

O STATISTICKÉ INFERENCI VE WEIBULLOVĚ ROZDĚLENÍ

Jan Hurt

Matematicko-fyzikální fakulta Univesity Karlovy
Sokolovská 86, 186 00 Praha 8

1. ÚVOD

Weibullovo rozdělení jistým způsobem zobecňuje exponenciální rozdělení. Weibullovo rozdělení s parametry $\theta > 0, \beta > 0$, které budeme označovat $W(\theta, \beta)$ má distribuční funkci

$$(1.1) \quad F(x) = \begin{cases} 1 - \exp(-(x/\theta)^\beta), & x > 0 \\ 0 & \text{jinak} \end{cases}$$

a hustotu

$$(1.2) \quad f(x) = \begin{cases} (\beta/\theta)(x/\theta)^{\beta-1} \exp(-(x/\theta)^\beta), & x > 0 \\ 0 & \text{jinak.} \end{cases}$$

Pro $\beta=1$ je tedy Weibullovo rozdělení totožné s exponenciálním rozdělením s parametrem θ . Weibullovo rozdělení $W(\theta, \beta)$ má střední hodnotu

$$(1.3) \quad \mu = \theta \Gamma(1+1/\beta)$$

a rozptyl

$$(1.4) \quad \sigma^2 = \theta^2 (\Gamma(1+2/\beta) - \Gamma^2(1+1/\beta)),$$

takže variační koeficient

$$(1.5) \quad v = \left(\frac{\Gamma(1+2/\beta)}{\Gamma^2(1+1/\beta)} - 1 \right)^{1/2}$$

je pouze funkcií β . Tohoto faktu se využívá pro konstrukci odhadu parametru β momentovou metodou.

Snad nejčastěji se Weibullovo rozdělení používá v teorii spolehlivosti. Odpovídající intenzita poruch

$$(1.6) \quad r(x) = \frac{f(x)}{1 - F(x)} = (\beta/\theta)(x/\theta)^{\beta-1}, \quad x > 0,$$

je totiž v závislosti na parametru β klesající ($\beta < 1$), konstantní ($\beta=1$) nebo rostoucí ($\beta > 1$). Přitom pro $\beta < 2$ je $r(x)$ konkávní a pro $\beta > 2$ konvexní. Ukazuje se, že díky těmto skutečnostem lze pomocí Weibullovova rozdělení dobře approximovat doby do poruchy (doby života) prvků (systémů) s monotónní intenzitou poruch. Na základě teoretických úvah je též možné dojít k Weibullově rozdělení jakožto limitnímu rozdělení (teorie extrémních hodnot).

Poznamenejme, že někdy je užitečná následující reparametrizace (1.1): položíme-li $\gamma = \ln \theta$, $\xi = 1/\beta$ a $Y = \ln X$, kde X má rozdělení $W(\theta, \beta)$, má Y rozdělení s distribuční funkcí

$$(1.7) \quad G(y) = 1 - \exp(-\exp(\frac{y-\gamma}{\xi})), \quad -\infty < y < \infty.$$

Rozdělení s touto distribuční funkcí je asymptotické rozdělení extrémních hodnot typu I. Parametry γ a ξ zde mají význam parametrů polohy a měřítka, tj., distribuční funkci $G(y)$ lze vyjádřit ve tvaru $G(y) = H((y-\gamma)/\xi)$, kde H je distribuční funkce nezávislá na parametrech, v našem případě $H(z) = 1 - \exp(-\exp(z))$. Tohoto faktu se využívá při odhadu neznámých parametrů polohy a měřítka.

2. CENZOROVÁNÍ

V aplikacích je často potřebné provádět statistické závěry na základě neúplných výběrů. Zmíníme se proto o prakticky nejdůležitějších uspořádání experimentů poskytujících tzv. cenzorované výběry. Předpokládejme, že náhodná veličina X je doba do poruchy sledovaného prvku. V čase $t = 0$ začneme pozorovat n prvků stejněho typu. Sledované prvky se časem porouchávají. V případě, že pozorování provádíme až do doby, kdy se všechny prvky porouchají, získáme úplný výběr X_1, \dots, X_n dobu do poruchy. Z časových či ekonomických důvodů však může být nereálné konat pozorování až do okamžiku selhání posledního sledovaného prvku a proto se experiment ukončí podle předem daného pravidla dříve.

Je-li předepsán časový okamžik T , ve kterém je experiment ukončen (bez ohledu na to, ke kolika poruchám skutečně došlo), jsou takto získaná pozorování časově cenzorovaná, přičemž toto cenzorování se nazývá cenzorování typu I. Druhý možný způsob cenzorování zvaný cenzorování typu II spočívá v tom, že experiment je ukončen v okamžiku, kdy se porouchá r -tý sledovaný prvek, kde r je předem dané číslo. Data jsou tzv. cenzorovaná poruchou.

