

Pokusíme se řešit úlohu odhadu regresní funkce v případě, že data máme s náhodným cenzorováním. S takovými daty se můžeme setkat při zkoumání životnosti či doby besporuchového provozu. Zmíníme se i o nenáhodném cenzorování, i o Coxově modelu popisujícím regresi jiným způsobem.

1. Model: Představme si náhodné veličiny Y_1, Y_2, \dots, Y_N odpovídající modelu

$$Y_i = r(x_i) + \varepsilon_i, \quad i=1, \dots, N,$$

kde x_i jsou známé hodnoty, (pro jednoduchost zůstaneme v R_1^1), ε_i jsou nezávislé, stejně rozložené náhodné veličiny, centrované, se spojitou distribuční funkcí F , $r(\cdot)$ je reálná funkce, buď neznámá, nebo známého typu s neznámými parametry $\theta \in R_m$. V praxi takový model může být vhodný pro logaritmus doby besporuchového provozu.

Náhodné cenzorování zprava spočívá v tom, že necenzorujeme hodnoty Y_i , ale

$$T_i = \min(Y_i, V_i) \quad \text{a} \quad \delta_i = I[Y_i \leq V_i],$$

kde $I[A]$ označuje indikátor množiny A a V_1, \dots, V_N jsou nějaké náhodné veličiny, které mohou také záviset na hodnotách regresoru x_i . Předpokládáme, že tato závislost je "rozumná", třeba že je jí možné také vyjádřit ve formě modelu

$$V_i = v(x_i) + \varepsilon_i,$$

s náhodnými veličinami ε_i nezávislými navzájem i na $\{\varepsilon_j\}$.

2. Odhad parametrů: Můžeme být v situaci, kdy známe tvar regresní funkce r_θ a potřebujeme odhadnout hodnotu parametru θ . Zcela uspokojivý postup jsem zatím nenašel. Miller (1975) a později Buckley a James (1979) se snažili dospět k odhadu s nejmenšími čtverci iteračním postupem. Je navržen v několika modifikacích a výsledky v předvedených simulovaných příkladech nejsou špatné. Ale postupy jsou početně složité, navíc iterace mohou někdy začít oscilovat. Ani konzistence není zaručena.

Šikovně je odvozena metoda odhadu v práci Kouřil et al. (1981), pro případ veličin V_i nezávislých na x_i a stejně rozložených se spojitou distribuční funkcí G . Te pak můžeme konsistentně odhadnout, při označení p_i pořadí hodnoty T_i .

$$\hat{\theta}(t) = 1 - \prod_{i=1}^N \left(\frac{N-p_i+1}{N-p_i+2} \right)^{(1-\delta_i) \cdot I[T_i \leq t]}.$$

Je to tedy variace na Kaplanův a Meierův (1958) "product limit estimator"-PLE.

Všimneme-li si, že pro pevné x je

$$E_x(\delta_i T_i) = \int_{-\infty}^{\infty} t [1-G(t)] dF[t-r_\theta(x)],$$

tj.
$$E_x \left\{ \frac{\delta_i T_i}{1-G(T_i)} \right\} = \int_{-\infty}^{\infty} t dF[t-r_\theta(x)] = r_\theta(x)$$

a vidíme, že náhodné veličiny $\delta_i T_i / [1-G(T_i)]$ splňují regresní model s toutéž regresní funkcí. Jsou v něm stále ještě nezávislé, ale již různě rozložené a s daleko většími rozptyly. To činí metodu velice citlivou a neuspokojivou, jak ukážou výsledky. V praxi samozřejmě použijeme hodnot odhadu $\hat{\theta}(T_i)$, parametr pak odhadneme metodou nejmenších čtverců s hodnot $W_i = \delta_i T_i / [1-G(T_i)]$. Autofi dokázali středně kvadratickou konsistenci takto získaných odhadů parametrů lineární funkce $r_\theta(x) = a + bx$. Z důležitých předpokladů kromě spojitosti distribučních funkcí uplatňujeme na základní, odlišující náhodné cenzorování od nenáhodného useknutí, že $\sup \{t: G(t) < 1\} \geq \sup \{t: F(t) < 1\}$.

Můžeme jej obhajit, použijeme-li vhodně "useknutého" výběru.

