

O SPOLEHLIVOSTI VÝVOJOVÝCH STROMŮ

MARTIN BETINEC

ABSTRACT. Evolutionary tree is a dendrogram that is used to assess genetical similarity of biological objects. Contribution concerns two methods that use bootstrap to verify the confidence of the built tree. First one – suggested by Felsenstein – is relatively simple and is based on counting of occurrences of agreeing parts of trees in a bootstrap sample to the original ones. Despite that it is reasonable by the arguments that will be presented. The second method can be considered as a refinement of the previous one reflecting more precisely the shape of the feature space of the model. Results are illustrated on a real data (malaria gene sequences).

Резюме. Эволюционное дерево представляет дендрограмму, которая используется для достижений генетического совпадения биологических объектов. Результаты моей работы касаются двух методов использующих “bootstrap” для построения доверительной вероятности совпадения построенных деревьев. Первый, предложенный Фелзенштейном, относительно прост и основан на подставе частей деревьев, согласованных с оригинальным деревом в выборке генерированной при помощи “bootstrap”. Приводятся аргументы в пользу этого метода. Второй метод может рассматриваться как уточнение предыдущего, более аккуратно учитывающая форму элементов модели. Результаты иллюстрированы используя реальные данные, точнее, последовательности генов малярии.

1. ÚVOD

Představme si na chvíli, že jsme zvědaví biologové, kteří si na prvního máje vyšli na procházku do Stromovky. Při kochání se pozorováním rozličných organismů vůkol nás pojednou zachvátí šířavá potřeba zjistit, jak probíhala jejich evoluce. Co si v tuto chvíli počít? Někoho napadne, že organismy navzájem si vývojově příbuznější by si měly být podobnější i po stránce genetické. Počneme tedy odebírat vzorky tkání a extrahovat DNA. Jak však z těchto nukleotidových shluků něco smysluplného vyčíst? Zde již naše statistické já vstupuje do myšlenkového experimentu: „Nyní nastal čas pro shlukovou analýzu“.

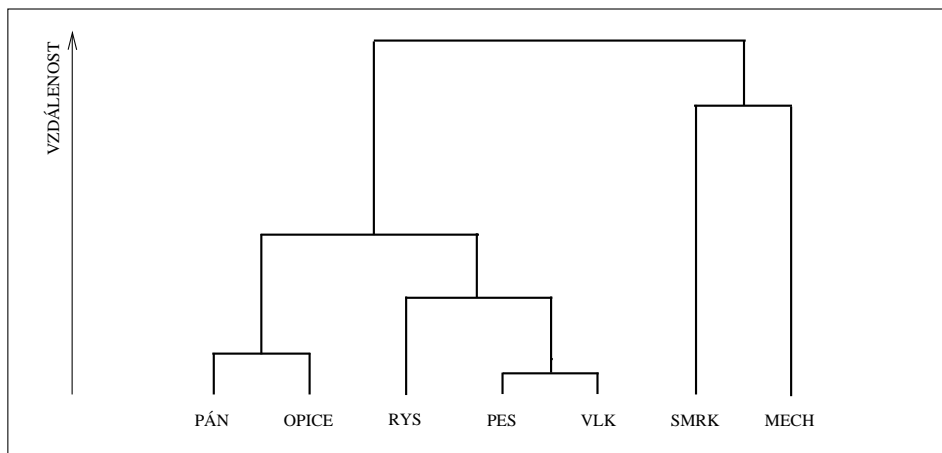
Upřesnění: *Nyní nastal čas pro vývojové stromy.*

Vývojové stromy, v anglické literatuře označované jako *Phylogenetic (Evolutionary) Trees*, představují způsob, jímž se v biologii popisuje evoluce sledovaných organismů. Z matematického hlediska se jedná o klasifikační stromy (dendrogramy). Tato metoda umožňuje nejen názorně vyjádřit „co s čím“ je si v evoluci blízké, ale zároveň také kvantifikovat, jak je ta která skupina vzdálena od ostatních, a to pomocí výšky větve spojující dotyčné skupiny.

