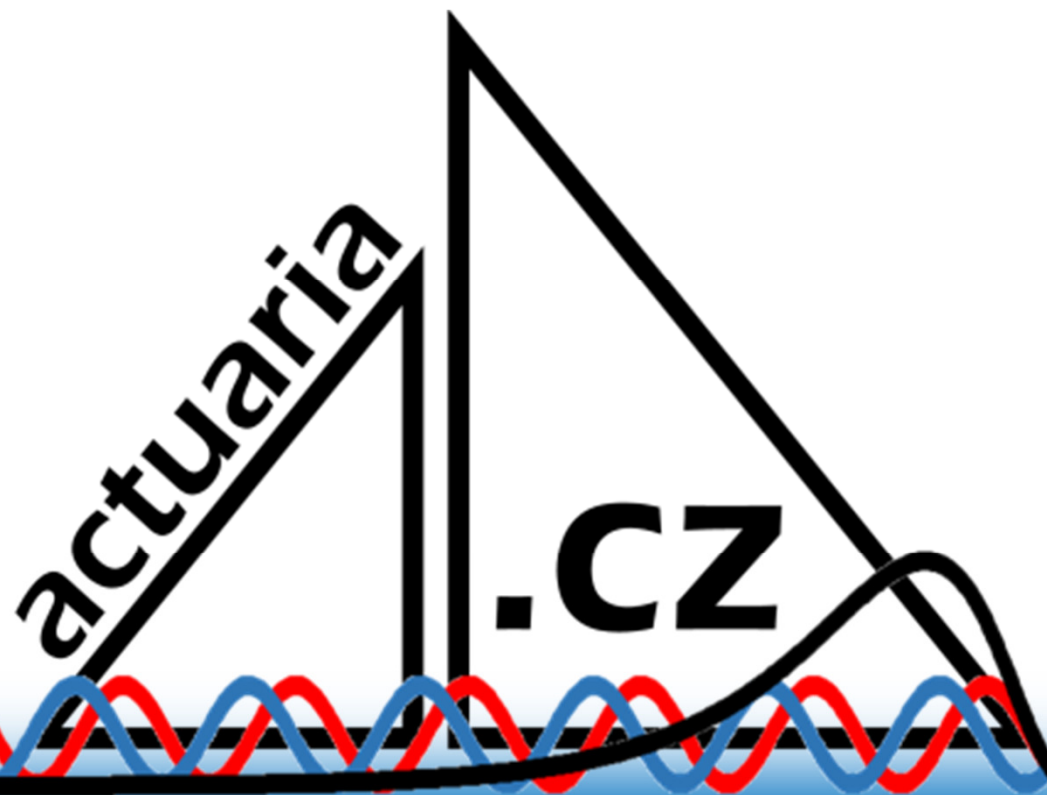


REGRESNÍ MODELY V POJIŠŤOVNICTVÍ



Seminář z aktuárských věd
2. prosince 2016

Kateřina Vlčková

REGRESNÍ MODELY V POJIŠŤOVNICTVÍ

1. PŘEDSTAVENÍ

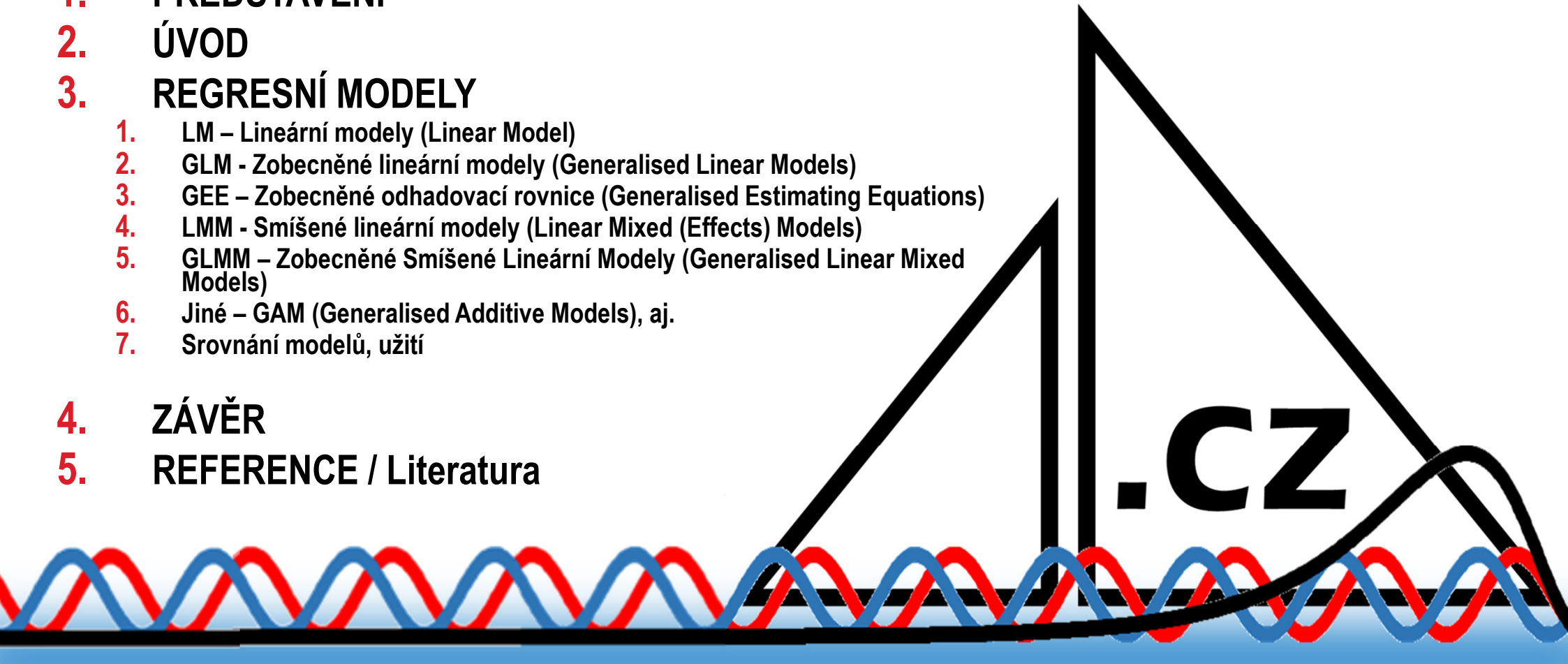
2. ÚVOD

3. REGRESNÍ MODELY

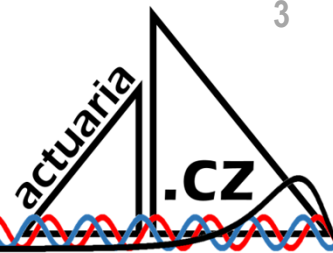
1. LM – Lineární modely (Linear Model)
2. GLM - Zobecněné lineární modely (Generalised Linear Models)
3. GEE – Zobecněné odhadovací rovnice (Generalised Estimating Equations)
4. LMM - Smíšené lineární modely (Linear Mixed (Effects) Models)
5. GLMM – Zobecněné Smíšené Lineární Modely (Generalised Linear Mixed Models)
6. Jiné – GAM (Generalised Additive Models), aj.
7. Srovnání modelů, užití

4. ZÁVĚR

5. REFERENCE / Literatura



REGRESNÍ MODELY V POJIŠŤOVNICTVÍ



▲ CO JE ««MODEL»»??

Model

- ▲ napodobenina předmětu postrádající některé původní vlastnosti
- ▲ objekt s charakteristickými vlastnostmi sloužící pro vytváření podobných objektů
- ▲ kategorie výrobků se společnými parametry

Matematický model

- ▲ Matematický model je abstraktní model používající matematický zápis k popisu chování soustavy (systému).



- ▲ "Remember that all models are wrong; the practical question is how wrong do they have to be to not be useful."

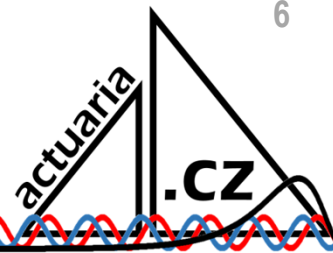
1987 , Empirical Model-Building and Response Surfaces

- ▲ „Essentially, all models are wrong, but some are useful.,“

- ▲ George Edward Pelham Box,
(*18. 10. 1919 – †28. 3. 2013)



REGRESNÍ MODELY V POJIŠŤOVNICTVÍ



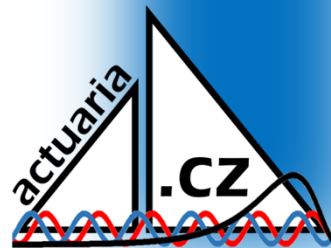
▲ CO JE ««REGRESE»»?

Co je REGRESE?

- ▲ Do statistiky zavedl pojem **REGRESE** britský učenec Francis Galton kolem roku 1880, a to jako „**regres(i) k průměru**“. Tím označil fakt, že např. synové vysokých rodičů jsou sice v průměru (statisticky) vyšší než průměrná populace, zároveň ale individuálně nedosahují extrémních hodnot předchozí generace. Jako kdyby se jedinci postupně "vraceli k průměru". Podobně je tomu i s jinými vlastnostmi, nejen u lidí.

REGRESNÍ ANALÝZA

- ▲ Regresní analýza je označení statistických metod, které umožňují odhadovat hodnotu jisté náhodné veličiny na základě znalosti jiných veličin



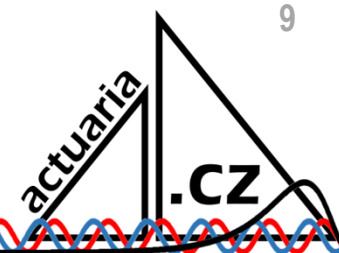


Lineární modely

Linear Models

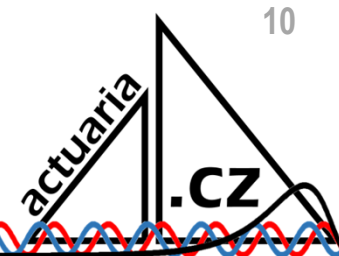
LM

LINEÁRNÍ MODELY



▲ CO JE «« LINEÁRNÍ »» ?

LINEÁRNÍ MODELY



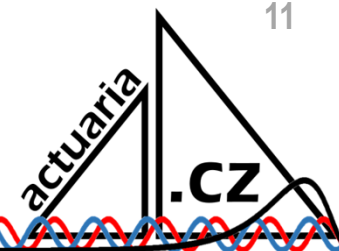
▲ CO JE «« LINEÁRNÍ »»

- ▲ Pomocí X ... regresorů / vysvětlujících proměnných (resp. jejich lineární kombinací $\eta = X\beta$) **modelujeme** $\mathbb{E}Y$

LINEÁRNÍ MODELY

$$Y = X\beta + \varepsilon$$

11



▲ LM (lineární model):

▲ $Y = X\beta + \varepsilon$,

▲ ε ...chybové členy jsou **nezávislé** náhodné veličiny,
t.ž. $\mathbb{E}\varepsilon = \mathbf{0}$, $\text{var}\varepsilon = \sigma^2$

▲ Obvykle předpokládáme Normální lineární model:

▲ $\varepsilon \sim N(0, \sigma^2 I)$ tj. $Y|X \sim N(X\beta, \sigma^2 I)$

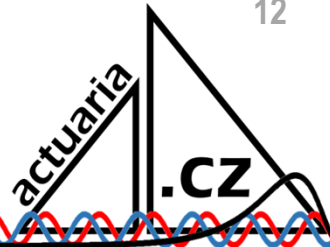
▲ Pozn:

▲ Pomocí X regresorů / vysvětlujících proměnných (resp. jejich lineární kombinací = $X\beta$) **modelujeme** $\mathbb{E}Y$

▲ odezvy Y jsou **homoskedastické** (mají stejný rozptyl)

LINEÁRNÍ MODELY

$$Y = X\beta + \varepsilon$$



▲ (N)LM (Normální lineární model):

▲ $Y = X\beta + \varepsilon$, nezávislé $\varepsilon \sim N_n(\mathbf{0}, \sigma^2 I)$