2. Neparametrické odhady: K-odhad hodnoty regresní funkce v daném bodě x je v podstatě vážený průměr z hodnot naměřených v okolních bodech, váhy potlačují vliv měření ve vzdálených bodech x_1 . Zvolíme nejjednodušší formu s "oknem" šířky $2 \cdot d_N$:

$$O_{d_N}(x) = \{s \in R_1 : |x-s| \leq d_N\}.$$

Musíme si ještě všimnout rozmístění bodů x_1 . Potřebujeme, aby ve zvoleném okolí bodu x byl dostatečný počet měření. Předpokládejme, že x_1 jsou realizace náhodné veličiny X definované v (R_1, β_1) , rozložené s hustotou h spojitou v bodě x , $h(x) > 0$. Označíme M_N počet bodů v $O_{d_N}(x)$. Pokud $d_N \rightarrow 0$ a $M_N \cdot d_N \rightarrow \infty$, už máme zajištěno, že aspoň ve středním kvadrátu

$$M_N / (2d_N \cdot N) \rightarrow h(x), \text{ čili } M_N \rightarrow \infty \text{ v pravděpodobnosti.}$$

Volí se běžně tvar $d_N = C \cdot N^{-\Delta}$, $\Delta \in (0, 1)$. Pohledem z druhé strany je volba d_N tak, aby byl zajištěn určitý počet bodů M_N -metoda M nejbližších bodů, např. Stone (1977). To spíše odpovídá přístupu v konkrétní praktické situaci.

Označíme-li $T_{k1} \leq \dots \leq T_{kM}$ ($M = M_N$) seřazené všechny výsledky v bodech $x_1 \in O_{d_N}(x)$, uděláme odhad (PLE) distribuční funkce veličiny $\{Y/x_1 \in O_{d_N}(x)\}$:

$$F_{N,x}^Y(t) = 1 - \prod_{j=1}^M \left(\frac{M-j}{M-j+1} \right)^{\delta_{kj}} \cdot I[T_{kj} \leq t]$$

Tato funkce má skoky $\Delta F_{N,x}^Y$ v bodech T_{kj} s $\delta_{kj} = 1$. Pokládáme $\delta_{kM} = 1$, aby $\sum \Delta F_{N,x}^Y = 1$. V práci Volf (1985) je dokázána pro $M \rightarrow \infty$ P konzistence odhadu

$$r_N(x) = \sum_{i=1}^M T_i \cdot \Delta F_{N,x}^Y(T_i) \cdot I[|T_i| \leq \Delta_N],$$

a to za následujících předpokladů:

1. Volíme $d_N = C \cdot N^{-\Delta}$, $C > 0$, $\Delta \in (\frac{2}{3}, 1)$!?
2. Body x_1 jsou rozmístěny tak, že $M_N \rightarrow \infty$.
3. Funkce r , v jsou spojitě v bodě x .
4. Veličiny e_1 jsou stejně rozděleny, s distribuční funkcí G (to není podstatné) a distribuční funkce F a G jsou stejnoměrně spojitě.
5. Označíme-li $\mathcal{P}_F = \sup\{t: F(t) < 1\}$, tak existují $a, b \in (0, 1)$ tak, že $a \cdot [1 - F(t)]^b \leq 1 - G(t)$ pro každé t z nějakého intervalu $\langle \mathcal{P}, \mathcal{P}_F \rangle$.
6. Δ_N je posloupnost čísel, $\Delta_N \rightarrow \infty$, $\Delta_N \cdot \left[\frac{\log(N \cdot d_N)}{(N \cdot d_N)^2} \right] \frac{1}{2+b} \rightarrow 0$.

Příkré omezení hodnot Δ nevypadá nejlépe, pokusím se ukázat, k čemu bylo potřeba. Označíme-li $\epsilon_1 = \min(\epsilon_1, e_1 + v(x) - r(x))$, dostanu se v důkazu k bodu, kdy vím, že

$$(1) \quad |r_N(x) - r(x) \cdot \sum \Delta F_{N,x}^Y(T_i) - \sum T_i \cdot \Delta F_{N,x}^Y(T_i)| = O(d_N) \text{ součet je přes } i \{T_i \leq \Delta_N$$

Přítom podmínka 5. zaručuje podle Rejts (1982) s.j. stejnoměrnou v t konzistenci odhadu $F_{N,x}$ pro F , konstruovaného jako PLE z hodnot $\{T_{kj}\}_{j=1}^{M_N}$. Podmínka 6. pak zajišťuje silnou konzistenci odhadu střední hodnoty ϵ_1 :

$$\sum T_i \cdot \Delta F_{N,x}^Y(T_i) \rightarrow E \epsilon_1 = 0,$$

to je vlastně vedlejší výsledek pro odhad parametru polohy. Pokud následující pravděpodobnost

$$(2) \quad \Pr \{ T_{k1} \leq T_{k2} \leq \dots \leq T_{kM_N} \text{ a zároveň } \delta_{kj} = I[\epsilon_{kj} \leq e_{kj} + v(x) - r(x)], j=1, \dots, M_N \}$$

roste k jedné s $N \rightarrow \infty$, bude tvrzení dokázáno, jak vidíme z (1), protože pak i

$$\Pr \{ \Delta F_{N,x}^Y(T_{kj}) = \Delta F_{N,x}^Y(T_{kj}), j=1, \dots, M_N \} \rightarrow 1.$$

