

Odhady základního rizika v regresních modelech oprav

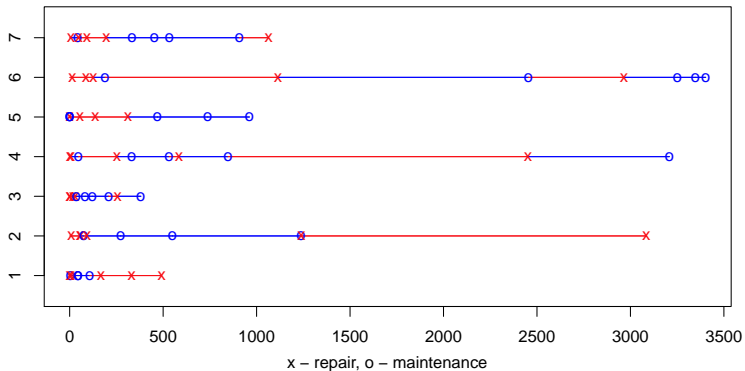
Petr Novák

KPMS MFF UK

21. ledna 2014

- Studujeme data o historii n zařízení, které podléhají opotřebení.
- Když se zařízení porouchá, je nutné provést opravu.
- Poruchám se snažíme předcházet preventivní údržbou.
- T_{i1}, \dots, T_{in_i} časy zásahů (oprav nebo údržeb).
- $\Delta_{i1}, \dots, \Delta_{in_i}$ indikátor, zda v j -tém čase byla provedena na i -tém zařízení oprava.
- $X_j(t)$ vysvětlující proměnné.

- Studujeme data o historii n zařízení, které podléhají opotřebení.
- Když se zařízení porouchá, je nutné provést opravu.
- Poruchám se snažíme předcházet preventivní údržbou.
- T_{i1}, \dots, T_{in_i} časy zásahů (oprav nebo údržeb).
- $\Delta_{i1}, \dots, \Delta_{in_i}$ indikátor, zda v j -tém čase byla provedena na i -tém zařízení oprava.
- $X_j(t)$ vysvětlující proměnné.
- Chceme rozumný popis vlivu oprav, údržby a dalších regresorů na životnost.



Data jako čítací procesy

- Zavedeme čítací procesy oprav a údržeb

$$N_{i\bullet}(t) = \sum_{j=1}^{n_i} I(T_{ij} \leq t, \Delta_{ij} = 1), \quad M_{i\bullet}(t) = \sum_{j=1}^{n_i} I(T_{ij} \leq t, \Delta_{ij} = 0).$$

- Označíme rizikovou funkcí

$$\lambda_i(t) = \lim_{h \rightarrow 0} P(N_{i\bullet}(t+h) - N_{i\bullet}(t) \geq 1 | \mathcal{H}(t)) / h$$

kde $\mathcal{H}(t)$ značí historii událostí do času t
a kumulativní rizikovou funkcí $\Lambda_i(t) = \int_0^t \lambda_i(s) ds$.

- Zavedeme čítací procesy oprav a údržeb

$$N_{i\bullet}(t) = \sum_{j=1}^{n_i} I(T_{ij} \leq t, \Delta_{ij} = 1), \quad M_{i\bullet}(t) = \sum_{j=1}^{n_i} I(T_{ij} \leq t, \Delta_{ij} = 0).$$

- Označíme rizikovou funkci

$$\lambda_i(t) = \lim_{h \rightarrow 0} P(N_{i\bullet}(t+h) - N_{i\bullet}(t) \geq 1 | \mathcal{H}(t)) / h$$

kde $\mathcal{H}(t)$ značí historii událostí do času t
a kumulativní rizikovou funkci $\Lambda_i(t) = \int_0^t \lambda_i(s) ds$.

- Předpokládejme, že oprava sice vrátí prvek jen do stavu těsně před poruchou, ale má vliv na rizikovou funkci.
- Rizikovou funkci vhodně parametrizujeme, chceme odhadnout parametry metodou max. věrohodnosti.

