informativními apriorními rozděleními je proveditelná s pomocí MCMC:

- detaily opět až po 22. hodině;
- softwarově implementováno v  balíčku `bayesSurv` (od verze 2.3, která zatím není na CRANu, ale časem snad bude).

Část V

Simulační studie

$J = 4$, $N = 500, 1\,000, 2\,000$.

$$\log(T_{(i,j)}) = 2.0 + 0.2 x_{(i,j),1} - 0.1 x_{(i,j),2} + b_i + \varepsilon_{(i,j)}.$$

$$x_{(i,j),1} \sim \mathcal{U}(0, 1), x_{(i,j),2} \sim \mathcal{A}(0.5).$$

$$\text{var}(b_i) + \text{var}(\varepsilon_{(i,j)}) = 0.1.$$

$$\sqrt{\frac{\text{var}(b_i)}{\text{var}(\varepsilon_{(i,j)})}} = \frac{\sigma_b}{\sigma_\varepsilon} = 0.5, 1, 2, 5.$$

g_b : (a) \mathcal{N} , (b) zřetelně bimodální dvousložková \mathcal{N} směs, (c) Gumbel.

$K_i = 10$ kontrol stavu události (v náhodných intervalech).

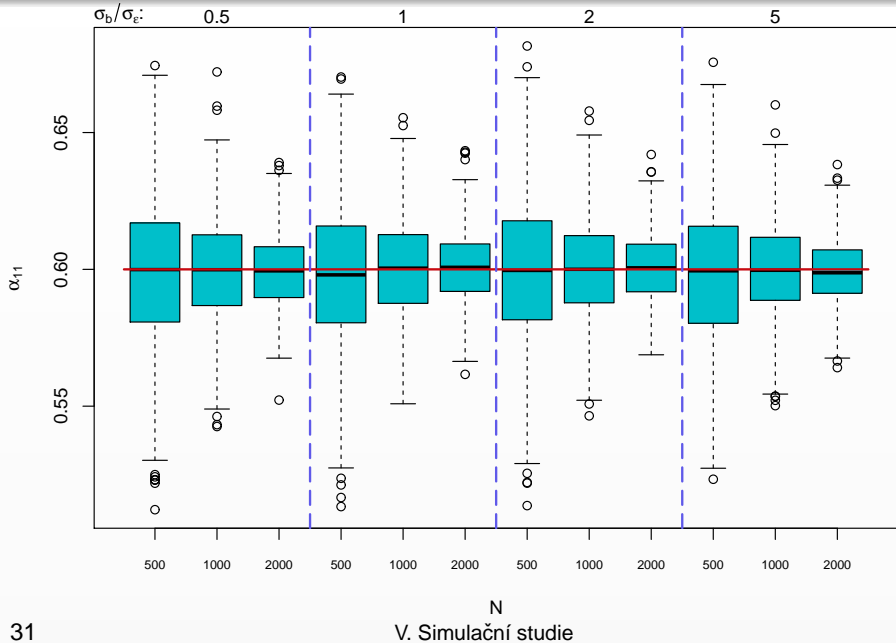
$Q = 5$ examinátorů náhodně přiřazovaných k jednotlivým kontrolám.

Sensitivity a specificity v rozmezí 0.60 – 0.96.

500 datových sad pro každý scénář.