Přítom T_{kj} se od T_{kj} liší o $r(x) + O(d_N)$ takže k (2) stačí, aby $\Pr \{ \min_{j=1, \dots, M_N} \frac{|T_{kj} - T_{kj}|}{d_N} > \epsilon \} \rightarrow 1$

pro libovolné ϵ . Není těžké ukázat, že při 4. podmínce je pro $\epsilon > 2$ a libovolné ϵ

$$\Pr \left\{ \min_{i \neq j} \frac{|r_{ki} - r_{kj}|}{M_N^{-\alpha}} > K \right\} \rightarrow 1.$$

Pak už stačí vsít $d_N = M_N^{-\alpha}$ při $\alpha > 2$, tj. $M_N^{-1} = C_1 (N^{1-\alpha})^{-\alpha}$, což vymezí λ na interval $(\frac{1}{3}, 1)$.

Výsledky příkladů se simulovanými daty ukazují, že toto omezení asi není nutné, zároveň je však pravda, že asymptotické vlastnosti se zde projevují velice pomalu. Takže i volbou C se dá vhodně určit rozměr okna pro značně velké rozpětí N . Příklady, v nichž bylo χ rozloženo rovnoměrně, ukazují, že asi 10 až několik desítek bodů pro rozsahy $N=50 - 1\ 000$ jsou nejvhodnější.

3. Znovu parametrický případ - a robustnost. Protože zatím v části 2 uvedené metody odhadu parametru regresní funkce r , nejsou uspokojivé, tento případ jsme zkusili řešit následujícími způsoby: V 1. kroku odhadneme neparametricky hodnoty regresní funkce, buď ve všech x_1 , nebo jen ve vybraných bodech. V 2. kroku pak odhadneme parametr θ z takto získaných hodnot $r_N(x_j)$, třeba metodou nejm. čtverců.

Metody robustního odhadování sem můžeme vnést dvěma způsoby. Buď do 1. kroku, kdy počítáme "průměr", můžeme použít některé z technik robustního odhadování, třeba trimování. Druhý krok můžeme řešit s trimováním pomocí regresních kvantilů místo prostou metodou nejm. čtverců. Tak bychom se mohli bránit proti vlivu neočekávaných odchylek v šumech ε_1 . Jenže se ukazuje, že to není vždy nutné. Samotný K -odhad sice není robustní, ale vyrovnává data a odchylky tedy zmenšuje. Přidání robustní metody již nepřinese výrazný pokrok, nejsou-li data příliš nečistá.

Pokud jde o pozorování v odlehlých bodech x_1 , zvláště na krajích oboru χ , mají nadále větší vliv, což je o to horší, že v nich jsou K -odhady počítány z méně bodů (a jen z jedné strany). Těže v 2. kroku by se měly použít běžné metody na oslabení vlivu pozorování v odlehlých bodech na odhad parametru. Stejně tak se musíme zmínit o dalších typických vlastnostech jednoduchých K -odhadů, že totiž podceňují konkávní regresní funkci a přeceňují konvexní.

O neparametrických odhadech regresní funkce je možné se dočíst v tomto sborníku v práci J. Antocha, dále např. v sborníku Gasser a Rosenblatt (1973).

4. Cenzorování nenáhodné. Zvláště v technických aplikacích se setkáváme s daty cenzorovanými v okamžiku jejich sběru, při současných počátcích pozorování. Pak vysoké hodnoty jsou useknuty, říká se tomu cenzorování časem. V jiném případě mohou být useknuty i malé hodnoty, jsou třeba mimo rozsah přístroje. Veličiny Y_1 pak pozorujeme přímo jen v nějakém intervalu $\langle \mathcal{P}_1, \mathcal{P}_2 \rangle$. To nám jistě připomene data usečená náhodně při trimování. Pokud předpokládáme symetricky rozložené šumy ε_1 , nic nám nebrání v parametrizovaném případě výběr symetricky trimovat pomocí regresních kvantilů, a zvolíme tak velké, abychom se zbavili všech cenzorovaných hodnot. Pokud se nemůžeme spolehnout na symetrii odchylek nebo je-li cenzorování příliš veliké, je třeba se uchýlit k regresnímu mediánu. O těchto postupech viz Antoch, Robust 84.