2000 *Mathematics Subject Classification*. Primary 62H30; Secondary 62P10.

Klíčová slova. Vývojové stromy, shlukování, klasifikace, bootstrap, DNA, malárie.

Tato práce vznikla za podpory grantů MSM 113200008 a GAČR č. 201/00/0769.



Obr. 1 *Evoluce ve Stromovce.*

Vraťme se k našemu příkladu; vypěstujeme-li vývojový strom, mohl by vypadat (dle mých přírodovědných znalostí) asi jako na obrázku 1. Otázkou zůstává, nakolik se na takto zvolený popis situace mohou spolehnout, tj. do jaké míry výsledný strom (tedy to, co spočtu z naměřených dat) popisuje skutečnou příbuznost sledovaných organismů.

V tomto článku představíme dvě metody využívající bootstrapu pro odhad výše zmíněné spolehlivosti a pokusíme se ukázat, proč lze právě ze vztahu mezi datovým souborem a „boot-populací“ tvrdit něco o souvislosti dat a reality.

Podotkneme jen, že metodikou pěstování evolučního stromu se podrobněji zabývat nebudeme, nýbrž k výpočtům použijeme některého z osvědčených a již dostatečně známých postupů shlukové analýzy.

Pro snadnější orientaci ve značení ještě poznamenejme, že metoda bootstrapu se dotýká tří oblastí reality. Jsou jimi **roviny**:

- teoretická:** týkající se odhadovaných vlastností populace, tj. neznámých parametrů;
- praktická:** odvozená od naměřených dat – odpovídající statistiky budeme označovat stříškou (odhad neznámého parametru μ budeme značit pomocí $\hat{\mu}$);
- virtuální:** vztahující se k populacím vzniklým z bootstrapu – příslušné veličiny označíme hvězdičkami ($\hat{\mu}^*$).

2. METODA RELATIVNÍCH ČETNOSTÍ

2.1. **Data.** Nejprve se budeme věnovat asi nejjednoduššímu postupu, jímž bychom mohli spolehlivost stromu odhadnout, který bývá někdy nazýván jako Felsensteinova metoda. Data, jichž využijeme pro názornější vysvětlení, nashromáždila skupina vědců (viz [4]) zkoumajících část genomu živočicha *Plasmodium falciparum*, parazita žijícího v krvi člověka, který způsobuje malárii. Data, jež jsou přístupná na adrese www.ncbi.nlm.nih.gov, tvoří kompletní sekvence genu kódujícího antigen AMA-1 o délce 1866 nukleotidů a pocházejí z pěti různých oblastí (Keňa, Indie, Thajsko a Venezuela – provincie Bolívar a Amazonas). K ilustraci použijeme 13 vzorků zmíněného genu, které pro zkrácení označíme počátečním písmenem oblasti, z níž pocházejí, a v rámci jedné oblasti je rozlišme čísly, tedy např. B6 označuje vzorek číslo 6 z provincie Bolívar, T3 třetí vzorek z Thajska apod., viz obrázek 2.

Pro $i = 1, \dots, 13$ označme jednotlivé vzorky:

$$\mathbf{x}_i = \underbrace{(A, T, G, \dots, C, T, A, T)}_{1866 \text{ nukleotidů}},$$

kde jsou jednotlivým nukleotidům přiřazena písmena: A – adenin, T – thymin, C – cytosin, G – guanin.

Důležitá vlastnost těchto vzorků spočívá v tom, že každý z nich začíná přesně na začátku genu a končí na jeho konci, přičemž žádné nukleotidy mezi nimi nechybí, jinými slovy – složky vektorů \mathbf{x}_i si navzájem odpovídají.