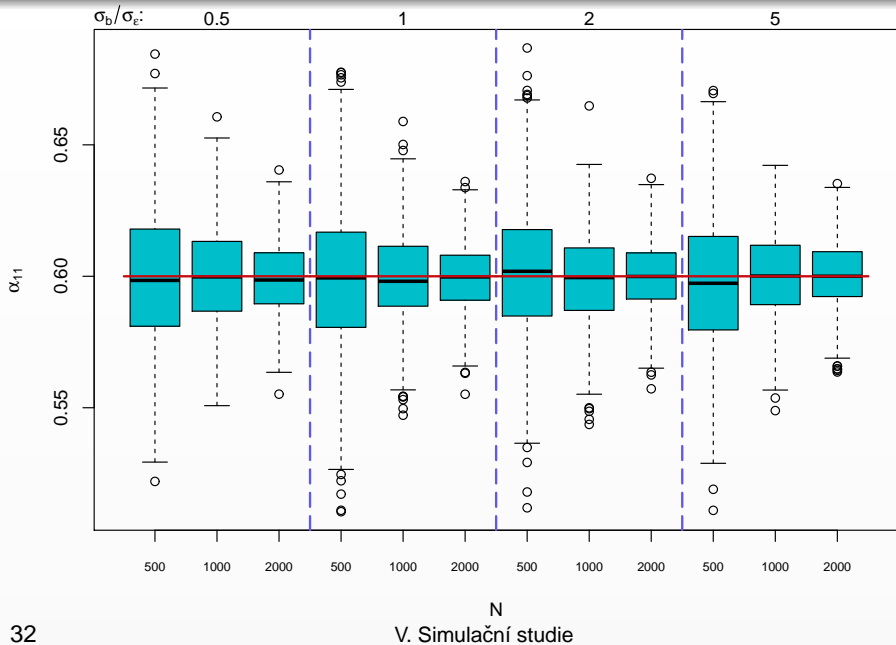
Sensitivita $\alpha_{(1,1)} = 0.60$

g_b : bimodální dvousložková \mathcal{N} směs



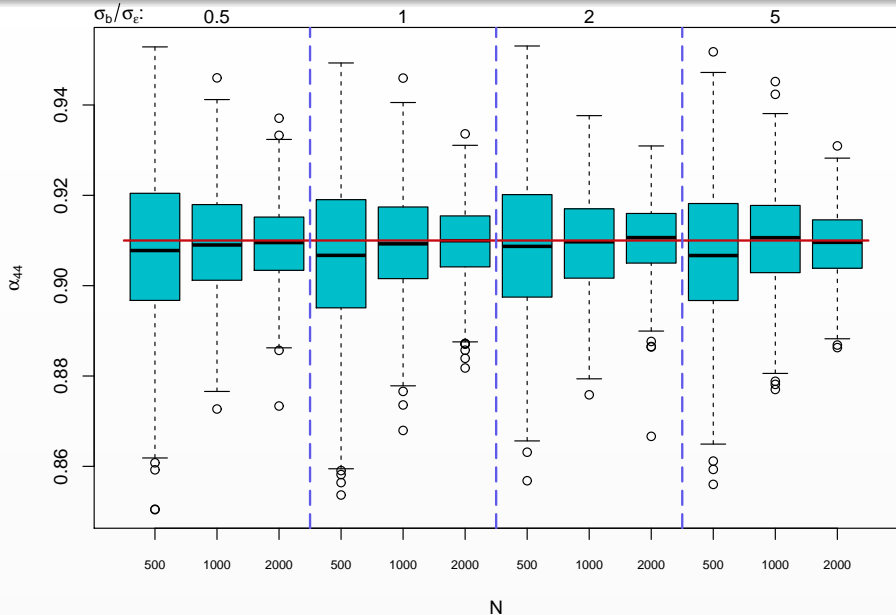
Sensitivita $\alpha_{(1,1)} = 0.60$

g_b : Gumbel



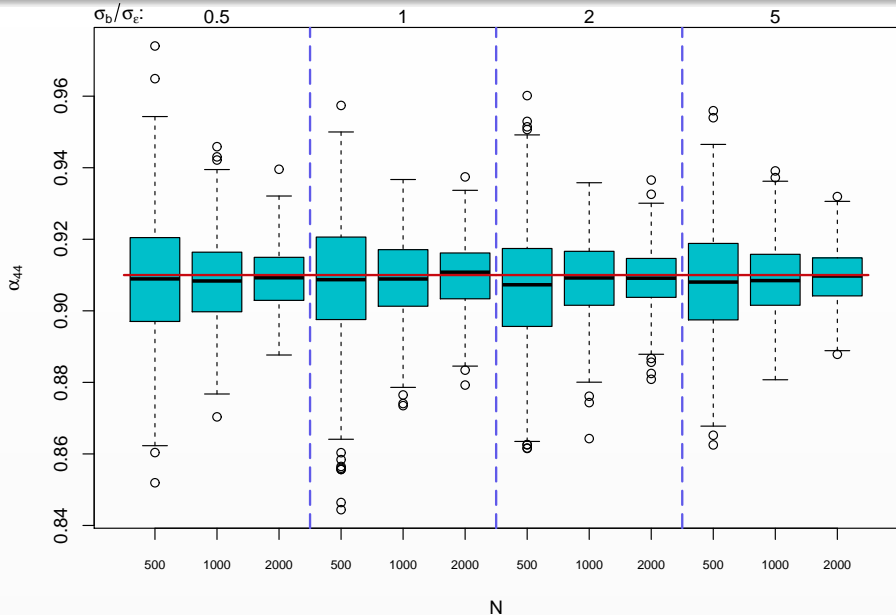
Sensitivita $\alpha_{(4,4)} = 0.91$

g_b : bimodální dvousložková \mathcal{N} směs



Sensitivita $\alpha_{(4,4)} = 0.91$

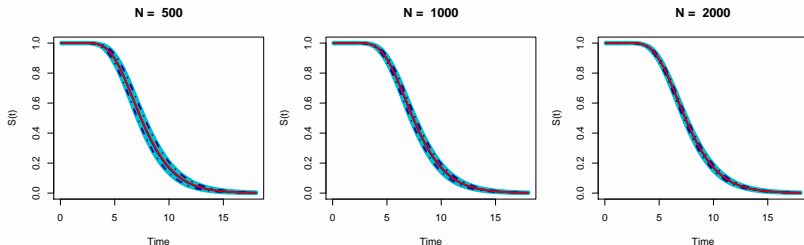
g_b : Gumbel



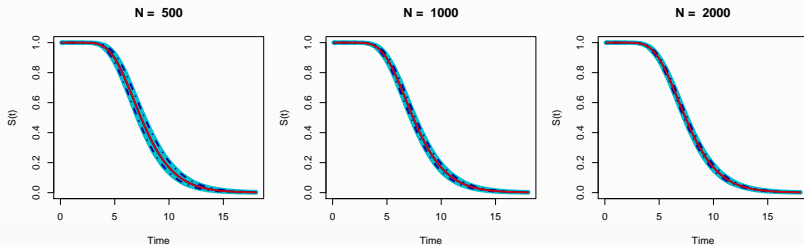
Funkce přežití pro $x_{(i,j),1} = 0.5$, $x_{(i,j),2} = 0$

$$\sigma_b / \sigma_\varepsilon = 0.5$$

g_b : bimodální dvousložková \mathcal{N} směs



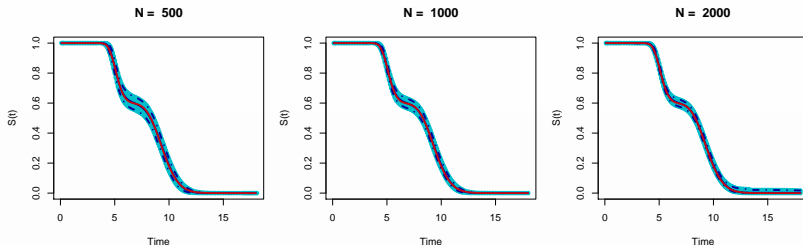
g_b : Gumbel



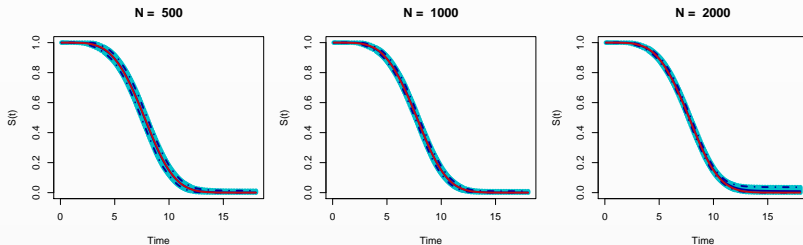
Funkce přežití pro $x_{(i,j),1} = 0.5$, $x_{(i,j),2} = 0$

$$\sigma_b / \sigma_\varepsilon = 5$$

g_b : bimodální dvousložková \mathcal{N} směs



g_b : Gumbel



DĚKUJI ZA POZORNOST.

29th International Workshop on Statistical Modelling

14. – 18. července/júla 2014

Göttingen

- dlouholeté působiště Carla Friedricha Gause,
- 388 km po silnici z Jetřichovic.

Termín pro zaslání příspěvků: až 31. ledna/januára 2014
(a stejně bude ještě o minimálně týden posunut).

Zvaní řečníci

- Antoine de Falguerolles (*Université de Toulouse*)
- Alejandro Jara (*Pontificia Universidad Católica de Chile*)
- Sophia Rabe-Hesketh (*University of California, Berkeley*)
- Gerhard Tutz (*Ludwig-Maximilians-Universität München*)
- Simon Wood (*University of Bath*)