

LOGARITMICKO LINEÁRNÍ MODELY PRO POSLOUPNOSTI ZÁVISLÝCH BINÁRNÍCH DAT

Martin Janžura, ÚTIA ČSAV, Praha

S posloupnostmi binárních dat se lze setkat vsude tam, kde jsou v pravidelných časových intervalech opakovaně pozorovány veličiny nabývající pouze dvě hodnoty (např. přítomnost - nepřítomnost příznaku, překročení - nepřekročení kritické úrovně apod.). Tyto hodnoty mohou být zcela libovolné, pro jednoduchost je účelné pracovat s dvoubodovou množinou $\{0,1\}$.

Při statistickém zpracování takových posloupností dat se obvykle využívá teorie markovských řetězců (libovolného řádu). Pravděpodobnosti přechodu však neposkytují dostatečně jemný pohled na skutečnou závislostní strukturu a současně se neshodují, je-li rozsah výběru nedostatečný vzhledem k řádu markovosti.

Z těchto důvodů půjdeme cestou zobecnění teorie logaritmicko lineárních modelů, které se osvědčily pro popis simultánní distribuce konečného systému náhodných veličin (viz Bishop, Fienberg a Holland (1975)). Přírodním zobecněním logaritmicko lineárních modelů pro náhodné posloupnosti obdržíme modely známé v oblasti statistické fyziky pod pojmem gibbovské náhodné posloupnosti.

1. Gibbovské náhodné posloupnosti

Vyděme od definice logaritmicko lineárního modelu pro simultánní rozdělení konečného systému binárních náhodných veličin.

Pro libovolnou množinu V označme $k(V) = \{A \subset V: 0 < |A| < \infty\}$ systém všech konečných neprázdných podmnožin (symbolem $|A|$ značíme kardinalitu množiny A).

Předpokládejme nejprve, že množina V je konečná, tedy $|V| < \infty$. Pravděpodobnostní míra μ^U definovaná na $\{0,1\}^V$ vyhovuje logaritmicko lineárnímu modelu, jestliže existuje systém reálných konstant

$$U = \{U_A\}_{A \in k(V)}$$

tak, že platí

$$(1) \quad \log \mu^U(x_V) = \sum_{A \in k(V)} U_A \prod_{s \in A} x_s - q(U)$$

pro každé $x_V \in \{0,1\}^V$.

kde

$$q(U) = \log \sum_{y_V \in \{0,1\}^V} \exp \left\{ \sum_{A \in k(V)} U_A \prod_{s \in A} y_s \right\}$$

je logaritmus příslušné normovací konstanty.

Ekvivalentní definice má tvar

$$(2) \quad \log \mu^U(x_{\{t\}} | x_{V \setminus \{t\}}) = \\ = \sum_{A \in k(V), A \ni t} U_A \prod_{s \in A} x_s - \log \left[1 + \exp \left\{ \sum_{A \in k(V), A \ni t} U_A \prod_{s \in A \setminus \{t\}} x_s \right\} \right]$$

pro každé $t \in V$, $x_V \in \{0, 1\}^V$.

Hodnoty U_A , $A \in k(V)$ se nazývají interakcemi, celý systém $U = \{U_A\}_{A \in k(V)}$ se ve statistické fyzice nazývá potenciál.

Samotný model je pak určen systémem nenulových interakcí. Pro $A \subset C \subset k(V)$ budeme psát $U \in M_A$ jestliže $U_A = 0$ pro každé $A \in k(V) \setminus A$. Přirazení $A \rightarrow M_A$ je zřejmě monotónní, tj. pro $B \subset A$ platí $M_B \subset M_A$.

Poznámka: V obecné definici logaritmicke lineárního modelu se uvažují libovolné funkce $U_A(x_A)$ namísto zde uvedeného speciálního tvaru $U_A \prod_{s \in A} x_s$. Pak je ale třeba v zájmu jednoznačnosti doplnit definici o nějaké normující vazby, aby zůstalo kýzených $2^{|V|} - 1$ volných parametrů. Tomu se zde užitím speciálního tvaru funkcí vyhneme při zachování plné obecnosti, tj. pro každou kladnou pravděpodobnostní míru μ na $\{0, 1\}^V$ existuje výše uvedené vyjádření ve formě logaritmicke lineárního modelu.

Nahradíme nyní konečnou množinu V prostorem celých čísel Z a budeme uvažovat pouze omezené potenciály

$$U = \{U_A\}_{A \in k(Z)}$$

$$\text{tj. } \sum_{A \in k(Z), A \ni t} |U_A| < \infty \text{ pro každé } t \in Z$$

Rekneme, že distribuce μ^U na $\{0, 1\}^Z$ je gibbovská vzhledem k potenciálu U , jestliže definice (2) platí pro $V = Z$ (zde jsme potřebovali předpoklad omezenosti, aby nekonečné součty konvergovaly).