2.2. Pěstování. Z řádkových vektorů $\mathbf{x}_1, \dots, \mathbf{x}_{13}$ sestavíme datovou matici

$$\mathbf{X}_{(13 \times 1866)} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{13} \end{pmatrix}$$

Klíčovou otázkou pro určení vzdáleností mezi vektory matice \mathbf{X} zůstává, krom volby metriky, kvantifikace nukleotidů. Zde je nutno uvážit jejich chemickou podobnost, pravděpodobnost bodových mutací, apod. Toto zakódování může zásadně ovlivnit podobu vývojového stromu. Já jsem nukleotidy překódoval následovně:

$$(1) \quad (A, G, C, T) = (1, 2, 5, 6).$$

Toto přiřazení zachovává příbuznost jak u purinů (A, G), tak i u pyrimidinů (C, T), zároveň odráží i různost obou skupin. Zde by jistě bylo možno vzdálit obě skupiny od sebe mnohem více, čímž bychom vypěstovali odlišný strom. Stejně tak lze u zvoleného kódování diskutovat o nesymetrii mezi dvojicemi komplementárních bazí (A, T), (G, C). Zde se otevírá možnost, pro konzultace s odborníky v oblasti genetiky a případné pokračování výzkumu. Faktem nicméně zůstává, že pro opačně zvolenou nesymetrii $(A, G, C, T) = (2, 1, 6, 5)$, a stejně tak i pro symetrii $(A, G, C, T) = (1, 2, 6, 5)$ vyrostly stromy velice podobné původnímu. Proto jsem v souladu s literaturou ([3]) zvolil výše zmíněný způsob.

Metriku jsem volil Eukleidovskou:

$$(2) \quad d_E(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{k=1}^{1866} (x_i^k - x_j^k)^2} \quad i, j = 1, \dots, 13,$$

kde x_i^k značí k -tou složku i -tého vzorku. Zřejmě lze vynechat konstatní sloupcové vektory matice \mathbf{X} , které nepřinášejí žádnou informaci o vzdálenosti řádků, neboť jim odpovídající sčítance jsou nulové. Tímto se sníží dimenze matice \mathbf{X} z původních (13×1866) na (13×73) .

Strom $\hat{\Psi}$ vypěstovaný z \mathbf{X} na základě principu nejbližšího souseda ukazuje obrázek 2.

2.3. Spolehlivost. Zřejmě nejjednodušší odhad hladiny spolehlivosti jednotlivých větví stromu $\hat{\Psi}$ bychom získali jako relativní četnosti výskytu jeho větví ve stromech vzniklých prostřednictvím bootstrapu. Pro pevně zvolenou větev v bychom mohli postupovat např. dle následujícího **algoritmu**:

```
count := 0
FOR b = 1 TO B
```

(1) náhodným výběrem s vracením ze sloupců matice \mathbf{X} vytvoříme boot-matici \mathbf{X}_b^* stejné dimenze (zde předpokládáme, že každý ze sloupců matice \mathbf{X} může být vybrán v každém kroku se stejnou pravděpodobností)