- Věrohodnost lze zapsat jako

$$\begin{aligned} L &= \prod_{i=1}^n \prod_{j=1}^{n_i} \left(\frac{f_i(T_{ij}^-)}{S_i(T_{i(j-1)})} \right)^{\Delta_{ij}} \left(\frac{S_i(T_{ij})}{S_i(T_{i(j-1)})} \right)^{1-\Delta_{ij}} \\ &= \prod_{i=1}^n \left(\prod_{j=1}^{n_i} \lambda_i(T_{ij}^-)^{\Delta_{ij}} \cdot S_i(T_{in_i}) \right) \end{aligned}$$

a log-věrohodnost má pak tvar

$$l = \sum_{i=1}^n \left(\sum_{j=1}^{n_i} \Delta_{ij} \log \lambda_i(T_{ij}^-) - \int_0^{T_{in_i}} \lambda_i(t) dt \right).$$

- V Coxově modelu působí regresory multiplikativně na rizikovou funkci.
- Předpokládáme, že každá oprava či údržba multiplikativně sníží nebo zvýší riziko, stejně tak případné regresory. Uvažujeme rizikovou funkci ve tvaru (Percy & Alkali 2005):

$$\lambda_i(t) = \lambda_0(t) e^{M_{i\bullet}(t)\rho + N_{i\bullet}(t)\sigma + X_i^T(t)\beta}.$$

- Pokud se hodnoty kovariáty mění jen v časech událostí, je možné snadno dosadit do logaritmické věrohodnosti a při parametrickém základním riziku maximalizovat.

- Můžeme také předpokládat, že každá oprava či údržba a regresory způsobí, že virtuální čas plyne pomaleji nebo rychleji (Accelerated Failure Time model, AFT). Využijeme transformaci času (Lin & Ying, 1995):

$$t \rightarrow \int_0^t e^{M_{i\bullet}(s)\rho + N_{i\bullet}(s)\sigma + X_i^T(s)\beta} ds =: h_i(t, \beta),$$

- kde jsme označili $\beta = (\rho, \sigma, \beta)^T$. Riziková funkce pak má tvar

$$\lambda_i(t) = \lambda_0(h_i(t, \beta)) e^{M_{i\bullet}(t)\rho + N_{i\bullet}(t)\sigma + X_i^T(t)\beta}.$$

- Pokud základní riziková funkce bude konstantní, oba modely splývají.

- Zavedeme čítací procesy pro j -té selhání či opravu i -tého zařízení a příslušný indikátor rizika

$$N_{ij}(t) = \Delta_{ij}I(T_{ij} \leq t),$$

$$M_{ij}(t) = (1 - \Delta_{ij})I(T_{ij} \leq t),$$

$$Y_{ij}(t) = I(T_{i,j-1} < t \leq T_{ij}).$$

Dostaneme

$$l = \sum_{i=1}^n \sum_{j=1}^{n_i} \int_0^{\infty} (\log \lambda_i(t^-) dN_{ij}(t) - Y_{ij}(t) \lambda_i(t^-) dt).$$

- Označíme $\mathbf{X}_i(t) = (N_{i\bullet}(t), M_{i\bullet}(t), X_i(t))^T$.
- Skóre získané dosazením rizikové funkce do logaritmické věrohodnosti a derivováním podle parametrů ale závisí na neznámé $\Lambda_0(t) = \int_0^t \lambda_0(s) ds$.
- Tu můžeme nahradit Nelson-Aalenovým odhadem

$$\hat{\Lambda}_0(t, \beta) = \int_0^t \frac{dN_{\bullet\bullet}(s)}{\sum_{ij} e^{\mathbf{x}_i^T(s^-)\beta} Y_{ij}(s)}.$$

- Po dosazení získáme skóre ve tvaru

$$U(\beta) = \sum_{ij} \int_0^\infty \left(\mathbf{x}_i(t^-) - \frac{\sum_{kl} \mathbf{x}_k(t^-) e^{\mathbf{x}_k^T(t^-)\beta} Y_{kl}(t)}{\sum_{ij} e^{\mathbf{x}_k^T(t^-)\beta} Y_{kl}(t)} \right) dN_{ij}(t)$$

a pro nalezení odhadů parametrů řešíme rovnice $U(\beta) = 0$.