▲ V Normálním LM platí:

▲ $Y \sim N_n(X\beta, \sigma^2 I)$

▲ $\hat{\beta} \sim N_p(\beta_0, \sigma^2 (X^T X)^{-1})$ - odhad

▲ $\hat{Y} \sim N_n(X\beta_0, \sigma^2 H)$ - odhad

▲ $u \sim N_n(\mathbf{0}, \sigma^2 (I - H))$ rezidua

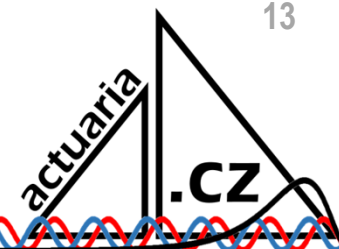
▲ $SS_e / \sigma \sim \chi_{n-p}^2$

▲ SS_e a $\hat{\beta}$ jsou nezávislé

LINEÁRNÍ MODELY

$$Y = X\beta + \varepsilon$$

13



▲ V Normálním LM platí:

$$\text{▲ } \hat{\beta} \sim N_p(\beta, \sigma^2(\mathbb{X}^T \mathbb{X})^{-1})$$

▲ $\hat{\beta}$ je **nestranný** odhad β (tj. $\mathbb{E}\hat{\beta} = \beta$)

▲ $\hat{\beta}$ je MLE – **maximálně věrohodný** - odhad β

▲ $\hat{\beta}$ je **konzistentní** odhad β



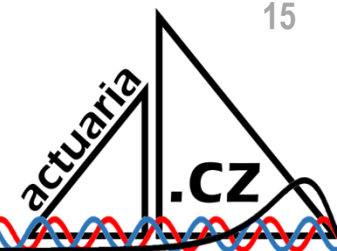
Zobecněné lineární modely

Generalised Linear Models

GLM

ZOBECNĚNÉ LINEÁRNÍ MODELY

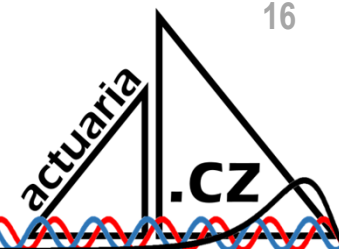
15



▲ CO JE «« ZOBECNĚNÝ »» ?

ZOBECNĚNÉ LINEÁRNÍ MODELY - GLM

16



▲ Lineární model

▲ $Y = X\beta + \varepsilon$, nezávislost, $\varepsilon \sim N(0, \sigma^2 I)$

▲ Lineární model má svá omezení: Co nám vadí ?

▲ Normalita \rightarrow EDF – **Rozdělení exponenciálního typu**

▲ „linearita“ \rightarrow **linková / spojovací funkce**

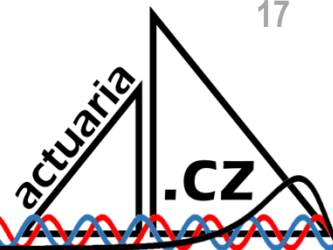
▲ Zobecnění lineárního modelu

▲ Nelder, John A; Wedderburn, Robert W (1972). "Generalized linear models". *Journal of the Royal Statistical Society, Series A. Royal Statistical Society.* 135 (3): 370–384.

▲ Zobecnění – definování společných předpokladů pro metody, které si již dříve používaly:
Logistická regrese, Poissonovská regrese, Gamma regrese, atd.

ROZDĚLENÍ EXPONENCIÁLNÍHO TYPU

17

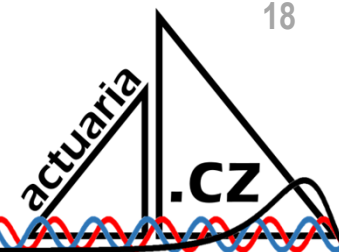


- ▲ Exponential Dispersion Family , EDF
- ▲ Rodina exponenciálních rozdělení
- ▲ Rozdělení s hustotou ve tvaru

$$f(y; \theta, \phi) = \exp \left\{ \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right\}, \quad y \in \mathbb{R}$$

- ▲ $Y \in EDF$

ROZDĚLENÍ EXPONENCIÁLNÍHO TYPU

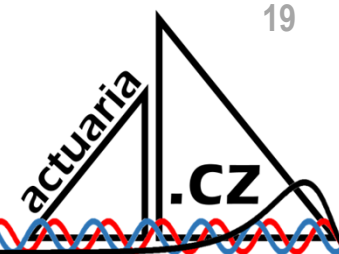


▲ EDF...hustota

$$f(\mathbf{y}; \boldsymbol{\theta}, \boldsymbol{\phi}) = \exp \left\{ \frac{\mathbf{y}\boldsymbol{\theta} - b(\boldsymbol{\theta})}{a(\boldsymbol{\phi})} + c(\mathbf{y}, \boldsymbol{\phi}) \right\}, \quad \mathbf{y} \in \mathbb{R}$$

- ▲ $\boldsymbol{\phi}$...disperzní (škálový) parametr, neznámý...společný pro celý model
- ▲ $\boldsymbol{\theta}$... kanonický parametr, neznámý... θ_i ... závisí na pozorování
- ▲ a, b, c ... funkce, známé
 - ▲ b ... kumulantová funkce, dvakrát spojitě diferencovatelná
 - ▲ a ... obvykle: $a(\boldsymbol{\phi}) = \boldsymbol{\phi}$ nebo $a(\boldsymbol{\phi}) = \boldsymbol{\phi}/w$, w ... váhy

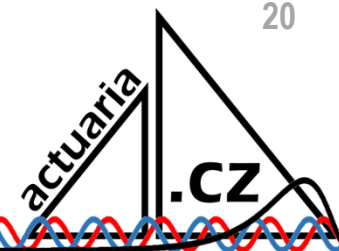
ROZDĚLENÍ EXPONENCIÁLNÍHO TYPU



- ▲ Normální rozdělení
- ▲ Gama rozdělení
- ▲ Inverzní Gaussovo rozdělení
- ▲ Poissonovo rozdělení
- ▲ Binomické rozdělení
- ▲ Negativně binomické rozdělení
- ▲ Geometrické rozdělení
- ▲ Alternativní rozdělení
- ▲ Exponenciální rozdělení
- ▲ Chí-kvadrát rozdělení
- ▲ Weilbullovo rozdělení
- ▲ Paretovo rozdělení

ROZDĚLENÍ EXPONENCIÁLNÍHO TYPU

20

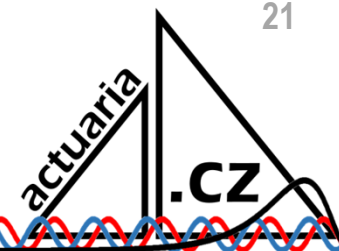


- ▲ Y_i ...náhodná veličina , $Y_i \in EDF$
 - ▲ Pro její momenty platí (b ... dvakrát spojitě diferencovatelná funkce):
- ▲ $\mathbb{E}Y_i = b'(\theta_i) = \mu$
- ▲ $\text{var } Y_i = a(\phi)b''(\theta_i) = \varphi b''(\theta_i) \quad \left(= \frac{\varphi}{w} b''(\theta_i) \right)$
- ▲ $V(\mu)$... rozptylová (varianční) funkce: $V(\mu) \equiv b''[(b')^{-1}(\mu)]$
- ▲ $\text{var } Y_i = a(\phi)V(\mu) = \varphi V(\mu)$

ROZPTYLOVÁ FUNKCE

$V(\mu)$

21



- ▲ $V(\mu)$... rozptylová funkce (*variance function*)
- ▲ definovaná vtahem: $V(\mu) = b''[(b')^{-1}(\mu)] = \varphi V(\mu)$
- ▲ určuje vztah mezi střední hodnotou a rozptylem
- ▲ $\text{var } Y = a(\phi) V(\mu) = \varphi V(\mu)$
- ▲ Jednoznačně identifikuje konkrétní rozdělení z EDF

EDF – SPOJITÁ ROZDĚLENÍ

▲ Normální:

▲ $V(\mu) = \mu$

▲ Gama:

▲ $V(\mu) = \mu^2$

▲ Inverzní Gaussovo:

▲ $V(\mu) = \mu^3$

EDF – DISKRÉTNÍ ROZDĚLENÍ

▲ Poissonovo:

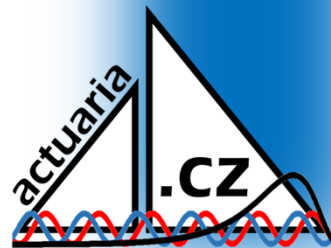
▲ $V(\mu) = \mu$

▲ Binomické:

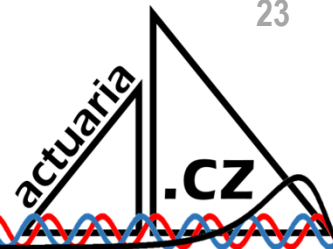
▲ $V(\mu) = n \pi (1 - \pi)$

▲ Negativně binomické:

▲ $V(\mu) = \mu (1 - \mu \kappa)$



ROZDĚLENÍ EXPONENCIÁLNÍHO TYPU

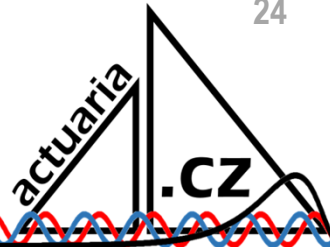


Rozdělení	θ	$b(\theta)$	$a(\phi) = \varphi$	$\mu = EY$	$V(\mu) = \text{var}Y/\varphi$
$B(n, \pi)$	$\ln \frac{\pi}{1-\pi}$	$n \ln(1 + e^\theta)$	1	$n\pi$	$n\pi(1 - \pi)$
$P(\mu)$	$\ln \mu$	e^θ	1	μ	μ
$N(\mu, \sigma^2)$	μ	$\frac{1}{2}\theta^2$	σ^2	μ	1
$G(\mu, \nu)$	$-\frac{1}{\mu}$	$-\ln(-\theta)$	$\frac{1}{\nu}$	μ	μ^2
$IG(\mu, \sigma^2)$	$-\frac{1}{2\mu^2}$	$-\sqrt{-2\theta}$	σ^2	μ	μ^3
$NB(\mu, \kappa)$	$\ln \frac{\kappa\mu}{1+\kappa\mu}$	$-\frac{1}{\kappa} \ln(1 - \kappa e^\theta)$	1	μ	$\mu(1 + \kappa\mu)$

Zdroj: P.de Jong, G.Z. Heller – Generalized Linear Models for Insurance Data

ZOBECNĚNÉ LINEÁRNÍ MODELY – od LM ke GLM

24



▲ LM (lineární model):

▲ $Y = X\beta + \varepsilon$, $\varepsilon \sim N(0, \sigma^2 I)$ tj. $Y|X \sim N(X\beta, \sigma^2 I)$

▲ GLM zobecňuje (Normální) Lineární model:

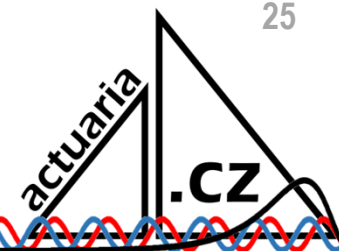
▲ Rozdělení Y ...odezvy / vysvětlované proměnné... nemusí být normální (ani se mu blížit)

▲ Pomocí X regresorů / vysvětlujících proměnných (resp. jejich lineární kombinací $\eta = X\beta$) **nemodelujeme** $\mathbb{E}Y$, ale její transformaci $g(\mathbb{E}Y)$

▲ Důsledek: odezvy Y mohou být (a obvykle jsou) **heteroskedastické**

ZOBECNĚNÉ LINEÁRNÍ MODELY – 3 pilíře GLM

25



Definice modelu GLM

▲ Rozdělení exponenciálního typu... $Y \dots f \in \text{EDF}$

$$\text{▲ } f(y; \theta, \varphi) = \exp \left\{ \frac{y\theta - b(\theta)}{\varphi} + c(y, \varphi) \right\}, \quad y \in \mathbb{R}$$

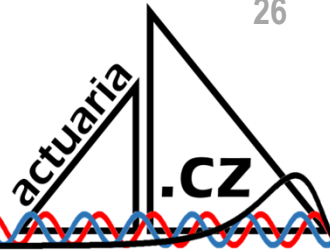
▲ Lineární prediktor... η (čti éta)... lineární kombinace regresorů

$$\text{▲ } \eta = X\beta$$

▲ Linková (spojovací) funkce... g ... striktně monotónní, dvakrát spojitě diferencovatelná

$$\text{▲ } g(\mu) = \eta = g(X\beta) = g(\mathbb{E}Y) \qquad \mu = \mathbb{E}Y = g^{-1}(\eta)$$

▲ Další předpoklady: 1. rozdělení Y závisí na X . 2a) (Y, X) nezávislé N.Ve. Nebo 2b) Y jsou nezávislé N.V. a X měřené konstanty.