(2) z matice \mathbf{X}_b^* vypěstuj strom $\hat{\Psi}_b^*$

(3) IF $\hat{\Psi}_b^* \ni v$ THEN count:=count+1

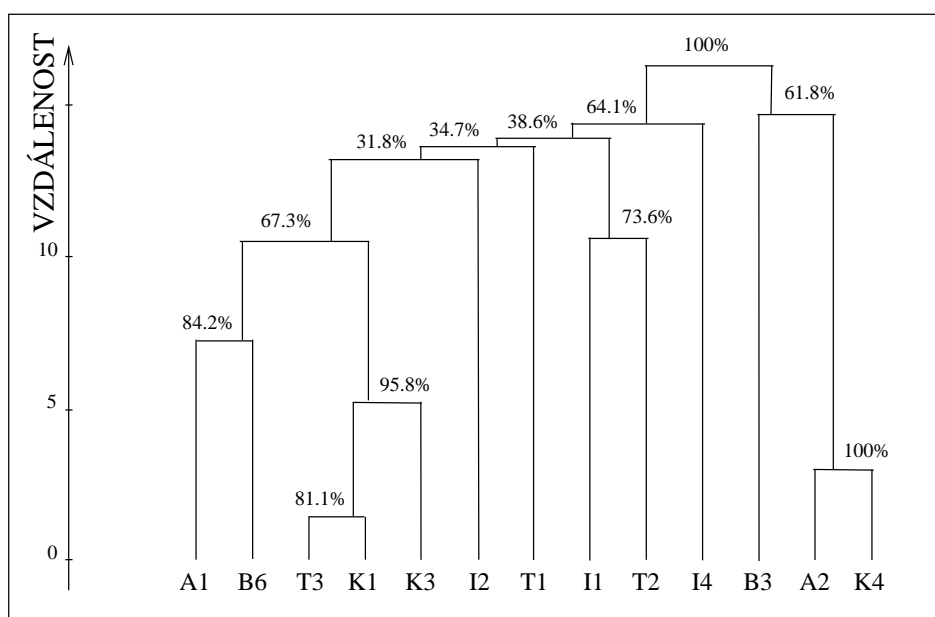
ENDFOR

Spolehlivostní hladinu větve v pak definujeme jako:

$$(3) \quad \alpha_F(v) \stackrel{df}{=} \frac{\text{count}}{B} = \frac{\#\{b : \hat{\Psi}_b^* \ni v\}}{B},$$

kde $\#$ značí „počet prvků“ a B je počet opakování při bootstrapu.

V praxi nás zpravidla zajímají spolehlivostní hladiny pro více než jen jednu větev. Přizpůsobení algoritmu je v tomto případě nasnadě.



Obr. 2 polehlivostní hladiny vývojového stromu bakterie *Plasmodium falciparum*. Jsou nadešány příslušným větvím.

V našem příkladu (kde jsem zvolil $B = 5000$) se větev skládající se právě jen ze vzorků označených jako $T3$ a $K1$ (viz obrázek 2) vyskytla ve 4055 případech. Odhad hladiny spolehlivosti pro zmíněnou větev tedy bude $\alpha_F = 4055/5000 = 0.811$. U původního stromu $\hat{\Psi}$ k této větvi následně přirostl vzorek $K3$, což se opakovalo jen u části z 4055 případů. Větev obsahující právě jen vzorky $T3$, $K1$ a $K3$ se mezi 5000 bootstrapovanými stromy objevila 4790 krát, tzn. její $\alpha_F = 4790/5000 = 0.958$. Ovšem pouze některé z těchto 4790 větví vznikly tak, že se nejdříve spojily vzorky $T3$ a $K1$, k nimž se poté přidal vzorek $K3$. Zbylé srůstaly v jiném pořadí. Větev složenou právě jen ze vzorků $A2$ a $K4$ mělo všech 5000 stromů, proto je její odhadnutá spolehlivost 100%. Její nadvětev, tj. větev, v níž se nacházejí právě jen prvky $A2$, $K4$ a $B3$, bylo možno nalézt už pouze u necelých 62% bootstrapovaných stromů.

V ostatních případech k původní větvi (tj. A2 a K4) buď přirostlo dříve něco jiného než B3, anebo se B3 už dříve spojilo s nějakým jiným vzorkem, a pak teprve takto vzniklý shluk přirostl k A2 a K4.

Podívejme se nyní podrobněji na vztah takto spočteného odhadu spolehlivosti ke skutečnosti.

2.4. Teoretické pozadí. O sloupcích $\mathbf{x}^1, \dots, \mathbf{x}^{73}$ matice

$$(4) \quad \mathbf{X} = (\mathbf{x}^1, \dots, \mathbf{x}^{73}), \quad \text{kde } x_l^k \in \{A, G, T, C\} \quad k = 1, \dots, 73, \\ l = 1, \dots, 13,$$

předpokládáme, že tvoří náhodný výběr o rozsahu $n = 73$ z nějakého pravděpodobnostního rozdělení na prostoru **všech možných nekonstantních** 13-ti rozměrných vektorů $\boldsymbol{\xi}$ nad abecedou $\{A, G, T, C\}$, tento prostor označme \mathcal{X} .