Definici (1) nelze přímo použít, neboť pro $V = Z$ by neměla dobrý smysl, a kdybychom tímto způsobem definovali systém konečně rozměrných distribucí $\{\mu_V^U\}_{V \in k(Z)}$, nebyl by tento systém nutně konzistentní.

Lze pouze definici rozšířit na systém všech konečně rozměrných podmíněných distribucí

$$\{\mu_V^U | Z \setminus V\}_{V \in k(Z)}$$

předpisem

$$\log \mu_V^U(x_V | x_{Z \setminus V}) = \sum_{A \in k(Z), A \cap V \neq \emptyset} U_A \prod_{s \in A} x_s - q_V^U(x_{Z \setminus V})$$

pro každé $V \in k(Z)$, kde

$$q_V^U(x_{Z \setminus V}) =$$

$$= \log \sum_{y_V \in \{0, 1\}^V} \exp \left\{ \sum_{A \in k(V), A \cap V \neq \emptyset} U_A \prod_{s \in A \cap V} y_s \prod_{s \in A \setminus (Z \setminus V)} x_s \right\}$$

je opět logaritmus příslušné normovací konstanty.

V dalším se budeme zabývat pouze stacionárními potenciály, tj. $U_A = U_{A+t}$ pro každé $A \in k(\mathbb{Z})$, $t \in \mathbb{Z}$. Označíme $k_+(\mathbb{Z}) = \{A \in k(\mathbb{Z}) : \min\{j \in A\} = 0\}$. Model je zde opět určen systémem nenulových interakcí, tj. pro $A \in k_+(\mathbb{Z})$ píšeme $U \in M_A$ jestliže $U_A = 0$ pro $A \in k_+(\mathbb{Z}) \setminus A$.

Jestliže $|A| < \infty$, určuje systém podmíněných distribucí $\{M_V^U | Z \setminus V\}_{V \in k(\mathbb{Z})}$ (který je ovšem jednoznačně dán systémem $\{M_{\{0\}}^U | Z \setminus \{0\}\}$ a jelikož je potenciál stacionární, tak dokonce stačí $M_{\{0\}}^U | Z \setminus \{0\}$ jednoznačným způsobem gibbsovskou distribuci μ^U na $\{0,1\}^{\mathbb{Z}}$. Tato distribuce je pak stacionární a ergodická. (To pro $|A| = \infty$ nemusí být pravda. Nejen že není zaručena jednoznačnost, ale může existovat i nestacionární distribuce.)

Označme $[n] = \{1, \dots, n\}$ pro každé $n \in \mathbb{Z}$. Pak platí

$$\lim_{n \rightarrow \infty} \frac{1}{n} q_U^{[n]}(x_{Z \setminus [n]}) = q(U)$$

stejně jako v $x \in \{0,1\}^{\mathbb{Z}}$ pro každý omezený stacionární potenciál U .

Nechť $|A| < \infty$. Potom q zúženou na M_A lze chápat jako silně konvexní reálnou analytickou funkci na R^A . Označíme-li

$$C_A = \left\{ x \in \{0,1\}^{\mathbb{Z}} : \prod_{i \in A} x_i = 1 \right\} \quad \text{pro } A \in k(\mathbb{Z}),$$

potom zejména platí

$$\nabla q(U) = \left\{ \frac{\partial q}{\partial U_A}(U) \right\}_{A \in A} = \left\{ M^U(C_A) \right\}_{A \in A}$$

•

$$\nabla^2 q(U) =$$

$$= \left\{ \frac{\partial^2 q}{\partial U_A \partial U_B}(U) \right\}_{A, B \in A} = \left\{ \sum_{k=-\infty}^{\infty} \left[M^U(C_A \cap T^{-k} C_B) - M^U(C_A) M^U(C_B) \right] \right\}_{A, B \in A}$$

kde $T : \{0,1\}^{\mathbb{Z}} \rightarrow \{0,1\}^{\mathbb{Z}}$ je posunutí definované předpisem $T(x)_n = x_{n+1}$ pro každé $x \in \{0,1\}^{\mathbb{Z}}$, $n \in \mathbb{Z}$. Přitom $\nabla q(U)$ i $\nabla^2 q(U)$ jsou spojité jako funkce $U \in M_A$ a $\nabla^2 q(U) > 0$ pro každé pevné $U \in M_A$.