GLM – KANONICKÝ LINK

- ▲ **Linková (spojovací) funkce**... g ... striktně monotónní, dvakrát spojitě diferencovatelná

$$\blacktriangle g(\boldsymbol{\mu}) = \boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} = g(\mathbb{E}Y)$$

$$\boldsymbol{\mu} = \mathbb{E}Y = g^{-1}(\boldsymbol{\eta})$$

- ▲ **Kanonický link**

$$\blacktriangle g \equiv (\mathbf{b}')^{-1} \implies g(\boldsymbol{\mu}) = \boldsymbol{\eta} = \boldsymbol{\theta},$$

- ▲ tj. lineární prediktor je roven kanonickému parametru

$$\blacktriangle \text{Platí : } g'(\boldsymbol{\mu}) = \mathbf{1}/V(\boldsymbol{\mu})$$

$$V(\boldsymbol{\mu}) \equiv \mathbf{b}''[(\mathbf{b}')^{-1}(\boldsymbol{\mu})]$$

- ▲ Kanonický link zjednodušuje vztahy při odhadu parametrů

Linková funkce

▲ Linková funkce:

▲ $g(\mu) = \eta = X\beta = g(\mathbb{E}Y)$;

▲ $\mu = \mathbb{E}Y = g^{-1}(\eta)$

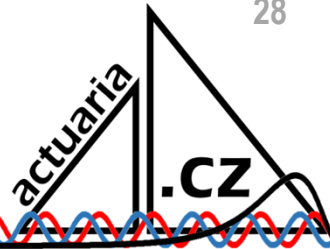
▲ Kanonický link:

▲ $g \equiv (b')^{-1} \Rightarrow$

▲ $g(\mu) = \eta = \theta$,

Příklady

Linková funkce	$g(\mu)$	Rozdělení y
identita	μ	Normální
log	$\ln(\mu)$	Poisson
inverze	$\frac{1}{\mu}$	Gama
mocnina	$\mu^p = \mu^{-1}$	Gama $p = -1$
mocnina	$\mu^p = \mu^{-2}$	Inverzní Gauss $p = -2$
odmocnina	$\sqrt{\mu}$	
logit	$\ln \frac{\mu}{1 - \mu}$	Binomické



GLM – VZTAHY MEZI PARAMETRY

$$X\beta = \eta \quad \eta = g(\mu) \quad \mu = b'(\theta)$$

$$\beta \longrightarrow \eta \longleftarrow \mu \longleftarrow \theta$$

$$g^{-1}(\eta) = \mu \quad (b')^{-1}(\mu) = \theta$$

η ... lineární prediktor

β ... regresní parametr

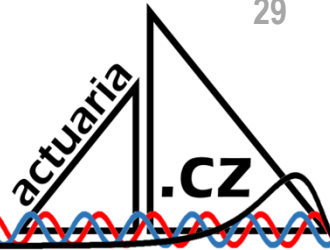
g ... linková funkce

μ ... střední hodnota EY ,

θ ... kanonický parametr

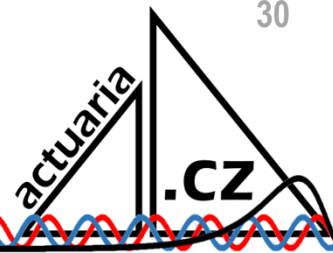
b ... kumulantová funkce

7 KROKŮ K MODELU GLM

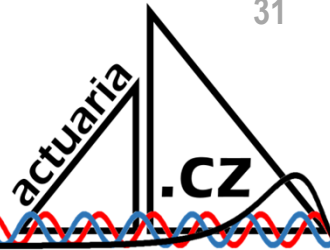


- ▲ Volba rozdělení odezvy / vysvětlované proměnné... $f(y)$
- ▲ Volba linkové funkce... $g(\mu)$
- ▲ Volba regresorů / vysvětlujících proměnných... x
- ▲ Výběr (nezávislých) dat: pozorování y_1, \dots, y_n a odpovídajících hodnot regresorů x_1, \dots, x_n
- ▲ Odhad regresních parametrů β a případně disperzního parametru φ
- ▲ Model fit: kvalita modelu, výběr podmodelů
- ▲ Kalibrace modelu, odhad predikční chyby modelu

GLM – 1) volba rozdělení odezvy ... $f(y)$



- ▲ Normální rozdělení
- ▲ Gama rozdělení
- ▲ Inverzní Gaussovo rozdělení
- ▲ Poissonovo rozdělení
- ▲ Binomické rozdělení
- ▲ Negativně binomické rozdělení
- ▲ Alternativní rozdělení
- ▲ Jiné \in EDF

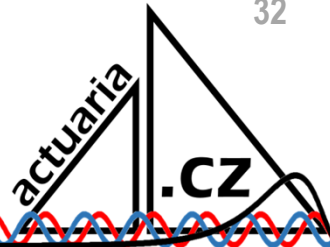


GLM – 2) volba linkové funkce... $g(\mu)$

- ▲ kanonický link
- ▲ jiná linková funkce

Linková funkce	$g(\mu)$	Rozdělení y
identita	μ	Normální
log	$\ln(\mu)$	Poisson
inverze	$\frac{1}{\mu}$	Gama
mocnina	$\mu^p = \mu^{-1}$	Gama $p = -1$
mocnina	$\mu^p = \mu^{-2}$	Inverzní Gauss $p = -2$
odmocnina	$\sqrt{\mu}$	
logit	$\ln \frac{\mu}{1 - \mu}$	Binomické

GLM – 3) volba regresorů ... x



▲ REGRESORY – Vysvětlující proměnné

- ▲ **Metrické proměnné:** věk, výše škody
- ▲ **Faktory:** pohlaví, typ vozidla, počet dětí
- ▲ **Interakce:** MxM, FxF, MxF

▲ Metrické proměnné:

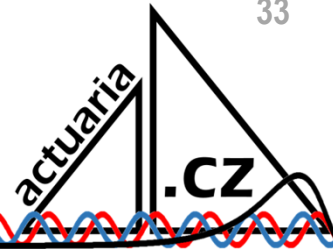
- ▲ model odhaduje vliv jednotlivých proměnných na odezvu
- ▲ β_j ... vyjadřuje změnu při jednotkové změně regresoru x
- ▲ Lze i polynomy

▲ Faktory:

- ▲ Diskrétní proměnné nebo kategorie
- ▲ Dummy proměnné (o 1 kategorii méně, referenční kategorie - typická)
- ▲ β_j ... vyjadřuje změnu oproti referenční kategorii

▲ Interakce: $\beta_{1,2}x_1x_2$

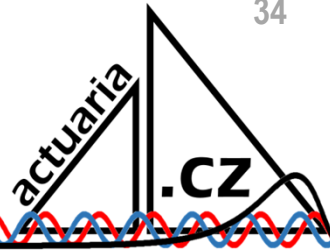
GLM – 3) volba regresorů - OFFSET



- ▲ Člen v lineárním prediktoru s pevně daným koeficientem
- ▲ Regresní koeficient roven 1, neodhaduje se
- ▲ Obvykle použit jako korekce modelu s ohledem na expozici v riziku
 - ▲ velikost skupiny,
 - ▲ různá doba pozorování
- ▲ Nejčastěji: pro logaritmický link; zde příklad pro expozici n_i . řádku:

$$\eta_i = \ln n_i + x_i^T \beta \quad \mu_i = e^{\eta_i} = n_i \cdot e^{x_i^T \beta}$$

- ▲ PRAXE: délka platnosti smlouvy, počet rizik



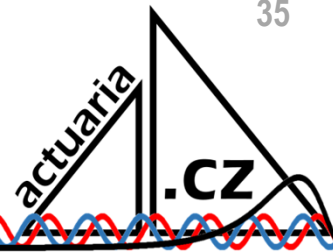
GLM – 3) volba regresorů – váhy pozorování w

- ▲ Do modelu je možné zahrnout apriorní váhy pro jednotlivá pozorování ... w
- ▲ Parametrizace v EDF: $a(\phi) = \phi/w$

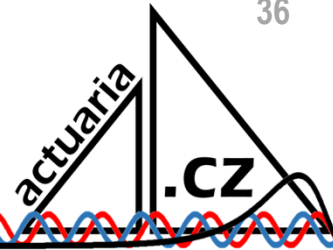
$$E[\mathbf{Y}] = b'(\theta), \quad \text{var}[\mathbf{Y}] = b''(\theta) \cdot \phi/w = V(\theta) \cdot \phi/w$$

- ▲ PRAXE: průměrné výše škody: w ... počet škod na dané smlouvě
- ▲ PRAXE: škodní frekvence: w ... délka expozice / délka platnosti smlouvy

GLM – 4) Výběr dat: $\dots y_1, \dots, y_n$



- ▲ nezávislá data
- ▲ náhodný výběr

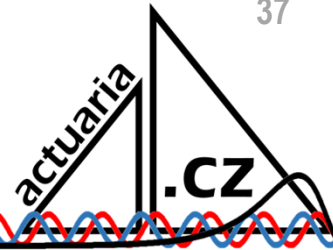


GLM – 5a) odhad regresních parametrů β

▲ ODHADY REGRESNÍCH PARAMETRŮ V GLM modelu – *MODEL FIT*

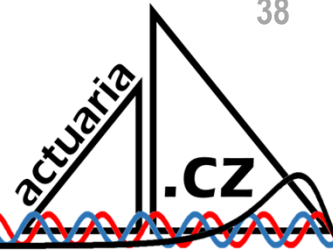
- ▲ Odhady a INFERENCE v GLM jsou založeny na [teorii Maximální věrohodnosti](#)
- ▲ Maximalizace vyžaduje iterativní řešení
 - ▲ Metoda iterativních vážených čtverců? IWLS X IRLS
 - ▲ Newton –Raphsonův iterační algoritmus
 - ▲ Fisherova metoda skóru
- ▲ MLE – Maximálně věrohodné odhady – MAXIMUM LIKELIHOOD
 - ▲ Vlastnosti MLE odhadů:
 - ▲ Asymptoticky nezkrácené, konzistentní, asymptoticky vydatné, invariantní vůči monotónní funkci, asymptotický normální
- ▲ IRLS algoritmus –Iteratively Re-Weighted Least Squares

GLM – 5b) odhad disperzního parametru φ



- ▲ Disperzní (škálový) parametr φ
- ▲ Obvykle není znám
- ▲ Pro MLE odhady regresních parametrů β_0 není nutné znát odhad skutečné hodnoty disperzního parametru φ_0 . MLE odhady jsou stejné v obou případech. Asymptotické vlastnosti platí.
- ▲ Odhad disperzního parametru φ_0 je však potřeba pro odhad asymptotického rozptylu.
- ▲ MLE odhad φ_0 není vždy možné vypočítat. Pro odhady se proto obvykle používá modifikovaná Momentová metoda.
- ▲ Odhad založen na Pearsonově χ^2 statistice

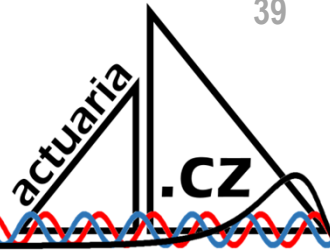
$$\widehat{\varphi}_n = \frac{1}{n-p} \sum_{i=1}^n \frac{(Y_i - \widehat{\mu}_i)^2}{V(\widehat{\mu}_i)} = \frac{1}{n-p} X^2$$