Při neparametrickém odhadování hodnoty regresní funkce můžeme opět odhad v každém bodě počítat jako trimovaný průměr, až jako medián v krajním případě. Stejně se dá postupovat při kombinovaném cenzorování (usečené kraje a ještě náhodné cenzorování mezi nimi).

5. Coxův regresní model. Duchovním otcem je D.R.Cox (1972), i když vlastně rozvíjí model závislosti parametru exponenciálního rozložení. Popisuje vývoj intenzity poruchy (hazard rate), což je veličina definovaná jako $\lambda(t) = -d \ln[1-F(t)]/dt$, když $F(t)$ je distribuční funkce. Předpokládá a modeluje tedy vývoj rozložení veličiny závislosti na x , nikoli vývoj střední hodnoty. V nejjednodušším případě má tvar

$$\lambda^y(t, x) = \lambda_0^y(t) \cdot \exp(\beta x).$$

Je vidět, β vesměs záporná znamená pokles intenzity poruch s růstem x , tedy větší životnost. Dobře se dá tento model začlenit do věrohodnostní funkce, takže je možné dospět k maximálně věrohodným odhadům parametru β a pak i k odhadu kumula-

tivní intenzity poruch definované jako $L(\tau) = \int_0^{\tau} \lambda(t) dt$. Maximalizací parciální věrohodnostní funkce

$$\mathcal{L}(\beta) = \prod_{i=1}^N \left\{ \frac{\exp(\beta \cdot x_i)}{\sum_{j=1}^N \exp(\beta x_j)} \cdot I[T_j \geq T_i] \right\}^{\delta_i}$$

dostaneme odhad $\hat{\beta}$. Nabízí se iterační postup, konečné řešení je silně konzistentní. Potom můžeme také odhadnout

$$\hat{L}_0(t) = \sum_{i=1}^N \frac{\delta_i \cdot I[T_i \leq t]}{\sum_{j=1}^N \exp(\beta x_j) \cdot I[T_j \geq T_i]}$$

Navíc asymptotická normalita odhadů $\hat{\beta}$ umožňuje testovat významnost regrese. Z literatury na toto téma uveďme článek Tsiatis (1981), variantu pro tříděná data včetně ukázky praktického použití v Prentice, Gloeckler (1978) i knihu Kalbfleisch, Prentice (1980), kde je i ukážka programu. Novější verze BMDP také již obsahují tuto proceduru. Někdy jsou obě pojetí regrese ekvivalentní, odpovídá-li log Y lineárnímu modelu s dvojně exponenciálními šumy, odpovídá Y Coxově modelu a je rozloženo Weibullovsky. Výhodou Coxova modelu je, že nevyžaduje žádné specifické vlastnosti od cenzorující náhodné veličiny V; může také záviset na x, může i veličinu Y uřezávat tj. $\sup\{V\} \leq \sup\{Y\}$, i tím spíše odpovídá praktické situaci s kombinovaným cenzorováním.

Literatura:

Antoch J., sborníky Robust 1984 a Robust 1986.

Buckley J., James I. (1979), Linear regression with censored data, Biometrika 66, 429.

Cox D.R. (1972), Regression models and life tables -with discussion, J.R.Statist.Soc. B 34, 187-220.

Gasser Th., Rosenblatt M. -editors (1979), Smoothing Techniques for Curve Estimation, Lecture Notes on Math. No 757, Springer-Verlag.

Kalbfleisch J.D., Prentice R.L. (1980), The Statistical Analysis of Failure Time Data, N.Y., Wiley.

Kaplan E.L., Meier P. (1958), Nonparametric estimation from incomplete observations, J.Amer.Stat.Assoc. 53, 457-81.

Koel H., Susarla V., Van Ryzin J. (1981), Regression analysis with randomly right-censored data, Ann.Stat. 9, 1276-88.

Miller R.G. (1976), Least squares regression with censored data, Biometrika 63, 449-64.

Prentice R.L., Gloeckler L.A. (1978), Regression analysis of grouped data with application to breast cancer data, Biometrics 34, 57-67.

Rejtő L. (1982), On the fixed censoring model and consequences for the stochastic case, 9th Prague Conf. on Inf. Theory ..., Vol B, Academia.

Stone C.J. (1977), Consistent nonparametric regression -with discussion, Ann.Stat. 5, No 4.

Tsiatis A.A. (1981), Large sample study of Cox's regression model, Ann.Stat. 9, 93-108.

Volf P. (1985), Estimation in linear model with censored data, posláno do Proc. of 5th Pannonian Symposium on Math. Stat., Budapest.