$$(5) \quad \mathcal{X} = \{\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_K\} \quad K = 4^{13} - 4 = 67\,108\,860.$$

Schématem

$$(6) \quad \mathbf{X} \implies \text{Shluková analýza} \implies \widehat{\Psi} \quad \text{resp.} \quad \mathbf{X}^* \implies \text{Shluková analýza} \implies \widehat{\Psi}^*,$$

lze popsat jak oblast praktickou, tak oblast virtuální, nicméně těžko rovinu teoretickou, jejíž vztah k oběm předešlým nás zajímá. Proto se přesuneme do prostoru \mathcal{P} , tedy prostoru **pravděpodobností zahrnutí** π_i vektorů $\boldsymbol{\xi}_i \in \mathcal{X}$ do výběru:

$$(7) \quad \mathcal{P} \stackrel{\text{df}}{=} \left\{ \boldsymbol{\pi} = (\pi_1, \dots, \pi_K)' : \pi_i = \mathbf{P}(\text{zahrnutí } \boldsymbol{\xi}_i \text{ do výběru}), \quad i = 1, \dots, K \right\}$$

Teoretickým pravděpodobnostem $\boldsymbol{\pi}$ by na rovině praktické odpovídaly **pozorované četnosti** výskytu vektorů $\boldsymbol{\xi}_i$ ve sloupcích datové matice \mathbf{X} , tj.

$$(8) \quad \boldsymbol{\Upsilon} = (\Upsilon_1, \dots, \Upsilon_K)', \quad \text{kde } \Upsilon_i \stackrel{\text{df}}{=} \#\{j : \mathbf{x}^j = \boldsymbol{\xi}_i\} \quad i = 1, \dots, K,$$

resp. **relativní četnosti**

$$(9) \quad \widehat{\boldsymbol{\pi}} = (\widehat{\pi}_1, \dots, \widehat{\pi}_K)', \quad \text{kde } \widehat{\pi}_i \stackrel{\text{df}}{=} \frac{1}{n} \#\{j : \mathbf{x}^j = \boldsymbol{\xi}_i\} = \frac{\Upsilon_i}{n}, \quad i = 1, \dots, K.$$

Analogicky bychom z bootstrapovaných matic $\mathbf{X}_b^* = (\mathbf{x}_b^{*1}, \dots, \mathbf{x}_b^{*73})$, $b = 1, \dots, B$, získali:

$$(10) \quad \boldsymbol{\Upsilon}_b^* = (\Upsilon_{b,1}^*, \dots, \Upsilon_{b,K}^*)', \quad \text{kde } \Upsilon_{b,i}^* \stackrel{\text{df}}{=} \#\{j : \mathbf{x}_b^{*j} = \boldsymbol{\xi}_i\}, \quad i = 1, \dots, K,$$

$$(11) \quad \widehat{\boldsymbol{\pi}}_b^* = (\widehat{\pi}_{b,1}^*, \dots, \widehat{\pi}_{b,K}^*)', \quad \text{kde } \widehat{\pi}_{b,i}^* \stackrel{\text{df}}{=} \frac{1}{n} \Upsilon_{b,i}^*.$$

Zřejmě platí:

$$(12) \quad \boldsymbol{\Upsilon} \sim \text{Multi}_K(n, \boldsymbol{\pi}) \quad \text{a} \quad \boldsymbol{\Upsilon}^* \sim \text{Multi}_K(n, \widehat{\boldsymbol{\pi}}),$$

kde $K = 67\,108\,860$ a $n = 73$.

Lze namítnout, že právě zmíněný přístup nás nutí pracovat s velmi dlouhými vektory obsahujícími navíc téměř výhradně samé nuly. Nicméně, pro teoretická odvození se tato reprezentace ukazuje být přínosnou. Krom toho, z důvodů, jež budou vysvětleny později, vyplyne, že při výpočtech se lze omezit jen na nenulové složky.