Zatímco při cenzorování typu I je počet skutečně pozorovaných dob do poruchy náhodná veličina a doba trvání experimentu předem dána, je při cenzorování typu II počet pozorovaných poruch předem dán a doba trvání experimentu je dána dobou do poruchy r -tého porouchaného prvku a je tedy náhodnou veličinou.

Při sledování složitých systémů se můžeme setkat s případem, kdy se časové cenzory systém od systému liší, a kdy ukončení pozorování u každého systému je více méně náhodné. Skončí-li pozorování systémů $1, \dots, n$ v okamžicích T_1, \dots, T_n , je ekonomicky neúnosné považovat získaná pozorování jakožto cenzorovaná v okamžiku $T = \min(T_1, \dots, T_n)$ a aplikovat metody cenzorování typu I. Místo toho je vhodné považovat časový censor též za náhodnou veličinu. V takovém případě mluvíme o náhodném cenzorování.

Mimo uvedené typy uspořádání experimentu se vyskytuje i další, často vnějšími podmínkami vynucená uspořádání.

3. METODY ODHADU

Omezíme se pouze na prakticky nejdůležitější metody, a to metodu maximální věrohodnosti a metodu nejmenších čtverců. Předpokládejme, že X_1, \dots, X_n je náhodný výběr z rozdělení s distribuční funkcí (1.1). Logaritmus sdružené hustoty výběru (věrohodnostní funkce) je

$$(3.1) \quad n \ln \beta - n \beta \ln \theta + (\beta - 1) \sum_1^n \ln X_i - \sum_1^n (X_i / \theta)^\beta .$$

Po derivování podle β a θ obdržíme soustavu věrohodnostních rovnic

$$(3.2) \quad n/\beta + \sum \ln X_i + \ln \theta (\theta^{-\beta} \sum X_i^\beta - n) - \theta^{-\beta} \sum X_i^\beta \ln X_i = 0,$$

$$(3.3) \quad \theta = (n^{-1} \sum X_i^\beta)^{1/\beta} .$$

Dosazením (3.3) do (3.2) dostaneme

$$(3.4) \quad \frac{\sum_i^{\beta} \ln X_i}{\sum_i^{\beta}} - \frac{1}{\beta} - \frac{1}{n} \sum \ln X_i = 0.$$

Řešení věrohodnostních rovnic pak spočívá v tom, že iteracním postupem získáme řešení $\hat{\beta}$ rovnice (3.4) a to dosadíme do (3.3). Velmi účinnou metodou pro řešení rovnice (3.4) se ukazuje Newton-Raphsonova metoda. Numerické experimenty ukazují, že konvergence metody je velmi rychlá, a že příliš nezávisí na zvoleném počátečním řešení. Je proto možné zvolit počáteční řešení rozumným odhadem. Dříve se doporučovalo užít za počáteční řešení odhad získaný momentovou metodou na základě vyjádření variačního koeficientu (1.5). Numerický proces je však náročný na přesnost - doporučuje se proto provádět výpočty ve dvojnásobné aritmetice.

Jednoduchá je adaptace věrohodnostních rovnic v případě cenzorovaných pozorování. Výsledkem experimentu při cenzorování typu I je prvních r hodnot pořádkových statistik $X_{(1)} \leq \dots \leq X_{(r)} \leq T$ a informace $X_{(r+1)} > T$. Je-li F distribuční funkce základního souboru a f odpovídající hustota, je věrohodnostní funkce

$$(3.5) \quad \frac{n!}{(n-r)!} \prod_1^r f(X_{(i)}) (1 - F(T))^{n-r}.$$

Konkrétně po zlogaritmování (3.5) a po úpravě dostaneme pro Weibullovo rozdělení soustavu věrohodnostních rovnic

$$(3.6) \quad \frac{\sum_1^r X_{(i)}^{\beta} \ln X_{(i)} + (n-r)T^{\beta} \ln T}{\sum_1^r X_{(i)}^{\beta} + (n-r)T^{\beta}} - \frac{1}{\beta} - \frac{1}{r} \sum_1^r \ln X_{(i)} = 0$$

$$(3.7) \quad \Theta = (r^{-1} (\sum_1^r X_{(i)}^{\beta} + (n-r)T^{\beta}))^{1/\beta},$$

které se řeší analogicky jako rovnice pro úplný výběr. Při cenzorování typu II je výsledkem experimentu prvních r hodnot pořádkových statistik $X_{(1)} \leq \dots \leq X_{(r)}$. Odpovídající věrohodnostní funkce je

$$(3.8) \quad \frac{n!}{(n-r)!} \prod_1^r f(X_{(i)}) (1 - F(X_{(r)}))^{n-r}.$$

Věrohodnostní rovnice pro cenzorování typu II jsou tedy zcela analogické rovnicím pro cenzorování typu I s tím, že místo T píšeme $X_{(r)}$.