- Pro každý prvek máme transformaci času $h_i(t, \beta)$. Zavedeme transformované procesy

$$N_{ij}^*(t, \beta) = \Delta_{ij} I(h_i(T_{ij}, \beta) \leq t), \quad M_{ij}^*(t, \beta) = (1 - \Delta_{ij}) I(h_i(T_{ij}, \beta) \leq t),$$

$$Y_{ij}^*(t, \beta) = I(h_i(T_{i,j-1}, \beta) < t \leq h_i(T_{ij}, \beta)), \quad X_i^*(t, \beta) = X_i(h_i^{-1}(t, \beta)).$$

- Přesné skóre má složitější tvar, je ale možné jej nahradit přibližným (Lin & Ying, 1995) a dosadit odhad

$$\hat{\Lambda}_0(t, \beta) = \int_0^t \frac{dN_{\bullet\bullet}^*(s, \beta)}{\sum_{ij} Y_{ij}^*(s, \beta)}.$$

- Získáme

$$\tilde{U}(\beta) = \sum_{ij} \int_0^\infty \left(\mathbf{X}_i^*(t^-, \beta) - \frac{\sum_{kl} \mathbf{X}_k^*(t^-, \beta) Y_{kl}^*(t, \beta)}{\sum_{kl} Y_{kl}^*(t, \beta)} \right) dN_{ij}^*(t, \beta)$$

a minimalizací $\|U(\beta)\|$ opět můžeme najít odhady parametrů.

Vlastnosti odhadů základního rizika

- Pomocí skóre získáme odhady $\hat{\beta}$.
- Chceme znát i vlastnosti $\hat{\Lambda}_0(t, \hat{\beta})$ a testovat hypotézy o skutečném tvaru.
- Uvažujme proces

$$W(t) = \sqrt{n}(\hat{\Lambda}_0(t, \hat{\beta}) - \Lambda_0(t)).$$

- $W(t)$ konverguje v obou modelech ke Gaussovskému procesu s nulovou střední hodnotou.
- Varianční funkce má složitý tvar - odhadneme pomocí resamplingu $W(t)$.

Resampling $W(t)$ v Coxově modelu

- Generujme G_1, \dots, G_n (*iid*) z $N(0, 1)$.
- Označme

$$U_G(\beta) = \sum_{ij} \int_0^\infty \left(\mathbf{x}_i(t^-) - \frac{\sum_{kl} \mathbf{x}_k(t^-) e^{\mathbf{x}_k^T(t^-)\beta} Y_{kl}(t)}{\sum_{ij} e^{\mathbf{x}_i^T(t^-)\beta} Y_{ij}(t)} \right) G_i dN_{ij}(t)$$

- Najdeme $\hat{\beta}_G$ jako řešení rovnice $U(\hat{\beta}_G) = U_G(\hat{\beta})$
- Označme $\hat{\Lambda}_{0G}(t, \beta) = \sum_{ij} \int_0^t \frac{G_i dN_{ij}(s)}{\sum_{kl} e^{\mathbf{x}_k^T(s^-)\beta} Y_{kl}(s)}$
- Potom

$$\hat{W}(t) = \sqrt{n} \left(\hat{\Lambda}_0(t, \hat{\beta}) - \hat{\Lambda}_0(t, \hat{\beta}_G) + \hat{\Lambda}_{0G}(t, \hat{\beta}) \right)$$

konverguje ke Gaussovskému procesu se stejnými vlastnostmi jako $W(t)$ v Coxově modelu.