Pěstování stromů, popsané schématem (6), lze nyní vyjádřit jako:

$$(13) \quad \widehat{\boldsymbol{\pi}} \implies \widehat{\Psi}, \quad \text{resp.} \quad \widehat{\boldsymbol{\pi}}_b^* \implies \widehat{\Psi}_b^*.$$

Podobně si lze představit, že i teoretickým pravděpodobnostem $\boldsymbol{\pi}$ odpovídá nějaký teoretický strom Ψ vystihující skutečné vztahy mezi sledovanými organismy.

Stojí za zmínku, že ze dvou různých pravděpodobnostních vektorů $\hat{\pi}^{(1)}$ a $\hat{\pi}^{(2)}$ může vzniknout tentýž strom. Prostor \mathcal{P} je tak rozložen na disjunkttní třídy ekvivalence $\mathcal{P} = \bigcup_i \mathcal{P}_i$, pro něž platí:

$$(14) \quad \{ \mathcal{P}_i \ni \hat{\pi}^{(1)} \ \& \ \mathcal{P}_j \ni \hat{\pi}^{(2)} \} \implies \{ \hat{\Psi}^{(1)} \equiv \hat{\Psi}^{(2)} \Leftrightarrow i = j \},$$

kde $\hat{\Psi}^{(j)}$ odpovídá stromu, který vznikl z $\hat{\pi}^{(j)}$, $j = 1, 2$.

Vraťme se nyní k původní otázce po spolehlivosti námi vypěstovaného stromu $\hat{\Psi}$. Tážeme se po pravděpodobnosti, s níž se $\hat{\Psi}$ bude shodovat s teoretickým vzorem Ψ . Zajímá nás tedy:

$$(15) \quad \alpha = \mathbf{P}(\hat{\Psi} \equiv \Psi) = \mathbf{P}(\pi \in \mathcal{P}_i | \hat{\pi} \in \mathcal{P}_i)$$

Tuto pravděpodobnost jsme odhadli pomocí (3) jako empirickou pravděpodobnost (relativní četnost) toho, že $\hat{\Psi}^*$ vypěstovaný nad \mathbf{X}^* má danou větev shodnou s původním stromem $\hat{\Psi}$, tedy:

$$(16) \quad \alpha_F = \hat{\mathbf{P}}_{\hat{\pi}^* | \hat{\pi}}(\hat{\pi}^* \in \mathcal{P}_i | \hat{\pi} \in \mathcal{P}_i).$$

Abychom si mohli udělat představu o chybě, jíž jsme se tímto dopustili, potřebovali bychom znat rozdělení $\hat{\pi}^*$ (resp. aposteriorní rozdělení π) při daném $\hat{\pi}$, tedy $\mathcal{L}(\hat{\pi}^* | \hat{\pi})$ (resp. $\mathcal{L}(\pi | \hat{\pi})$).

První z nich známe, neboť ze vztahů (11) a (12) plyne, že:

$$(17) \quad \mathcal{L}(n\hat{\pi}^* | \hat{\pi}) = \text{Multi}_K(n, \hat{\pi}).$$

Co se týče aposteriorního rozdělení $\mathcal{L}(\pi | \hat{\pi})$, lze je vyjádřit přes výběrové rozdělení $\hat{\pi}$ při daném π , pro které díky uvažovanému modelu platí:

$$(18) \quad \mathcal{L}(n\hat{\pi} | \pi) = \text{Multi}_K(n, \pi).$$

Multinomické rozdělení tvoří konjugovaný systém spolu s rozdělením Dirichletovým, a tudíž lze snadno odvodit (viz též např. [5]), že platí:

Tvrzení 2.1. *Řídí-li se apriorní rozdělení π principem neurčitosti, pak*

$$(19) \quad \mathcal{L}(\pi | \hat{\pi}) = \text{Dir}_K(n\hat{\pi}).$$

Tvrzení 2.2. *$\mathcal{L}(\hat{\pi}^* | \hat{\pi})$ a $\mathcal{L}(\pi | \hat{\pi})$ mají následující vlastnosti:*