V případě úplných výběrů je možno provádět statistické závěry na základě známých asymptotických vlastností maximálně věrohodných odhadů. Asymptotické rozptyly a kovariance maximálně věrohodných odhadů parametrů Weibullova rozdělení v případě cenzorování typu II odvodili HARTER a MOORE. Odvodit vlastnosti (byť i asymptotické) maximálně věrohodných odhadů v případě cenzorování typu I tradičními metodami zdá se doposud být neschůdné. Často se můžeme setkat s tvrzením, že výsledky platné pro cenzorování typu II je možné použít i pro cenzorování typu I. Z intuitivního hlediska se však toto tvrzení zdá být nesprávné proto, že ignoruje informaci $X_{(r+1)} > T$. Numerická

Časový censor T	1	2
Počet simulací, pro něž $r=21$	340	375
MC průměr $\bar{\theta}$	0.705	1.333
Výběrový rozptyl $\hat{\theta}$	0.0061	0.0242
MC průměr $\bar{\beta}$	0.5452	0.5044
Výběrový rozptyl $\hat{\beta}$	0.0127	0.0084

Tabulka

ilustrace ukazuje, že tomu tak skutečně je. V jedné studii metodami Monte Carlo (MC) byly studovány odhady parametrů Weibullova rozdělení s parametry $\theta = 1$, $\beta = 0.5$, rozsah základního souboru $n=30$ pro cenzorování typu I s časovými cenzory $T=1$ a $T=2$. Pro každý censor bylo provedeno $N=3000$ simulací náhodných výběrů, řešeny věrohodnostní rovnice a výsledky registrovány. V tabulce jsou uvedeny výsledky pro marginální rozdělení s $r=21$. Je patrné, že statistická indukce by v obou případech byla velice rozdílná.

Metodu nejmenších čtverců pro odhad parametrů polohy a měřítka je možné aplikovat na Weibullovu rozdělení po reparametrizaci (1.7), tj. převedením na model extrémních hodnot typu I. Je-li $X_{(1)}, \dots, X_{(n)}$ uspořádaný náhodný výběr z rozdělení s distribuční funkcí $G(x) = H((x-\gamma)/\xi)$, nezávisí rozdělení

$$(3.9) \quad Y_{(i)} = (X_{(i)} - \gamma)/\xi, \quad i=1, \dots, n$$

na neznámých parametrech. Tato skutečnost inspiruje k sestrojení regresního modelu

$$(3.10) \quad X_{(i)} = \gamma + \xi EY_{(i)} + u_{(i)}, \quad i=1, \dots, n,$$

kde $u_{(i)}$ jsou poruchy s nulovými středními hodnotami a kovariancemi

$$(3.11) \quad \text{cov}(u_{(i)}, u_{(j)}) = \xi^2 \text{cov}(Y_{(i)}, Y_{(j)}).$$

Kovariance na pravé straně (3.11) nezávisí na neznámých parametrech. Výhodou této formulace je, že z n rovnic (3.10) můžeme využít k odhadu jen některé. Parametry γ a ξ nyní můžeme odhadnout jednak obyčejnou metodou nejmenších čtverců, tj. za předpokladu nekorelovanosti poruch, jednak zobecněnou metodou nejmenších čtverců s využitím znalosti kovariancí $\text{cov}(Y_{(i)}, Y_{(j)})$ pro konkrétní rozdělení. V obou případech jsou odhady ve tvaru lineárních kombinací pořádkových statistik, přičemž koeficienty těchto kombinací jsou pro nejužívanější rozdělení tabelovány. Zajímavou a pravděpodobně dosud neřešenou otázkou je robustnost odhadů obyčejnou metodou nejmenších čtverců.

Odvodit jiné než asymptotické vlastnosti odhadů získaných výše uvedenými metodami je prakticky neproveditelné. Proto i konstrukce konfidenčních intervalů a testování hypotéz se provádí na základě asymptotických vlastností.

Existují i další, obvykle méně efektivní metody odhadu (zjednodušené lineární, kvantilové aj.).

4. LITERATURA

BURY,K.V.: Statistical models in applied science, Wiley, New York, 1975.

JOHNSON,N.L.-KOTZ,S.: Continuous univariate distributions-1, Houghton Mifflin, Boston, 1970.

MANN,N.R.-SCHAFFER,R.E.-SINGPURWALLA,N.D.: Methods for statistical analysis of reliability and life data, Wiley, New York, 1974.