Platí také

$$q(U) = r^{-1} \log \lambda_{\max}(Q_U),$$

kde $r = r(A) = \max_{A \in A} \max\{j \in A\}$ a Q_U je tzv. transfer matice definovaná předpisem

$$Q_U(x_{[r]}, z_{[r]}) = \exp \left\{ \sum_{A \in A} U_A \sum_{l=1}^r \prod_{k \in (A+l) \cap [r]} x_k \cdot \prod_{k \in (A+l-r) \cap [r]} z_k \right\}$$

pro každé $x_{[r]}, z_{[r]} \in \{0,1\}^{[r]}$. Matice Q_U je čtvercová, dimenze $2^r \times 2^r$ a skládá se z nenulových prvků, je tedy jednoznačně daná její v absolutní hodnotě největší vlastní číslo $\lambda_{\max}(Q_U)$, které je reálné kladné. Tímto způsobem se v praxi hodnoty funkce q vyčísľují.

Většinu výsledků týkajících se gibbovských náhodných postupností lze nalézt v klasické monografii Kulliba (1969).

2. Odhady interakcí a testování submodelu

Předpokládáme, že pozorovaná postupnost binárních dat

$$x_{[N]} = (x_1, \dots, x_N) \in \{0, 1\}^N$$

je generována nějakou gibbovskou distribucí μ^{U^0} s neznámým potenciálem $U^0 \in M_A$, $A \in k_+(V)$, $|A| < \infty$. Naším prvním úkolem je odhadnout tento potenciál chápaný jako vektorový parametr.

Pro $U \in M_A$ a $\theta \in R^d$ označme

$$F_A(U, \theta) = q(U) - \sum_{A \in A} U_A \theta_A.$$

Z vlastností funkce q na M_A plyne

- i) $U^1 = U^2$ právě když $\nabla q(U^1) = \nabla q(U^2)$
 ii) $F_A(U^0, \theta^0) = \min_{U \in M_A} F_A(U, \theta^0)$ právě když $\nabla q(U^0) = \theta^0$.

Odhad U^N parametru U^0 tedy získáme minimalizací funkce $F_A(U, \theta^N)$, kde za θ^N vezmeme odhad vektoru $\nabla q(U^0) = \left\{ \mu^{U^0}(C_A) \right\}_{A \in A}$. Tedy

$$\theta_A^N = (N - r(A))^{-1} \sum_{i=1}^{N-r(A)} \prod_{j \in A} x_{j+i} \quad \text{pro každé } A \in A$$

(zde $r(A) = \max\{k \in A\}$). Pokud minimum neexistuje, dodefinujeme U^N libovolně.

Věta: Odhad U^N je konzistentní a asymptoticky normální, tj.

$$U^N \rightarrow U^0 \quad \text{s. j.} \quad \left[\mu^{U^0} \right]$$

$$N^{\frac{1}{2}}(U^N - U^0) \rightarrow N(0, (\nabla^2 q(U^0))^{-1}) \quad \text{v distribuci} \quad \left[\mu^{U^0} \right].$$

Důkaz: Věta 3.2, Janžura (1986).

Rekneme, že máme podezření, že bychom mohli napozorovaná data vysvětlit pomocí jednoduššího modelu $B \subset A$. Budeme tedy testovat hypotézu

$$H^0 : U^0 \in M_B$$

proti alternativě

$$H^1 : U^0 \in M_A \setminus M_B.$$

Zavedme statistiku

$$G^N = \min_{U \in M_B} F_A(U, \theta^N) - \min_{U \in M_A} F_A(U, \theta^N),$$

kteřá je definovaná, pokud existuje $U_A^N \in M_A$, v němž nabývá minima funkce $F_A(\cdot, \theta^N)$. Potom se totiž v nějakém $U_B^N \in M_B$ nabude i minima funkce $F_A(\cdot, \theta^N)$ zúžené na M_B .

Můžeme také psát

$$G^N = H \left(\mu^A \left| \mu^B \right. \right)^{U^N}$$

kde $H(\cdot|\cdot)$ je I-divergence d'Annunzio pro obecné μ, ν výrazem

$$H(\mu|\nu) = \lim_{n \rightarrow \infty} n^{-1} \int \log \frac{\mu(x^{[n]})}{\nu(x^{[n]})} d\mu(x),$$

pokud jsou integrály definovány a limita existuje (pro gibbovské μ, ν je $H(\mu|\nu)$ konečná).

Pokud se minima nenabývá, můžeme G^N dodefinovat libovolně.

Věta: Pro každý $U^0 \in M_B$ platí

$$2 \cdot N \cdot G^N \rightarrow \chi_f^2 \text{ v distribuci } \left[\mu^{U^0} \right],$$

kde $f = |A \setminus B|$.

Důkaz: Věta 4.2, Janžura (1988b).

Hypotézu tedy zamítáme, jestliže statistika 2.N.G^N přesáhne příslušný kvantil rozdělení χ^2 s $|A \setminus B|$ stupni volnosti.