Resampling $W(t)$ v AFT modelu

- Generujeme G_1, \dots, G_n (*iid*) z $N(0, 1)$.
- Označme

$$\tilde{U}_G(\beta) = \sum_{ij} \int_0^\infty \left(\mathbf{x}_i^*(t^-) - \frac{\sum_{kl} \mathbf{x}_k^*(t^-) Y_{kl}^*(t)}{\sum_{ij} Y_{kl}^*(t)} \right) G_i dN_{ij}^*(t)$$

- Najdeme $\hat{\beta}_G$ jako řešení rovnice $\tilde{U}(\hat{\beta}_G) = \tilde{U}_G(\hat{\beta})$
- Označme $\hat{\Lambda}_{0G}(t, \beta) = \sum_{ij} \int_0^t \frac{G_i dN_{ij}^*(s)}{\sum_{kl} Y_{kl}^*(s)}$
- Potom

$$\hat{W}(t) = \sqrt{n} \left(\hat{\Lambda}_0(t, \hat{\beta}) - \hat{\Lambda}_0(t, \hat{\beta}_G) + \hat{\Lambda}_{0G}(t, \hat{\beta}) \right)$$

konverguje ke Gaussovskému procesu se stejnými vlastnostmi jako $W(t)$ v AFT modelu.

- Replikujeme mnohokrát $\hat{W}(t)$.
- Empiricky odhadneme rozptyl $W(t)$.
- Spočítáme bodové konfidenční intervaly kumulovaného rizika:

$$\hat{\Lambda}_0(t, \hat{\beta}) \pm u_{1-\alpha/2} n^{-1/2} \sqrt{\widehat{\text{var}} \hat{W}(t)},$$

případně pomocí log-transformace

$$\hat{\Lambda}_0(t, \hat{\beta}) \exp \left(\pm u_{1-\alpha/2} n^{-1/2} \sqrt{\widehat{\text{var}} \hat{W}(t) / \hat{\Lambda}_0(t, \hat{\beta})} \right),$$

- kde $u_{1-\alpha/2}$ je příslušný kvantil $N(0, 1)$.

Konfidenční pás pro Λ_0

- Pro supremový test potřebujeme konfidenční pás.
- Najdeme $q_{1-\alpha}$ výběrový $1 - \alpha$ kvantil generovaných hodnot

$$\sup_{[\tau_1, \tau_2]} \left| \frac{\hat{W}(t)}{\widehat{\text{var}} \hat{W}(t)} \right|$$

kde $[\tau_1, \tau_2]$ pokrývá zkoumanou část časového intervalu.

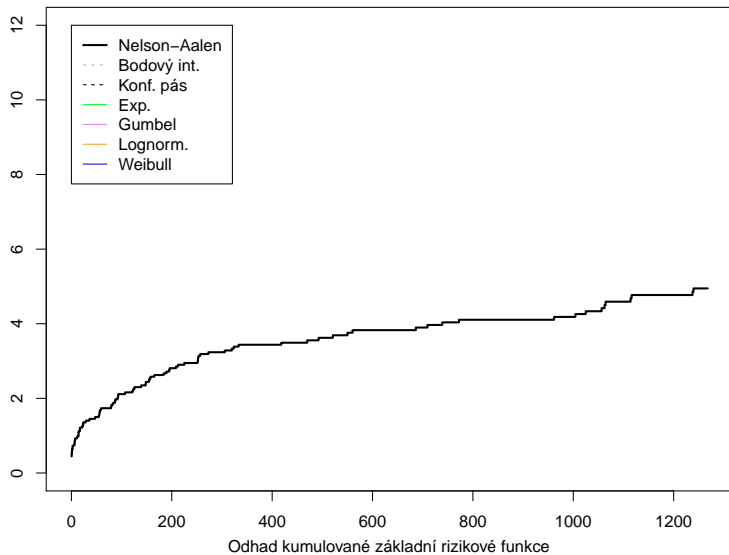
- Spočítáme konfidenční pás pro kumulované riziko:

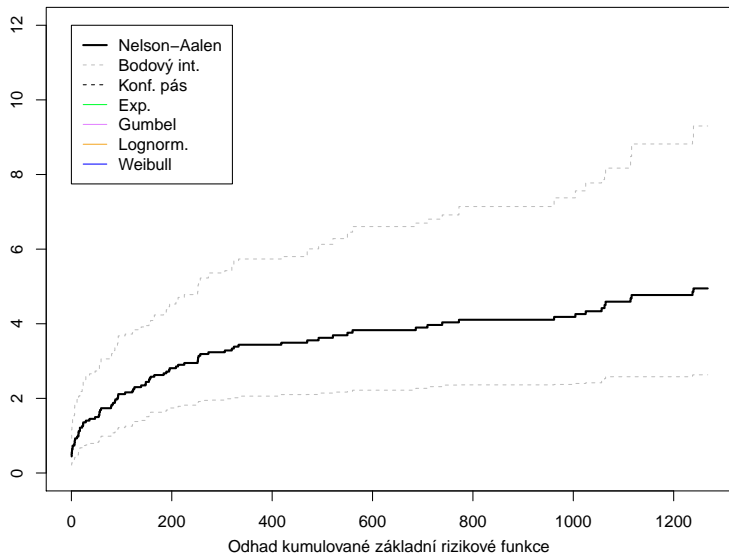
$$\hat{\Lambda}_0(t, \hat{\beta}) \pm q_{1-\alpha} n^{-1/2} \sqrt{\widehat{\text{var}} \hat{W}(t)},$$

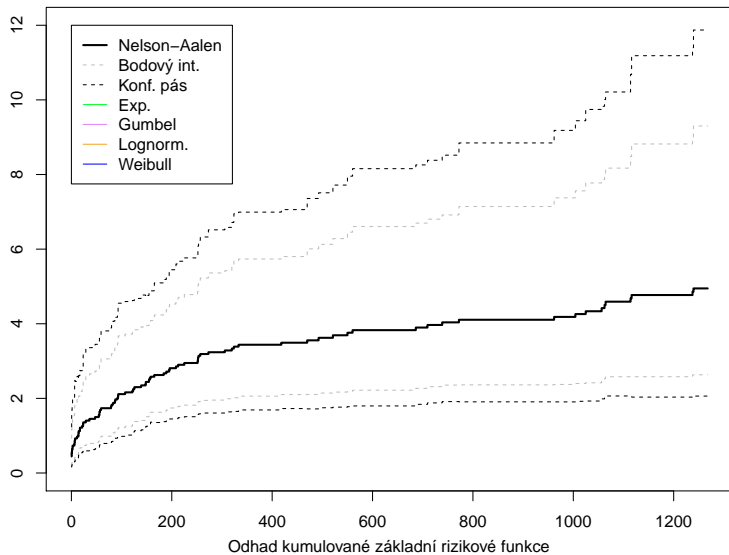
případně pomocí log-transformace

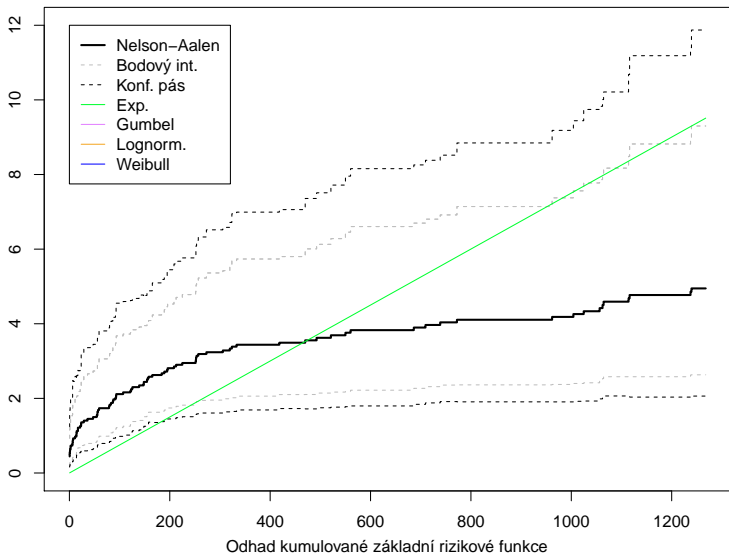
$$\hat{\Lambda}_0(t, \hat{\beta}) \exp \left(\pm q_{1-\alpha} n^{-1/2} \sqrt{\widehat{\text{var}} \hat{W}(t) / \hat{\Lambda}_0(t, \hat{\beta})} \right).$$