GLM – 6) Výběr modelu

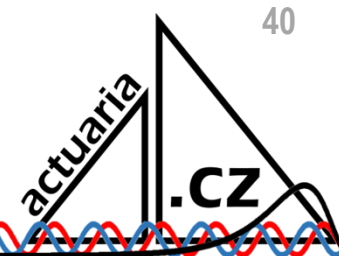
- ▲ Kvalita modelu a testy
- ▲ Saturovaný model:
 - ▲ Počet parametrů = počet pozorování
 - ▲ Prakticky se nepoužívá = perfect fit
 - ▲ Teoretická aplikace = je v něm dosaženo maximální možné věrohodnosti
- ▲ Deviance
 - ▲ $D = 2 * [l(Y) - l(\hat{\beta})]$ (škálová deviance)
 - ▲ Analogie indexu determinace v lineárním modelu
- ▲ Waldovy testy
- ▲ Testy poměrem věrohodnosti

GLM – 6) Výběr modelu – INFORMAČNÍ KRITÉRIA



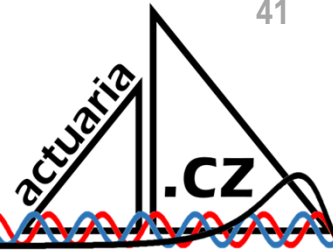
- ▲ Čím více vysvětlujících proměnných v modelu, tím lepší „fit“ modelu, ale: tím více parametrů a horší přesnost jejich odhadů (zvyšuje se rozptyl)
- ▲ => hledáme kompromis (vysoká věrohodnost × malý počet parametrů)
- ▲ **Akaikeho informační kritérium:**
- ▲ $AIC = -2 (l(Y; \hat{\beta}, \hat{\varphi}) - m)$
- ▲ m ... počet parametrů β (penalizace)
- ▲ **Bayesovské informační kritérium BIC**
- ▲ $BIC = -2 l(Y; \hat{\beta}, \hat{\varphi}) + \ln(n) \cdot \dim(\beta)$
- ▲ Vybíráme modely s nízkým AIC a BIC.
- ▲ BIC více penalizuje za počet parametrů, v porovnání s AIC vybírá modely s menším počtem vysvětlujících proměnných (někdy až příliš chudé???)
- ▲ Porovnávané modely musí být založeny na stejných pozorováních

GLM – 6) Výběr modelu – metody výběru (pod)modelů



- ▲ **Best Subset Selection** – regresní model sestaven pro daný počet regresorů a všechny možné jejich kombinace, vybrán nejlepší model pro zvolené kritérium (AIC, BIC, CD), výpočetně náročné až nereálné
- ▲ **Kroková regrese** (Stepwise Regression) – postupné přidávání/ubírání regresorů, postup ve směru největšího poklesu hodnoty kritéria, výpočetně jednodušší
- ▲ **Regularizovaná regrese** – ve věrohodnostní funkci se penalizuje nárůst regresních parametrů
- ▲ **Sekvenční výběr proměnných** – manuální verze krokové regrese

GLM – 7) Kalibrace modelu, odhad predikční chyby



- ▲ **Validace** : rozdělení dat na Trénovací část (kalibrace - odhad modelu) a Testovací část (stanovení predikční chyby)
- ▲ **Cross-validace**: data rozdělena na několik částí: na všech kromě jedné odhadujeme model, na poslední testujeme; opakujeme pro všechny kombinace
- ▲ **Nová data** – pro stanovení predikční chyby

ZOBECNĚNÉ LINEÁRNÍ MODELY

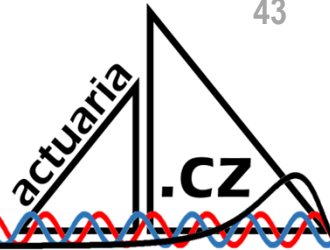
42



▲ CO SE DO «« GLM »»
NEVEŠLO

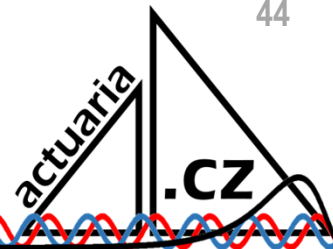
CO DÁL? - ROZŠÍŘENÍ / ZOBECNĚNÍ GLM

43



- ▲ V praxi však nebývají všechny předpoklady splněny...
- ▲ Je možné rozšíření teorie GLM na případy, kdy jsou předpoklady porušeny?
 - ▲ Rozdělení závislých proměnných / odezvy není exponenciálního typu
 - ▲ Rozdělení odezvy není blíže specifikováno
 - ▲ Známe první 2 momenty
 - ▲ Známe vztah mezi $\mathbb{E}Y_i$ a $var Y_i$
 - ▲ Data nejsou nekorelovaná / nezávislá
 - ▲ Regresní parametry / efekty β nejsou pevné, ale náhodné
 - ▲ Data vykazují nadměrnou disperzitu oproti teoretickým hodnotám

CO DÁL? - ROZŠÍŘENÍ / ZOBECNĚNÍ GLM



▲ OVERDISPERSION / NADMĚRNÁ DISPERZITA

- ▲ Pozorovaná data mají vyšší variabilitu/disperzitu, než by se očekávalo při platnosti zvoleného modelu
- ▲ Průměr je většinou možné parametricky upravit, aby odpovídal teoretické hodnotě.
- ▲ Vyšší momenty se však (obzvláště u malých výběrů) upravují těžko
- ▲ S nadměrnou disperzí se často setkáváme u modelů četností (Poissonovo rozdělení, či binomické), obzvláště u malých výběrů a heterogenních populací
 - ▲ Příklad: nadměrná disperze v Binomických datech -> Beta-binomické rozdělení
 - ▲ Příklad: nadměrná disperze v Poissonovských datech -> Poisson-Gama rozdělení

OVERDISPERSION / NADMĚRNÁ DISPERZITA

45



▲ Příklad: Poissonovo rozdělení

▲ $Y_1, \dots, Y_n \sim \mathbf{Po}(\lambda_0)$ nezávislá pozorování

▲ $\text{var}[Y_i] = \mathbb{E}[Y_i] = \lambda_0, \quad Y_i \in EDF$

▲ Parametr $[\lambda_i]$ považujeme za náhodnou veličinu (nikoli parametr),
 $\mathbb{E}[\lambda_i] = \lambda_0$

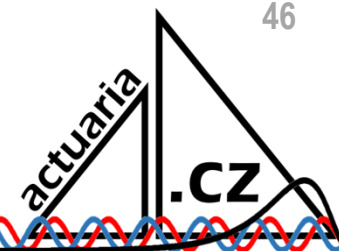
▲ $Y_i | \lambda_0 \sim \mathbf{Po}(\lambda_0), \quad \mathbb{E}[\lambda_i] = \lambda_0 \quad \text{var}[\lambda_i] = \sigma^2_\lambda$

▲ $\mathbb{E}[Y_i] = \lambda_0 \quad \text{var}[Y_i] = \lambda_0 + \sigma^2_\lambda$

▲ Předpokládejme: $\lambda_i \sim \Gamma(a, a\lambda_0)$

▲ -> Poisson-Gama rozdělení - není EDP, GLM neumí řešit

▲ Speciální případ: negativně-binomické rozd., geometrické



CO DÁL? - ROZŠÍŘENÍ / ZOBECNĚNÍ GLM

▲ KVAZIVĚROHODNOST / QUASI-LIKELIHOOD

▲ Neznáme rozdělení, ale známe rozptylovou funkci $V(\cdot)$, která určuje vztah mezi prvními dvěma momenty,

$$\triangle V(\mu) = b''[(b')^{-1}(\mu)] \quad \rightarrow \quad \text{var } Y_i = \varphi V(\mu) = \varphi V(\mathbb{E}Y_i)$$

▲ MODEL:

▲ Mějme n náhodných vektorů (Y_i, \mathbf{X}_i) , $i = 1, \dots, n$, kde

$\mathbf{X}_i = (X_{i1}, X_{i2}, \dots, X_{ip})^T$ jsou regresory a Y_i jsou odezvy.

▲ Předpokládejme:

▲ 1. Y_1, \dots, Y_n jsou nezávislé

▲ 2. $\mu_i = \mathbb{E}[Y_i]$, $g(\mu_i) = \eta_i = \mathbf{X}_i^T \boldsymbol{\beta}_0$, $\eta_i \dots$ lineární prediktor, $g \dots$ linková funkce

▲ 3. $\text{var } Y_i = \varphi V(\mu_i)$, $\varphi \dots$ disperzní parametr, $V \dots$ varianční funkce

▲ GLM je parametrický model, toto semi-parametrický model: rozdělení není dáno, pouze vztah mezi momenty

▲ Metodu maximální věrohodnosti (MLE) lze použít pro odhady pouze pro parametrické modely.

▲ Odhady regresních parametrů – „Metoda maximální kvazivěrohodnosti“ (?MQLE)

$$\triangle \frac{1}{\sqrt{n}} U_n(\boldsymbol{\beta}_0) \xrightarrow{d} N_p(\mathbf{0}, \mathbb{I}(\boldsymbol{\beta}_0)) \quad \sqrt{n}(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) \xrightarrow{d} N_p(\mathbf{0}, \mathbb{I}^{-1}(\boldsymbol{\beta}_0))$$

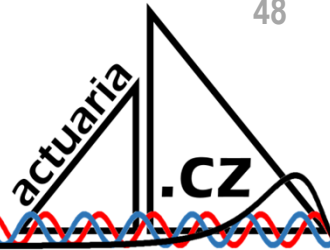
▲ Je možné aplikovat Waldovy a skórové testy; nikoli však devianci či AIC.



Zobecněné odhadovací rovnice Generalised Estimating Equations

GEE

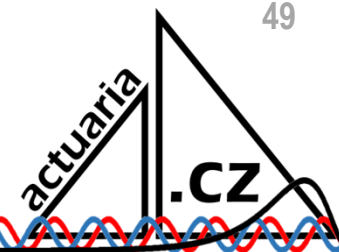
GEE - Zobecněné odhadovací rovnice



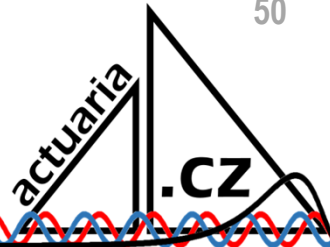
▲ GEE – Generalised Estimating Equations – Zobecněné odhadovací rovnice

- ▲ *rozšíření GLM – korelace, skupinově závislá data*
- ▲ GEE model navržen v článku z r. 1986 – K.Y. Liang a S.L. Zeger: *Longitudinal data analysis using generalized linear models*
- ▲ Umožňuje použít postupy z GLM i v případě, kdy nejsou splněny předpoklady pro GLM model
- ▲ **Data nejsou nekorelovaná**

GEE - DEFINICE MODELU pro skupinově závislá data



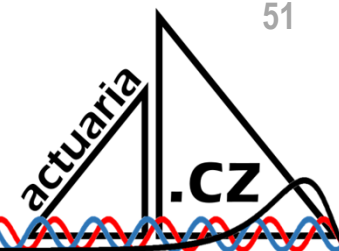
- ▲ **Data:** Y_1, \dots, Y_K nezávislé náhodné vektory $Y_i = (Y_{i1}, \dots, Y_{in_i})^T, i = 1 \dots K, \sum_{i=1}^K n_i = n$
- ▲ Data rozdělena do K nezávislých skupin (shluků, subjektů, jedinců)
- ▲ V každé skupině je různý počet (n_i) vzájemně korelovaných pozorování
- ▲ Tj. data jsou: **závislá v rámci jedné skupiny a nezávislá mezi jednotlivými skupinami** (př. zuby pacienta, mláďata z jednoho vrhu, škody na jedné pojistné smlouvě, budovy v jedné obci)
- ▲ **Shluková data** (bez uspořádání), **opakovaná měření** (uspořádání v čase), **longitudinální data** (se záznamem o čase), **panelová data** (ekonomická data)
- ▲ Ke každému pozorování (odezvě, závisle/vysvětlované proměnná) Y_{ij} přísluší vektor regresorů (nezávislých/vysvětlujících proměnných) $X_{ij} = (X_{ij1}, \dots, X_{ijp})^T, i = 1 \dots K, j = 1, \dots, n_i, \sum_{i=1}^K n_i = n.$