3. Použití

Ve statistických úlohách se systém A (určující model) považuje za známý, přičemž teoreticky může být zcela libovolný konečný. V praxi však existují dvě závažná omezení. V první řadě je zde omezení na "řád" $r = r(A)$, které je dáno výpočetními možnostmi, neboť je třeba opakovaně počítat největší vlastní číslo transfer matice Q_U , která je řádu $2^r \times 2^r$. Je tudíž třeba předpokládat zhruba $r < 10$.

Další omezení je na kardinalitu množin v systému A , tedy $C(A) = \max_{A \in \mathcal{A}} |A|$. Toto omezení je dáno rozsahem výměru N . Jestliže by totiž nebyl rozsah výběru N dostatečně velký vzhledem k $2^{|A|}$ možných konfigurací $x_A \in \{0,1\}^A$, byl by odhad θ_A^N pravděpodobně velmi nepřesný. Zde se ukazuje účelné (a v mnoha praktických úlohách také zcela postačující) uvažovat pouze párové interakce, tj. $A =$

$$= \{A_t = \{0, t\}\}_{t=0, \dots, r}$$

Potom již můžeme $\theta_{A_t}^N$ odhadnout s rozumnou přesností i při poměrně malém rozsahu výběru N a bez dalšího omezení na r . Jedná se pak vlastně o odhad jakýchsi "kovariancí" $\mu^{U^0}(x_0=1, x_t=1)$.

Uvažujme například úlohu predikce. Máme napozorovaná data $(x_1, \dots, x_N) \in \{0,1\}^N$ a předpokládáme, že generující náhodná posloupnost μ je markovská řádu R . Pro optimální predikci (ve smyslu minimální pravděpodobnosti chyby) je třeba umět vyjádřit podíl

$$L_{\mu}^{\mathcal{X}[\mathcal{R}]}) = \frac{\mu(1|\mathcal{X}[\mathcal{R}])}{\mu(0|\mathcal{X}[\mathcal{R}])} \text{ pro každé } \mathcal{X}[\mathcal{R}] \in \{0,1\}^R.$$

Věta: Pro každé $U \in M_{\mathcal{A}}$, $r(\mathcal{A}) = R$, platí

$$L_{\mu}^{\mathcal{X}[\mathcal{R}]}) = \left[\lambda_{\max}(Q_U) \right]^{\frac{1}{R}} \frac{r_U(\mathcal{X}[\mathcal{R}])}{r_U(\vec{\mathcal{X}}[\mathcal{R}])} - 1,$$

kde r_U je levý vlastní vektor příslušný vlastnímu číslu $\lambda_{\max}(Q_U)$ a konfigurace $\vec{\mathcal{X}}[\mathcal{R}]$ je dána tak, že $\vec{x}_1 = 0$, $\vec{x}_i = x_{i-1}$ pro $i = 2, \dots, R$.

Důkaz: Důsledek věty 1. Janžura (1988a).

Jestliže tedy N je "malé" ve srovnání s 2^{R+1} , budeme postupovat

takto:

1. Předpokládáme, že μ je gibbovské s nejvýše párovými interakcemi, tj. $\mu = \mu^U$, $U \in M_{\mathcal{A}_R}$, $\mathcal{A}_R = \left\{ \{0,1\} \right\}_{l=0, \dots, R}$.
2. Odhadneme U postupem uvedeným v části 2.
3. Spočteme $L = L_{\mu}^{\mathcal{X}[\mathcal{R}]}) = L_U(x_N, \dots, x_{N-R+1})$. Jestliže $L > 1$, predikujeme $x_{N+1} = 1$.

Tato metoda funguje poměrně spolehlivě. Jsou-li předpoklady o párových interakcích alespoň zhruba splněny, jsou získané výsledky zjevně lepší, než kdybychom se pokoušeli z dat odhadovat přímo přechodové pravděpodobnosti $\mu(1|\mathcal{X}[\mathcal{R}])$ (viz Janžura (1988a)).

Literatura

- Bishop, Y. M. M., Fienberg, S. E., Holland, P. W. (1975) Discrete Multivariate Analysis: Theory and Practice. MIT Press, Cambridge, Mass.
- Janžura, M. (1986) Estimating interactions in binary data sequences. Kybernetika 22, 277-284.
- Janžura, M. (1988a) Prediction in zero-one random sequences. Problems of Control and Information Theory, Vol. 17(1), 15-22.
- Janžura, M. (1988b) Test for submodel in Gibbs-Markov binary random sequences (v rukopisu).
- Ruelle, D. (1969) Statistical Mechanics. Rigorous Results. Benjamin, New York.