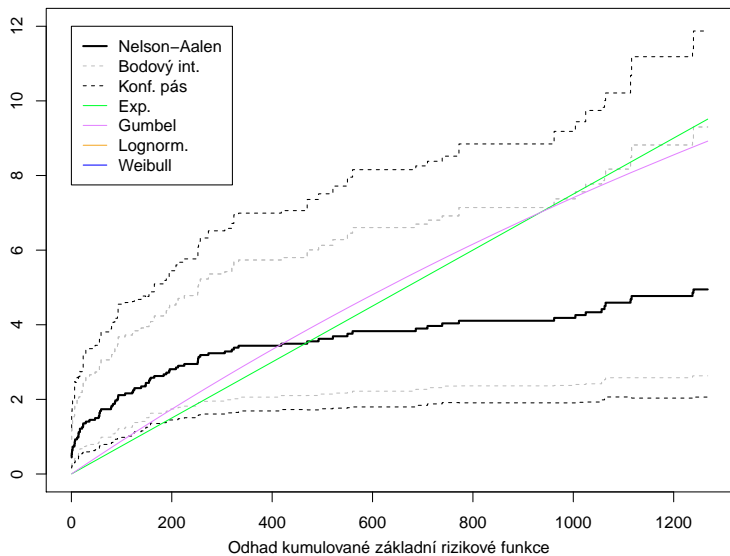
- Hypotézu zamítáme, pokud testovaná kumulovaná základní riziková funkce neleží v konfidenčním pásu.

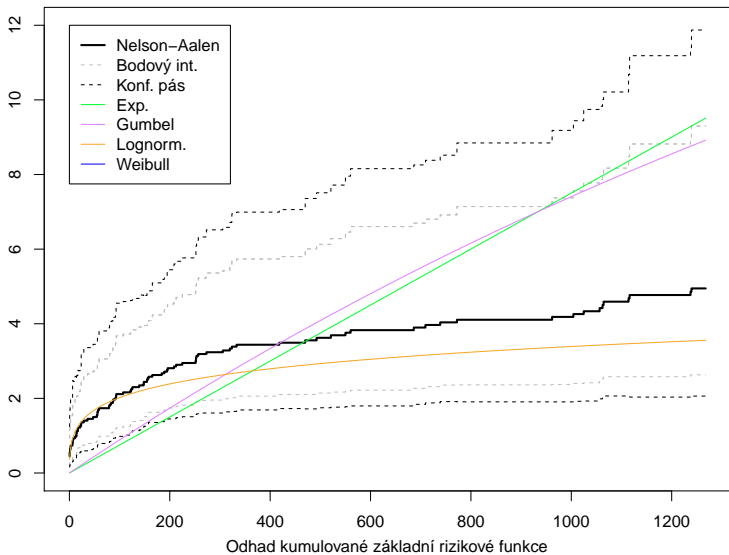


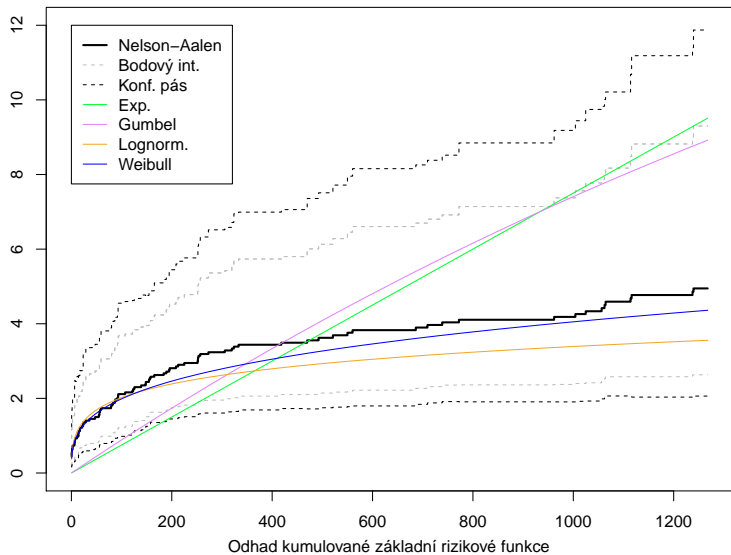












Generovaná data:

- Z Coxova i AFT modelu, $n = 20$ a $n = 50$, $n_i = 10$ s různými základními rizikovými funkcemi a parametry.
- Oprava zvýšila riziko resp. zrychlila čas ($\sigma = 0.1$) a údržba naopak ($\rho = -0.1$).
- Základní rozdělení:
 - Weibullovo $\lambda(t) = a\lambda^a t^{a-1}$,
 - useknuté Gumbelovo $\lambda(t) = \lambda a^t$,
 - Log-normální.
- 500 simulací

Testované hypotézy

- Testovali jsme, zda je základní rozdělení Exp., Weib., Gumb. nebo LN. s parametry odhadnutými metodou maximální věrohodnosti původního modelu považovanými za pevné.
- Sledovali jsme podíl zamítnutých hypotéz na intervalu mezi 5% a 95% kvantilem generovaných dat, $\alpha = 0.05$.

Simulační studie - výsledek

Generované rozdělení					Testované rozdělení - podíl zamítnutí			
Model	λ_0	λ	a	n	Exp.	Weibull	Gumbel	Ln
Cox	Weibull	1/10	5	20	0.908	0	0	0.216
				50	1	0	0	0.276
	Weibull	1/10	1/2	20	0.934	0	0.916	0.036
				50	1	0	1	0.134
	Gumbel	1/10	1.2	20	0.096	0.008	0	0.272
				50	0.642	0.02	0	0.64
	LN	$\mu=2$	$\sigma=2$	20	0.992	0	1	0
				50	1	0.082	1	0
AFT	Weibull	1/10	5	20	0.912	0.006	0.008	0.066
				50	1	0	0	0.492
	Weibull	1/10	1/2	20	0.904	0.002	0.796	0.088
				50	1	0	1	0.644
	Gumbel	1/10	1.2	20	0.372	0.014	0.008	0.592
				50	0.99	0.05	0	0.994
	LN	$\mu=2$	$\sigma=2$	20	0.996	0.142	0.878	0.092
				50	1	0.022	1	0

- Testy jsou s vyšším počtem pozorovaných zařízení přesnější, tj. nezamítají původní a zamítají ostatní základní rozdělení.
- U dat s Weibullovým základním rozdělením záleží, zda je Λ_0 konvexní či konkávní, podle toho je spíše zaměnitelné s Gumbelovým nebo lognormálním rozdělením.
- Exponenciální rozdělení je zamítáno skoro vždy - zde se nabízí srovnání s parametrickým testem zda $a = 1$ ve Weibullově rozdělení.

- Zkoumali jsme metody pro testování hypotéz o tvaru základního rizika při modelování vlivu údržby a oprav na životnost sledovaného zařízení.
- Pro data z Coxova modelu i modelu zrychleného času jsme představili asymptotický test a na simulovaných datech zkoumali jeho vlastnosti.

Outlook:

- Zohlednění variability testovaných parametrických odhadů - zkoumání $\sqrt{n}(\hat{\Lambda}_0(t, \hat{\beta}) - \Lambda_0(t, \hat{\gamma}))$.
- Testy dobré shody.

Děkuji za pozornost

Reference

- [1] Lin D.Y., Ying Z.: *Semiparametric inference for the accelerated life model with time-dependent covariates*, Journal of Statistical Planning and Inference 44, 47-63, 1995.
- [2] Percy D.F., Alkali B.M.: *Generalized proportional intensities models for repairable systems*, IMA Journal of Management Mathematics 17, 171-185, 2005.