GEE - DEFINICE MODELU pro skupinově závislá data

- ▲ **Data:** Y_1, \dots, Y_K nezávislé náhodné vektory $Y_i = (Y_{i1}, \dots, Y_{in_i})^T, i = 1 \dots K, \sum_{i=1}^K n_i = n$
- ▲ Ke každému pozorování (odezvě) Y_{ij} přísluší vektor regresorů $X_{ij} = (X_{ij1}, \dots, X_{ijp})^T, i = 1, \dots, K, j = 1, \dots, n_i, \sum_{i=1}^K n_i = n.$
- ▲ **Cíl:** popsat závislost střední hodnoty pozorování $\mu_{ij} = \mathbb{E}Y_{ij}$ na regresorech X_{ij} pomocí regresního modelu.
- ▲ **Předpoklad:** vztah je definován pomocí linkové funkce $g(\mu_{ij}) = X_{ij}^T \beta_0$, stejně jako v GLM
- ▲ g ... linková funkce: striktně monotónní, dvakrát spojitě diferencovatelná
- ▲ $\beta_0 = (\beta_{01}, \dots, \beta_{0p})^T$ neznámý vektor regresních parametrů (skutečná hodnota)
- ▲ Předpokládejme :
 - ▲ $\mathbb{E}Y_i = \mu_i = (\mu_{i1}, \mu_{i2}, \dots, \mu_{in_i})^T = (g^{-1}(X_{i1}^T \beta_0), \dots, g^{-1}(X_{in_i}^T \beta_0))^T$
 - ▲ $var Y_i$ nespecifikováno, bez předpokladů o rozptylu či kovarianci

GEE - ODHADY REGRESNÍCH PARAMETRŮ v modelu pro skupinově závislá data



▲ $\mathbb{X}_i = (\mathbf{X}_{i1}^T, \dots, \mathbf{X}_{ip}^T)^T$... regresní matice ($n_i \times p$)

$$\left(\frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} \right)_{(p \times n_i)} = \mathbb{X}_i^T \begin{pmatrix} g'(\mu_{i1}) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & g'(\mu_{in_i}) \end{pmatrix}^{-1} \quad \dots \text{matice partiálních derivací}$$

▲ Připomenutí:

skórová funkce GLM modelu: $U_i(\boldsymbol{\beta}) = \frac{\partial \mu_i}{\partial \boldsymbol{\beta}} \frac{1}{\phi V(\mu_i)} (Y_i - \mu_i)$, kde $V(\mu_i)$ je rozptylová funkce.

▲ Zobecnění pro vícerozměrné \mathbf{Y}_i :

pseudo-skórová funkce příslušející i . skupině pozorování: $U_i(\boldsymbol{\beta}) = \frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} \mathbb{Q}_i^{-1}(\boldsymbol{\mu}_i) (Y_i - \boldsymbol{\mu}_i)$.

▲ $\mathbb{Q}_i(\boldsymbol{\mu}_i) = \boldsymbol{\phi} \mathbb{V}_i^{1/2}(\boldsymbol{\mu}_i) \mathbb{R}_i \mathbb{V}_i^{1/2}(\boldsymbol{\mu}_i)$... „pracovní kovarianční matice“, která reprezentuje naši představu (guess) o rozptylu $\text{var } \mathbf{Y}_i$.

▲ \mathbb{V}_i ... diagonální matice, na diagonále členy $V(\mu_{i1}), \dots, V(\mu_{in_i})$, které představují „pracovní rozptyl“ pro pozorování $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^T$.

▲ \mathbb{R}_i ... „pracovní korelační matice“ ($n_i \times n_i$).

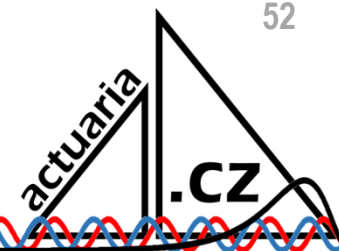
▲ Nepředpokládáme, že \mathbb{R}_i a $V(\mu_{ij})$, jsou známé či správně odhadnuté. (V opačném případě bychom rovnou použili $\text{cor } \mathbf{Y}_i$ a $\text{var } Y_{ij}$ místo \mathbb{R}_i a $\phi V(\mu_{ij})$.)

▲ **GEE odhad regresních parametrů $\boldsymbol{\beta}$** : $\widehat{\boldsymbol{\beta}}_K$ je definován jako řešení **soustavy „zobecněných odhadovacích rovnic“**:

$$U_K(\widehat{\boldsymbol{\beta}}_K) = \sum_{i=1}^K U_i(\widehat{\boldsymbol{\beta}}_K) = \mathbf{0} \quad (\text{GEE})$$

GEE - asymptotické vlastnosti odhadů regresních parametrů

52



▲ Označme: $\mathbb{D} = -\mathbb{E} \left[\frac{\partial}{\partial \beta^T} U_i(\boldsymbol{\beta}_0) \right] = \mathbb{E} \left[\left(\frac{\partial \mu_i}{\partial \beta} \right) \mathbb{Q}^{-1}(\boldsymbol{\mu}_i) \left(\frac{\partial \mu_i}{\partial \beta} \right)^T \right]$

▲ Pro $K \rightarrow \infty$ platí:

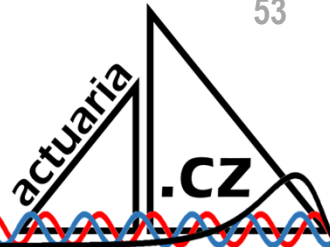
▲ (i) $\widehat{\boldsymbol{\beta}}_K \xrightarrow{p} \boldsymbol{\beta}_0$ (konzistence)

▲ (ii) $\frac{1}{\sqrt{K}} U_K(\boldsymbol{\beta}_0) \xrightarrow{d} N_p(\mathbf{0}, \boldsymbol{\Sigma})$

▲ (iii) $\sqrt{K}(\widehat{\boldsymbol{\beta}}_K - \boldsymbol{\beta}_0) \xrightarrow{d} N_p(\mathbf{0}, \mathbb{D}^{-1} \boldsymbol{\Sigma} \mathbb{D}^{-1})$. (asymptotická normalita)

▲ Asymptotika funguje pro velká $K \rightarrow \infty$. **Potřebujeme velký počet nezávislých skupin.**

▲ Počet (korelovaných) pozorování v jednotlivých skupinách není podstatný.



GEE - odhad asymptotického rozptylu

▲ Asymptotický rozptyl $\mathbb{D}^{-1}\Sigma\mathbb{D}^{-1}$ odhadujeme pomocí sendvičového odhadu (sandwich estimator / White estimator):

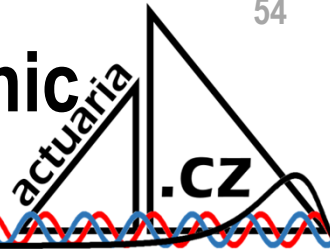
▲ $\mathbb{D}^{-1}\widehat{\Sigma}\mathbb{D}^{-1} = \widehat{\mathbb{D}}^{-1}\widehat{\Sigma}\widehat{\mathbb{D}}^{-1}$

▲ kde:
$$\widehat{\mathbb{D}} = \frac{1}{K} \sum_{i=0}^n \left(\frac{\partial \hat{\mu}_i}{\partial \beta} \right) \mathbb{Q}^{-1}(\hat{\mu}_i) \left(\frac{\partial \hat{\mu}_i}{\partial \beta} \right)^T$$

▲ a
$$\widehat{\Sigma} = \frac{1}{K} \sum_{i=1}^K U_i(\hat{\beta}_K)^{\otimes 2}.$$

- ▲ $\mathbb{Q}(\mu_i) = \varphi V_i^{1/2}(\mu_i) \mathbb{R}_i V_i^{1/2}(\mu_i) \dots$ „pracovní kovarianční matice“, která reprezentuje naši představu (guess) o rozptylu $\text{var } Y_i$.
- ▲ $V_i \dots$ diagonální matice, na diagonále členy $V(\mu_{i1}), \dots, V(\mu_{in_i})$, které představují „pracovní rozptyly“ pro pozorování $Y_i = (Y_{i1}, \dots, Y_{in_i})^T$.
- ▲ $\mathbb{R}_i \dots$ „pracovní korelační matice“ ($n_i \times n_i$).
- ▲ Nepředpokládáme, že \mathbb{R}_i a $V(\mu_{ij})$, jsou známe či správně odhadnuté.

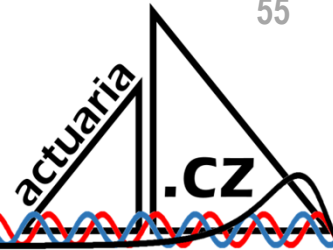
GEE - řešení soustavy zobecněných odhadovacích rovnic



$$U_K(\widehat{\boldsymbol{\beta}}_K) = \sum_{i=1}^K U_i(\widehat{\boldsymbol{\beta}}_K) = \mathbf{0} \quad (\text{GEE})$$

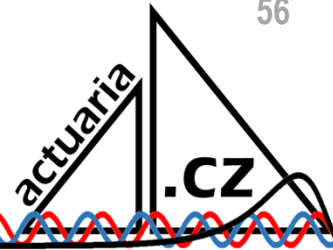
- ▲ Soustava (GEE) se řeší iteračně.
- ▲ GEE odhad $\widehat{\boldsymbol{\beta}}_K$ hledáme pomocí modifikované metody IWLS (Iterative Re-Weighted Least Squares).
- ▲ Iterujeme: $\widehat{\boldsymbol{\beta}} = \left[\sum_{i=0}^n \left(\frac{\partial \widehat{\boldsymbol{\mu}}_i}{\partial \boldsymbol{\beta}} \right) \mathbf{Q}^{-1}(\widehat{\boldsymbol{\mu}}_i) \left(\frac{\partial \widehat{\boldsymbol{\mu}}_i}{\partial \boldsymbol{\beta}} \right)^T \right]^{-1} \left[\sum_{i=0}^n \left(\frac{\partial \widehat{\boldsymbol{\mu}}_i}{\partial \boldsymbol{\beta}} \right) \mathbf{Q}^{-1}(\widehat{\boldsymbol{\mu}}_i) \widehat{\mathbf{Z}}_i \right]$,
- ▲ kde: $\widehat{\mathbf{Z}}_i = (\widehat{Z}_{i1}, \dots, \widehat{Z}_{in_i})^T$ a $\widehat{Z}_{ij} = \frac{\widehat{\eta}_{ij} + (Y_{ij} - \widehat{\mu}_{ij})g'(\widehat{\mu}_{ij})}{g'(\widehat{\mu}_{ij})}$

GEE – volba korelační struktury



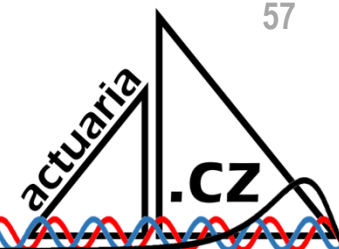
- ▲ Rozptylová funkce $V(\mu)$ v GLM vyjadřuje vztah mezi rozptylem a střední hodnotou.
- ▲ Pracovní kovarianční matice $\mathbb{Q}(\mu_i) = \varphi V_i^{1/2}(\mu_i) \mathbb{R}_i V_i^{1/2}(\mu_i)$,
resp. pracovní korelační matice \mathbb{R}_i ,
reprezentuje v GEE rovněž naši představu o rozptylu.
- ▲ Jak zvolit vhodnou pracovní korelační matici?