- (1) *Obě jsou soustředěna jen na takových ξ_i , pro něž $\pi_i > 0$, tj. na sloupcích matice \mathbf{X} .*
- (2) *Mají stejnou střední hodnotu.*
- (3) *Pro kovarianční matice platí:*

$$(20) \quad \left(\text{Cov}(\hat{\pi}_i^*, \hat{\pi}_j^* | \hat{\pi}) \right)_{i,j=1}^K = \frac{n+1}{n} \cdot \left(\text{Cov}(\pi_i, \pi_j | \hat{\pi}) \right)_{i,j=1}^K.$$

Z tvrzení je patrné, že odhad hladin spolehlivosti větví vývojového stromu (pomocí relativních četností výskytu těchto větví v bootstrapových populacích) definovaný vztahem (3) má krom předností, jež spočívají především v jednoduchosti, snadné implementovatelnosti a v rychlosti výpočtu i své omezení. Problém nastává v okamžiku, kdy nelze splnit předpoklady tvrzení 2.1 o apriorním rozdělení. Potom nemusí platit závěry tvrzení 2.2 a rozdělení $\mathcal{L}(\hat{\pi}^* | \hat{\pi})$ a $\mathcal{L}(\pi | \hat{\pi})$ mohou být natolik odlišná, že odhadnutá hladina spolehlivosti α_F může být značně vzdálena od skutečné hladiny α . Vzhledem ke složitosti prostoru \mathcal{P} není příliš zřejmé, jak tyto situace rozeznat. Pokusme se tedy nalézt lepší přístup.

3. KTERAK POSTUPOVAT LÉPE

3.1. Teorie. Podobně jako v odstavci 2 vyjdeme z modelu (7). Index $i = 1$ vyhradíme pro tu část prostoru $\mathcal{P} = \dot{\bigcup}_i \mathcal{P}_i$ při dělení dle pravidla (14), která obsahuje $\hat{\pi}$, zatímco $\mathcal{P}_1^C = \dot{\bigcup}_{i \neq 1} \mathcal{P}_i$ budiž její doplněk.

Jedna z cest, která se nyní nabízí, je vést v patrnosti nejen skutečnost, že $\hat{\pi} \in \mathcal{P}_1$, jako jsme činili v předchozím, ale zohlednit i vzdálenost vektoru $\hat{\pi}$ od hranice oblasti \mathcal{P}_1 . Tato snaha je motivována úvahou, že vektor π , z něhož je vektor pozorovaných relativních četností $\hat{\pi}$ dle předpokládaného modelu (9) nagenеровán, by se měl nacházet „někde poblíž“ vektoru $\hat{\pi}$, a proto by hladina spolehlivosti

$$(21) \quad \alpha = \mathbf{P}(\pi \in \mathcal{P}_1 \mid \hat{\pi} \in \mathcal{P}_1),$$

měla být vyšší, pokud by $\hat{\pi}$ leželo „někde uprostřed“ \mathcal{P}_1 , než by tomu bylo v případě $\hat{\pi}$ umístěného v blízkosti hranice oblasti \mathcal{P}_1 .

Označme π_0 jako nejbližší hraniční bod k $\hat{\pi}$, tj.

$$(22) \quad \pi_0 \stackrel{df}{=} \arg \min_{\tilde{\pi} \in \partial \mathcal{P}_1} d(\hat{\pi}, \tilde{\pi}),$$

kde $\partial \mathcal{P}_1$ značí hranici \mathcal{P}_1 a $d(\cdot, \cdot)$ je vhodná metrika na \mathcal{P} , tj. nezáporná funkce splňující navíc pro libovolné $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{P}$, že $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$, dále že $d(\mathbf{x}, \mathbf{y}) = 0$ právě tehdy, když $\mathbf{x} = \mathbf{y}$, a $d(\mathbf{x}, \mathbf{y}) \leq d(\mathbf{x}, \mathbf{z}) + d(\mathbf{z}, \mathbf{y})$.