GEE – volba korelační struktury



- ▲ **pracovní nezávislost / working independence**
- ▲ **metoda parametrizované korelace / parametrized correlation**
 - ▲ pásová korelace 1. řádu
 - ▲ pásová korelace m. řádu:
 - ▲ exchangeable correlations
 - ▲ AR(1) korelace a další korelace na bázi časových řad

GEE – volba korelační struktury pracovní nezávislost



- ▲ Nejjednodušší volba, analyzujeme **data**, jako by byla **nezávislá**.
- ▲ Volíme pracovní korelační matici : $\mathbb{R}_i \equiv \mathbb{I}_{n_i}$ (čtvercová jednotková matice $(n_i \times n_i)$).

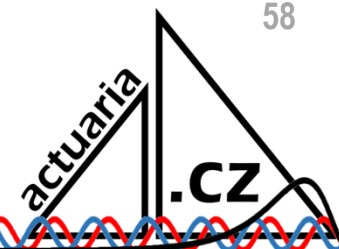
- ▲ Pracovní kovarianční matice je diagonální: $\mathbb{Q}(\boldsymbol{\mu}_i) = \varphi \mathbb{V}_i(\boldsymbol{\mu}_i) = \begin{pmatrix} \varphi V(\mu_{i1}) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \varphi V(\mu_{in_i}) \end{pmatrix}$.

- ▲ Odhady regresních parametrů $\boldsymbol{\beta}_0$ jako při nezávislých datech, pomocí standardního IWLS algoritmu.
- ▲ Rozptyl odhadů je poté upraven pomocí sendvičového odhadu, zohlednění možných korelací a případné nevhodné volby rozptylové funkce $V(\cdot)$.
- ▲ Odhady jsou konzistentní. V případě velké korelace mezi daty nejsou eficientní.

GEE – volba korelační struktury

metoda parametrizované korelace

58



▲ Volíme **kovarianční strukturu**, která není nezávislá.

▲ Zavedeme nový vektorový parametr $\alpha \in \mathbb{R}^m$.

▲ Položíme:

▲ pracovní korelační matice: $\mathbb{R}_i \equiv \mathbb{R}_i(\alpha)$

▲ pracovní korelační matice: $\mathbb{Q}_i(\mu_i) \equiv \mathbb{Q}_i(\mu_i, \alpha) = \varphi V_i^{1/2}(\mu_i) \mathbb{R}_i(\alpha) V_i^{1/2}(\mu_i)$

▲ odhad pracovní korelační matice: $\widehat{\mathbb{Q}}_i(\mu_i) \equiv \mathbb{Q}_i(\mu_i, \hat{\alpha}) = \hat{\varphi} V_i^{1/2}(\mu_i) \mathbb{R}_i(\hat{\alpha}) V_i^{1/2}(\mu_i)$

▲ kde $\hat{\alpha}$ je \sqrt{K} -konzistentní odhad parametru α , např. momentový odhad.

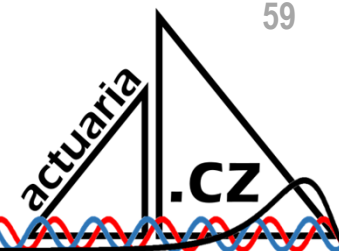
▲ Dostáváme modifikované skóre: $\widehat{U}_i(\beta) = \frac{\partial \mu_i}{\partial \beta} \widehat{\mathbb{Q}}_i^{-1}(\mu_i) (Y_i - \mu_i)$.

(! Ztráta nezávislosti vektorů)

▲ Odhady regresních parametrů $\widehat{\beta}$ jsou řešením soustavy $U_K(\widehat{\beta}) = \sum_{i=1}^K \widehat{U}_i(\widehat{\beta}) = \mathbf{0}$.

GEE – volba korelační struktury

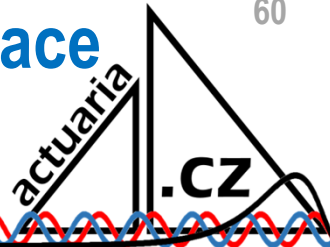
obecný postup odhadu parametru α



- ▲ Obecný postup při odhadování parametru α : (založen na reziduích):
- ▲ Odhad β za předpokladu pracovní nezávislosti (Volíme: $\mathbb{R}_i = \mathbb{I}_{n_i}$)
- ▲ Spočteme Personova rezidua:
$$r_{ij}^P = \frac{Y_{ij} - \hat{\mu}_{ij}^{(I)}}{\sqrt{V(\hat{\mu}_{ij}^{(I)})}}$$
- ▲ Pokud $\hat{\mu}_{ij}^{(I)}$ správně odhadují střední hodnotu $\mathbb{E}Y_{ij}$, potom:
- ▲ $\mathbb{E}[r_{ij}^P] \approx 0$, $\text{var}[r_{ij}^P] \approx \varphi$, $\mathbb{E}[r_{ij}^P, r_{ik}^P] \approx \varphi \text{cor}(Y_{ij}, Y_{ik}) = \varphi (R_{ijk}(\alpha))$
- ▲ Hledáme momentové odhady parametru α spočtené ze součinů Personových reziduí $r_{ij}^P \cdot r_{ik}^P$ y dat ze stejné skupiny.

GEE – volba korelační struktury metoda parametrizované korelace

Příklady volby korelační struktury



▲ Pásová korelace 1. řádu:

$$\text{▲ } R_i(\alpha) = \text{cor}(Y_{ij}, Y_{ik}) = \begin{pmatrix} 1 & \alpha & 0 & \dots & \dots & 0 \\ \alpha & 1 & \alpha & \ddots & & \vdots \\ 0 & \alpha & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \alpha & 0 \\ \vdots & & \ddots & \alpha & 1 & \alpha \\ 0 & \dots & 0 & 0 & \alpha & 1 \end{pmatrix}$$

▲ Konzistentní odhad parametru α :

$$\hat{\alpha} = \frac{1}{\hat{\varphi}} \frac{1}{n-K-p} \sum_{i=1}^K \sum_{j=1}^{n_i-1} r_{ij}^P \cdot r_{i,j+1}^P$$

GEE – volba korelační struktury metoda parametrizované korelace

Příklady volby korelační struktury



▲ Pásová korelace m. řádu: (m=2)

$$\text{▲ } R_i(\alpha) = \text{cor}(Y_{ij}, Y_{ik}) = \begin{pmatrix} 1 & \alpha_1 & \alpha_2 & 0 & \cdots & 0 \\ \alpha_1 & 1 & \alpha_1 & \ddots & & \vdots \\ \alpha_2 & \alpha_1 & \ddots & \ddots & \ddots & 0 \\ 0 & \ddots & \ddots & \ddots & \alpha_1 & \alpha_2 \\ \vdots & & \ddots & \alpha_1 & 1 & \alpha_1 \\ 0 & \cdots & 0 & \alpha_2 & \alpha_1 & 1 \end{pmatrix}$$

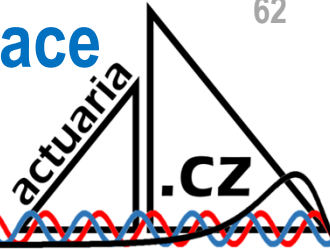
▲ matici rozšíříme na více pásem, m nesmí být moc vysoké

▲ Konzistentní odhad parametru α :

$$\hat{\alpha} = \frac{1}{\hat{\varphi}} \frac{1}{n - K * l - p} \sum_{i=1}^K \sum_{j=1}^{n_i - 1} r_{ij}^P \cdot r_{i,j+1}^P$$

GEE – volba korelační struktury metoda parametrizované korelace

Příklady volby korelační struktury



▲ Exchangable correlation

$$\text{▲ } R_i(\alpha) = \text{cor}(Y_{ij}, Y_{ik}) = \begin{pmatrix} 1 & \alpha & \alpha & \cdots & \cdots & \alpha \\ \alpha & 1 & \alpha & \ddots & & \vdots \\ \alpha & \alpha & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \alpha & \alpha \\ \vdots & & \ddots & \alpha & 1 & \alpha \\ \alpha & \cdots & \alpha & \alpha & \alpha & 1 \end{pmatrix}$$

$$\text{▲ Konzistentní odhad parametru } \alpha : \quad \hat{\alpha} = \frac{1}{\hat{\varphi}} \frac{1}{\dots} \sum_{i=1}^K \sum_{j=1}^{n_i} \Sigma r_{ij}^P \cdot r_i^P$$

GEE – volba korelační struktury metoda parametrizované korelace

Příklady volby korelační struktury

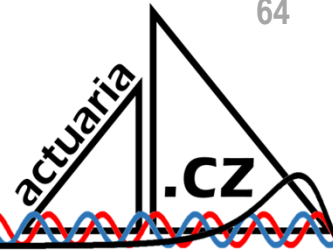
63



▲ AR(1) correlation

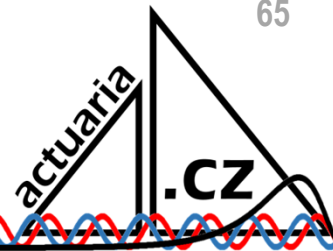
▲ Y_{i1}, \dots, Y_{in_i} jsou z AR(1) časové řady

▲ $cor(Y_{ij}, Y_{ik}) = \alpha^{|j-k|}$



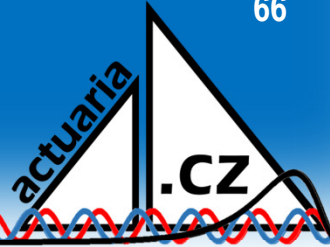
GEE - shrnutí

- ▲ GEE metoda vhodná pro regresní analýzu dat rozdělených do K nezávislých skupin. Uvnitř skupin jsou data vzájemně korelovaná.
- ▲ Není nutně přesně určit rozdělení dat Y , ani jejich rozptyl či korelační strukturu uvnitř skupin.
- ▲ Regresní parametry β mají „population-average“ interpretaci



GEE - shrnutí

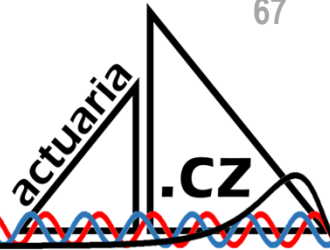
- ▲ GEE metoda není založena na věrohodnosti, pouze na kvazi-věrohodnosti. Pro analýzu modelu proto nelze použít statistiky na bázi věrohodnosti (tj. Deviance, AIC, BIC).
- ▲ Existují analogické statistiky na bázi kvazi-věrohodnosti (QIC (Quasilikelihood under the Independence model Information Criterion), QIC_{HH} , CIC (Correlation Information Criterion), CIC_{HH}), více viz Hudecová&Pešta, , podrobněji Hardin&Hilbe. Je rovněž možné použít skorové testy. (dj-h)
- ▲ GEE – ilustrace
 - ▲ Příklad: Vehicle Insurance Claims – DeJong, Heller: Logistická regrese s korelovanými pozorováními, longitudinální data, navrhovaná korelační struktura AR(1).
- ▲ Využití GEE metod pro rezervování škod. Možnost modelování závislosti vývoje škodních trojúhelníků v jednotlivých letech – viz. Hudecová&Pešta, Gerthofer



Lineární smíšené modely Linear Mixed (Effect) Models

LMM

LLM



▲ Odezvy Y_1, \dots, Y_K splňují **jednourovňový model LMM**, pokud platí:

▲ $Y_i = X_i \beta_i + Z_i b_i + \varepsilon_i, i = 1, \dots, K$

▲ b_i (náhodné efekty)...nezávislé vektory t.ž. $b_i \sim N_q(\mathbf{0}, \mathbb{D}),$

▲ ε_i ... náhodné vektory, t.ž.: $\varepsilon_i \sim N_{n_i}(\mathbf{0}, \sigma_e^2 \mathbb{I}_{n_i}),$

▲ X_i ... regresní matice pro pevné efekty β_i

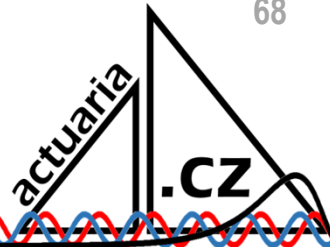
▲ Z_i ... regresní matice pro náhodné efekty b_i

▲ pevná složka: $X_i \beta_i$

▲ náhodná složka: $Z_i b_i + \varepsilon_i, \mathbb{E}[Z_i b_i + \varepsilon_i] = \mathbf{0}, \text{var}[Z_i b_i + \varepsilon_i] = Z_i \mathbb{D} Z_i^T + \sigma_e^2 \mathbb{I}_{n_i}$

▲ (neznámé) parametry modelu MLE: β, σ_e^2 a \mathbb{D} (kovariantní matice, symetrická a pozitivně definitní)

▲ Matice \mathbb{D} se nahrazuje maticí Δ , t.ž. $\Delta^T \Delta = \sigma_e^2 \mathbb{D}^{-1}$



▲ Odezvy Y_1, \dots, Y_K splňují **jednoúrovňový model LMM**, pokud platí:

▲ $Y_i = X_i \beta_i + Z_i b_i + \varepsilon_i, i = 1, \dots, K$

▲ b_i (náhodné efekty)...nezávislé vektory t.ž. $b_i \sim N_q(\mathbf{0}, \mathbb{D}),$

▲ ε_i ... náhodné vektory, t.ž.: $\varepsilon_i \sim N_{n_i}(\mathbf{0}, \sigma_e^2 \mathbb{I}_{n_i}),$

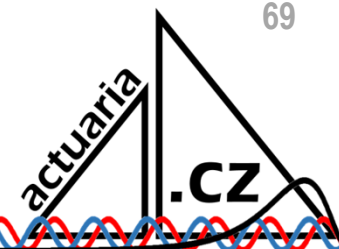
▲ **Marginální tvar** - alternativní zápis modelu LMM: Marginální tvar:

▲ $Y_i \sim N_{n_i}(X_i \beta, \sigma_e^2 \Sigma_i),$ kde $\Sigma_i = Z_i \mathbb{D} Z_i^T / \sigma_e^2 + \mathbb{I}_{n_i}$

▲ nebo

▲ $Y \sim N_n(X \beta, \sigma_e^2 \Sigma)$

LLM – odhady parametrů



▲ Marginální tvar modelu LMM:

▲ $\mathbf{Y}_i \sim N_{n_i}(\mathbb{X}_i \boldsymbol{\beta}, \sigma_e^2 \boldsymbol{\Sigma}_i),$

▲ kde $\boldsymbol{\Sigma}_i = \mathbb{Z}_i \mathbb{D} \mathbb{Z}_i^T / \sigma_e^2 + \mathbb{I}_{n_i}$ modelu je funkcí parametrů $\boldsymbol{\beta}$ a θ .

▲ 1. Marginální věrohodnost:

▲ Logaritmická věrohodností funkce: $l_n(\boldsymbol{\beta}, \theta, \sigma_e^2)$

▲ Minimalizujeme $l_n(\boldsymbol{\beta}, \theta, \sigma_e^2)$ vzhledem k $\boldsymbol{\beta}$ pro dané θ :

▲ řešení metodou vážených nejmenších čtverců

▲ řešení: $\hat{\boldsymbol{\beta}}(\theta) = (\sum_{i=1}^K \mathbb{X}_i^T \boldsymbol{\Sigma}^{-1} \mathbb{X}_i)^{-1} \sum_{i=1}^K \mathbb{X}_i^T \boldsymbol{\Sigma}^{-1} \mathbf{Y}_i$

▲ Minimalizujeme $l_n(\hat{\boldsymbol{\beta}}(\theta), \theta, \sigma_e^2)$ vzhledem k σ_e^2 pro dané θ :

▲ řešení: $\hat{\sigma}_e^2(\theta) = \frac{1}{n} \sum_{i=1}^K (\mathbf{Y}_i - \mathbb{X}_i \hat{\boldsymbol{\beta}}(\theta))^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y}_i - \mathbb{X}_i \hat{\boldsymbol{\beta}}(\theta))$

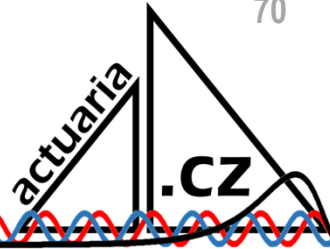
▲ Maximalizujeme profilovou věrohodnost: $l_n^*(\theta) = l_n(\hat{\boldsymbol{\beta}}(\theta), \theta, \hat{\sigma}_e^2(\theta))$ přes θ :

▲ řešení: $\hat{\theta}$

▲ Toto je však těžko řešitelné, aplikujeme jiný přístup využívající strukturu náhodných efektů a dekompozici věrohodnosti

LLM – odhady parametrů

Hendersonovy rovnice pro smíšený model



▲ Obecný tvar pro smíšený model:

$$\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \mathbb{Z}\mathbf{b} + \boldsymbol{\varepsilon}$$

$$\text{▲ kde } \mathbf{b} \sim N_{q^*}(\mathbf{0}, \mathbb{D}_*), \boldsymbol{\varepsilon} \sim N_n(\mathbf{0}, \boldsymbol{\Lambda}_*), \text{cov}(\mathbf{b}, \boldsymbol{\varepsilon}) = 0, \boldsymbol{\Sigma} = \text{var}(\mathbf{Y}) = \mathbb{Z}\mathbb{D}_*\mathbb{Z}^T + \boldsymbol{\Lambda}_*.$$

▲ jednoúrovňový model LMM je speciálním případem. $q^* = Kq$, $n = \sum_{i=1}^K n_i, \dots, \boldsymbol{\Lambda}_* = \sigma_e^2 \mathbb{I}_n$

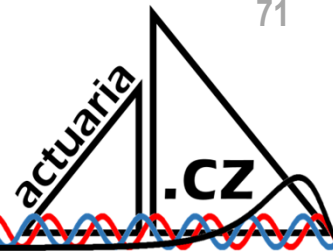
▲ Sdruženou hustotu $(Y, b) = f(y, b; \boldsymbol{\beta}) = f(y|b; \boldsymbol{\beta})f(b)$ považujeme za věrohodnostní funkci neznámých parametrů $(\boldsymbol{\beta}, b)$, maximalizujeme současně přes $\boldsymbol{\beta}$ i b , odhady $(\hat{\boldsymbol{\beta}}, \hat{b})$, jako řešení Hendersonových rovnic pro smíšený model

$$\text{▲ } \hat{\boldsymbol{\beta}} = (\mathbb{X}^T \boldsymbol{\Sigma}^{-1} \mathbb{X})^{-1} (\mathbb{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{Y})$$

$$\text{▲ } \hat{b} = \mathbb{D}_* \mathbb{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbb{X} \hat{\boldsymbol{\beta}})$$

LLM – odhady parametrů

Hendersonovy rovnice pro smíšený model



▲ Odhady regresních parametrů pomocí Hendersonových rovnic pro smíšený model:

$$\blacktriangle \hat{\beta} = (\mathbb{X}^T \Sigma^{-1} \mathbb{X})^{-1} (\mathbb{X}^T \Sigma^{-1} Y)$$

$$\blacktriangle \hat{b} = \mathbb{D}_* Z^T \Sigma^{-1} (Y - \mathbb{X} \hat{\beta})$$

▲ Pro jednoúrovňový LMM model dostáváme

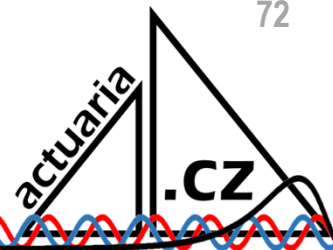
$$\blacktriangle \hat{b} = \mathbb{D} Z_i^T (Z_i \mathbb{D} Z_i^T + \sigma_e^2 \mathbb{I}_{n_i})^{-1} (Y_i - \mathbb{X}_i \hat{\beta})$$

▲ $\hat{\beta}$ je BLUE (nejlepší lineární nestranný odhad) a konzistentní odhad parametru β bez ohledu na rozdělení Y

▲ \hat{b} je BLUP (nejlepší lineární nestranný prediktor) pro b .

LLM – odhady parametrů

Hendersonovy rovnice pro smíšený model



- ▲ Odhady rozptylů regresních parametrů pomocí Hendersonových rovnic pro jednoúrovňový smíšený model:

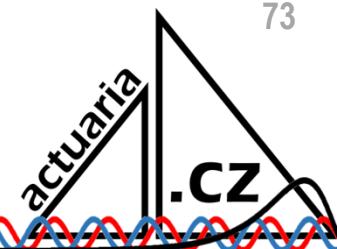
- ▲ $var(\hat{\beta}) = \left(\sum_{i=1}^K \mathbb{X}_i^T \boldsymbol{\Sigma}_i^{-1} \mathbb{X}_i \right)^{-1}$

- ▲ $var(\hat{b}_i) = \mathbb{D} \mathbb{Z}_i^T \left[\boldsymbol{\Sigma}_i^{-1} - \boldsymbol{\Sigma}_i^{-1} \mathbb{X}_i \left(\sum_{i=1}^K \mathbb{X}_i^T \boldsymbol{\Sigma}_i^{-1} \mathbb{X}_i \right)^{-1} \mathbb{X}_i^T \boldsymbol{\Sigma}_i^{-1} \right] \mathbb{Z}_i \mathbb{D}$

- ▲ $var(\hat{b}_i - \hat{b}) = \mathbb{D} - var(\hat{b}_i)$

- ▲ Odhady rozptylů regresních parametrů lze řešit i pomocí REML metody (Restricted Maximum Likelihood Estimators), výsledky se mírně liší od věrohodnostních odhadů.
- ▲ Pro velké počty nezávislých skupin K je rozdíl zanedbatelný. Asymptoticky jsou výsledky shodné.

GEE x LMM



▲ LMM:

- ▲ poskytuje detailní model pro $\text{var } Y$
- ▲ Předpokládá Normalitu b_i a ε_i
- ▲ Poskytuje informaci o struktuře rozptylu
 - ▲ Dekompozice
 - ▲ Odhady složek rozptylu
 - ▲ Náhodné efekty
 - ▲ Testy hypotéz o struktuře rozptylu
- ▲ Pokud struktura rozptylu NENÍ dobře specifikovaná, závěry o pevných efektech β (testy, intervaly spolehlivosti) jsou NEPLATNÉ
- ▲ Vyžaduje velké K (počet subjektů, nezávislých skupin)

▲ GEE:

- ▲ používá „pracovní model“ pro $\text{var } Y$, o kterém se však nepředpokládá, že je správný
- ▲ Neklade žádný předpoklad o rozdělení Y
- ▲ Neposkytuje dostatek informací o struktuře rozptylu
- ▲ Pokud je K (počet subjektů, nezávislých skupin) dost velké, závěry o β jsou platné i když pracovní struktura rozptylu není správná
- ▲ Vyžaduje velké K (počet subjektů, nezávislých skupin), pro malé selhává

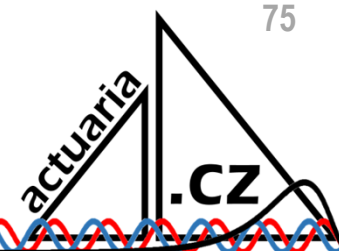


Zobecněné lineární smíšené modely

Generalised Linear Mixed Models

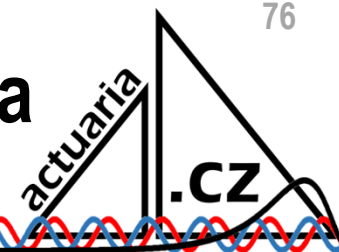
GLMM

GLMM



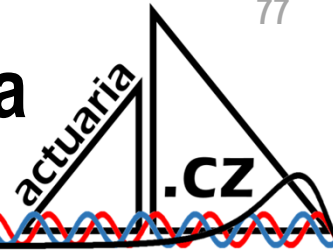
- ▲ GLMM - další možnost jak analyzovat data, u nichž je **porušen předpoklad nezávislosti / nekorelovanosti**.
- ▲ GLMM lze pohlížet jako na **zobecnění Lineárních smíšených modelů (LMM)**, kdy rozšíříme skupinu rozdělení, z nichž pochází odezva/vysvětlovaná proměnná Y .
- ▲ GLMM lze pohlížet jako na **zobecnění Zobecněných lineárních modelů GLM**, kdy do modelu kromě pevných regresních parametrů β (fixed effects, pevné efekty), **zavedeme další prvek náhodnosti pomocí náhodných efektů b** (random effects).

GLMM - DEFINICE MODELU pro skupinově závislá data



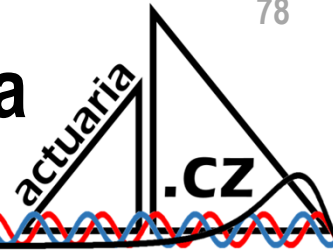
- ▲ 1. **Data:** $\mathbf{Y}_1, \dots, \mathbf{Y}_K$ nezávislé náhodné vektory $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^T, i = 1 \dots K, \sum_{i=1}^K n_i = n$
- ▲ 2. náhodné efekty: $\mathbf{b}_i \dots$ nezávislé q -rozměrné vektory s hustotou $h(\mathbf{b}; \psi)$
- ▲ Složky Y_{i1}, \dots, Y_{in_i} vektoru \mathbf{Y}_i jsou při daném \mathbf{b}_i podmíněně nezávislé, podmíněná hustota $\in EDF$
- ▲ $f(\mathbf{y}|\mathbf{b}_i) = \exp\left\{\frac{y\theta_{ij}-b(\theta_{ij})}{\varphi} + c(\mathbf{y}, \varphi)\right\}$
- ▲ 3. $\theta_{ij} \dots$ kanonický parametr, závisí na:
 - ▲ (p -rozměrných) regresorech pro pevné efekty... X_{ij}
 - ▲ pevných regresních parametrech $\beta \in \mathbb{R}$
 - ▲ (q -rozměrných) regresorech pro náhodné efekty... $Z_{ij}, q \leq p$
 - ▲ náhodných efektech \mathbf{b}_i prostřednictvím
 - ▲ Lineárního prediktoru $\eta_{ij} = \mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i$
- ▲ 4. $g \dots$ linková funkce, platí $g(\mu_{ij}) = \eta_{ij}$

GLMM - DEFINICE MODELU pro skupinově závislá data



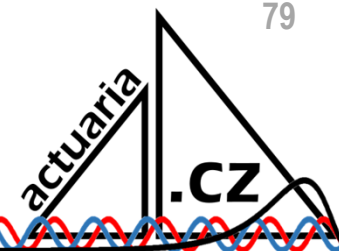
- ▲ Data: $\mathbf{Y}_1, \dots, \mathbf{Y}_K$ nezávislé náhodné vektory $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^T, i = 1 \dots K, \sum_{i=1}^K n_i = n$
- ▲ Náhodné efekty: $\mathbf{b}_i \dots$ nezávislé q -rozměrné vektory s hustotou $h(\mathbf{b}; \psi)$
- ▲ Podmíněná hustota $f(y|\mathbf{b}_i) = \exp \left\{ \frac{y\theta_{ij} - b(\theta_{ij})}{\varphi} + c(y, \varphi) \right\}$
- ▲ Lineární prediktor $\eta_{ij} = \mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i$
- ▲ Linková funkce $g(\mu_{ij}) = \eta_{ij}$
- ▲ $\mathbb{E}[Y_{ij} | \mathbf{b}_i] = \mu_{ij} = b'(\theta_{ij}) = g^{-1}(\eta_{ij}) = g^{-1}(\mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i)$
- ▲ $\text{var}[Y_{ij} | \mathbf{b}_i] = \varphi b''(\theta_{ij}) = \varphi V(\mu_{ij})$

GLMM - DEFINICE MODELU pro skupinově závislá data



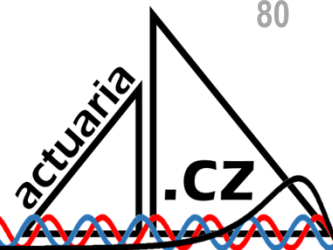
- ▲ Při daném \mathbf{b}_i , Y_{ij} podmíněně splňují **GLM** model,
 - ▲ podmíněná hustota $f(y|\mathbf{b}_i) \in EDF$
- ▲ stejně jako **LMM** modelu, má zavedení náhodných efektů \mathbf{b}_i do všech lineárních prediktorů za následek korelaci mezi Y_{i1}, \dots, Y_{in_i} .
- ▲ $\mathbb{E}[Y_{ij}|\mathbf{b}_i] = g^{-1}(\mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i)$
- ▲ $\mathbb{E}[Y_{ij}] = \mathbb{E}[\mathbb{E}[Y_{ij}|\mathbf{b}_i]] = \mathbb{E}[g^{-1}(\mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i)]$

GLMM – odhady parametrů



- ▲ Vzhledem k předpokladu podmíněného rozdělení odezvy/závislých proměnných, je možné pro odhady parametrů použít metodu maximální věrohodnosti
- ▲ Věrohodností funkce $l(\beta, \varphi | Y)$ však nemá explicitní řešení
- ▲ Je nutné řešit aproximačními metodami
- ▲ jedním z možných přístupů je použití Laplaceovy aproximace, která zjednoduší tvar věrohodnostní funkce a umožní odhady parametrů následujícím způsobem:
 - ▲ Odhad β : modifikovaný IWLS algoritmus
 - ▲ Odhad \mathbf{b}_i : Hendersonovy rovnice
 - ▲ Odhad φ a ϕ : momentová metoda

GLMM



▲ podmíněný model:

- ▲ $\mathbb{E}[Y_{ij}|b_i] = g^{-1}(\mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i)$
- ▲ $\boldsymbol{\beta}$ vyjadřuje podmíněný efekt \mathbf{X}_{ij} na $\mathbb{E}[Y_{ij}]$ vzhledem k \mathbf{b}_i
- ▲ tzv: subject-specific efekty
- ▲ Popisují vliv na $\mathbb{E}[Y_{ij}|b_i]$, když daný subjekt/jedinec mění hodnotu \mathbf{X}_{ij}
- ▲ Parametry $\boldsymbol{\beta}$ obecně nemohou porovnat dva různé subjekty, které se liší v hodnotě \mathbf{X}_{ij}

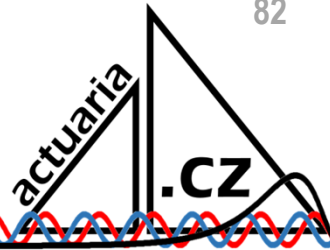
▲ marginální model:

- ▲ $\mathbb{E}[Y_{ij}] = \mathbb{E}[\mathbb{E}[Y_{ij}|b_i]] = \mathbb{E}[g^{-1}(\mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i)]$
- ▲ Nepodmíněný efekt \mathbf{X}_{ij} na $\mathbb{E}[Y_{ij}]$ vzhledem k \mathbf{b}_i
- ▲ tzv: population-average efekty
- ▲ $\mathbb{E}[Y_{ij}]$ obecně nespĺňuje předpoklady GLM, parametry $\boldsymbol{\beta}$ (v podmíněném modelu) obecně nemají population-average interpretaci.



GEE x GLMM

GEE x GLMM



- ▲ GLM analyzuje nezávislá data
- ▲ V případě korelovaných dat GLM nedávají dobré výsledky
- ▲ Řešením může být použití modelů GLMM nebo GEE, která pracují s korelovanými daty
- ▲ Jak zvolit mezi GLMM a GEE?
- ▲ Dle požadované interpretace parametrů

GEE x GLMM



- ▲ GLMM: Conditional Model
- ▲ GEE: Marginal Model

- ▲ GLMM: Subject Specific interpretace:
- ▲ GEE: Population Average interpretace

- ▲ GLMM: regresní koeficienty se vztahují na každého jednotlivce (subjekt), nikoli však nutně na celou populaci,
- ▲ GEE: regresní koeficienty se vztahují na celou populaci, nemusí však platit pro každého jednotlivce

- ▲ GLMM odhaduje jiné parametry než GEE. Pokud oba modely mají stejnou linkovou funkci -> aspoň jeden model není správný (výjimky)
- ▲ Pokud je GLMM dobrý model, obvykle je i GEE dobrý, má však výrazně odlišné odhady parametrů

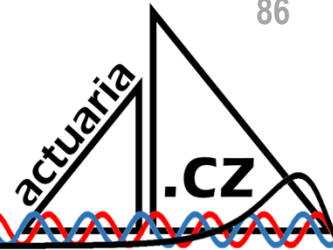
GEE x GLMM



- ▲ GLMM odhaduje jiné parametry než GEE. Pokud oba modely mají stejnou linkovou funkci -> aspoň jeden model není správný (výjimky)
- ▲ Pokud je GLMM dobrý model, obvykle je i GEE dobrý, má však výrazně odlišné odhady parametrů.
 - ▲ log-lineární model: liší se pouze v absolutním členu (intercept)
 - ▲ Logistický model: liší se i směrnice (slope)
- ▲ GLMM: závislost modeluje pomocí náhodných efektů přidaných do GLM modelu
- ▲ GEE: nepředpokládá žádné konkrétní rozdělení odezvy, stejně jako GLM vyžaduje specifikaci: linkové funkce a lineárního prediktoru. Místo konkrétního rozdělení proměnných však stačí specifikovat vztah mezi Střední hodnotou a kovariancí (analogie rozptylové funkce, v případě kvazivěrohodnosti v GLM), resp. Specifikovat „pracovní varianci“
- ▲ GEE: závislostní struktura je modelována pomocí „pracovní“ kovarianční matice, která nemusí odpovídat skutečné závislosti. Doporučuje se používat jednoduchou korelační strukturu



A CO DÁL?



REFERENCE

- ▲ P. de Jong, G. Z. Heller: *Generalized Linear Models for Insurance Data*. Cambridge University Press, 2008.
- ▲ E. Ohlsson, B. Johansson: *Non-Life Insurance Pricing with Generalized Linear Models*. EAA Series, Springer-Verlag Berlin Heidelberg, 2010.
- ▲ W. N. Venables, B. D. Ripley: *Modern Applied Statistics with S*. 4th edition. Springer 2002
- ▲ S. Wood: *Generalized Additive Models, An Introduction with R*, Chapman & Hall/CRC Press, 2006
- ▲ C.E. McCulloch, S.R. Searle: *Generalized, Linear, and Mixed Models*. Wiley Series in Probability and Statistics, Wiley 2001
- ▲ J. W. Hardin, J. M. Hilbe: *Generalized Estimating Equations*. Chapman & Hall/CRC, 2003
- ▲ A. Agresti: *An Introduction to Categorical Data Analysis*. Wiley, 2007
- ▲ M. Branda: *Zobecněné lineární modely v pojišťovnictví*. MFF UK 2013
- ▲ Š. Hudecová, M. Pešta: *Modeling Dependencies in Claims Reserving with GEE*. MFF UK 2003
- ▲ M. Gerthofer: *Claims reserving within the panel data Framework*. MFF UK 2015, diplomová práce
- ▲ Poznámky k přednášce: *Pokročilé regresní modely / Advanced Regression Models (NMST432)*, 2015, MFF UK, přednášející Doc. Mgr. Michal Kulich, Ph.D.
- ▲ Poznámky k přednášce: *Matematika neživotního pojištění (NMFM402)*, 2014, MFF UK, přednášející RNDr. Lucie Mazurová, Ph.D.
- ▲ Poznámky k přednášce: *Vybraný software pro finance a pojišťovnictví / Selected Software Tools for Finance and Insurance (NMFM404)*, 2014, MFF UK, přednášející RNDr. Michal Pešta, Ph.D.
- ▲ Poznámky ze semináře České společnosti aktuárů: *Zobecněné lineární modely (GLM) v pojišťovnictví*, 2012, přednášející Ing. Pavel Zimmermann, Ph.D.
- ▲ Poznámky ze semináře České společnosti aktuárů: *Aplikované modely storen*, 2015,
- ▲ <https://onlinecourses.science.psu.edu/statprogram/stat504> - Analysis of Discrete Data