V předchozím případě jsme boot-populaci obdrželi prostřednictvím četností výskytu vygenerovaných z $\hat{\pi}$, viz (12), nyní je nagenеровujeme z π_0 , tedy pro $b = 1, \dots, B_2$

$$(23) \quad \mathbf{Y}_b^{**} = (\mathbf{Y}_{1b}^{**}, \dots, \mathbf{Y}_{Kb}^{**})', \quad \text{kde} \quad \mathbf{Y}_{ib}^{**} \stackrel{df}{=} \#\{j : \mathbf{x}_b^{**j} = \boldsymbol{\xi}_i\}, \quad i = 1, \dots, K,$$

$$(24) \quad \mathbf{Y}_b^{**} \sim \text{Multi}_K(n, \pi_0) \quad \text{a} \quad \hat{\pi}_b^{**} = \frac{1}{n} \mathbf{Y}_b^{**}.$$

Chceme-li odhadnout hladinu spolehlivosti (21), můžeme uvažovat takto: čím blíže budou prvky boot-populace $\hat{\pi}_b^{**}$ vygenerované z hraničního bodu π_0 ležet k \mathcal{P}_1^C ve srovnání s polohou $\hat{\pi}$ (tj. čím menší bude jejich odhadnutý rozptyl), tím menší je pravděpodobnost, že by vektor skutečných pravděpodobností π ležel až v \mathcal{P}_1^C , neboť předpokládáme, že právě z π byl nagenеровán i vektor pozorovaných relativních četností $\hat{\pi}$, a to stejným způsobem jako byly nagenерованы vektory $\hat{\pi}_b^{**}$ z π_0 , srov. (12) a (24). Tím větší je tedy α . Hladinu spolehlivosti zavedeme následovně:

$$(25) \quad \alpha_E \stackrel{df}{=} \hat{\mathbf{P}}_{\hat{\pi}^{**} \mid \pi_0} (d(\hat{\pi}^{**}, \mathcal{P}_1^C) \leq d(\hat{\pi}, \mathcal{P}_1^C) \mid \hat{\pi} \in \mathcal{P}_1).$$

Takto definovaná hladina nám zaručuje, že $1 - \alpha_E$ je obdobou dosažené hladiny testu (tzv. p -hodnoty) v klasickém testování hypotéz, což lze nahlédnout např. z následujícího příkladu.

Příklad 3.1. Necht $Z_1, \dots, Z_n \stackrel{iid}{\sim} N(\mu, 1)$.

Testujme hypotézu:

$$(26) \quad \text{H: } \underbrace{\mu \leq \mu_0}_{\hat{\pi} \in \mathcal{P}_1^C} \quad \text{vs.} \quad \text{K: } \underbrace{\mu > \mu_0}_{\hat{\pi} \in \mathcal{P}_1}$$

Zde $\pi_0 \in \partial \mathcal{P}_1$ odpovídá hraničnímu bodu hypotézy a alternativy, tj. hodnotě $\mu = \mu_0$. Provedením experimentu spočteme odhad $\hat{\mu} = \hat{\mu}(Z_1, \dots, Z_n)$ – analogii $\hat{\pi}$. Při známé hodnotě $\hat{\mu}$, a tedy i známé $\mathcal{P}_1 \ni \hat{\mu}$, chceme potvrdit, že také $\mu \in \mathcal{P}_1$, tedy zamítnout

opak. Proto $H \equiv \mathcal{P}_1^C$. Označme dále $Z_1^{**}, \dots, Z_n^{**} \stackrel{iid}{\sim} N(\mu_0, 1)$ bootstrapovaný výběr z hypotézy H a příslušný odhad $\hat{\mu}^{**} = \hat{\mu}^{**}(Z_1^{**}, \dots, Z_n^{**})$ (analogie $\hat{\pi}^{**}$), potom

$$(27) \quad p = \mathbf{P}_{\mu_0}(\hat{\mu}^{**} > \hat{\mu}) = \mathbf{P}_{\mu_0}(\hat{\mu}^{**} - \mu_0 > \hat{\mu} - \mu_0).$$

Zároveň platí:

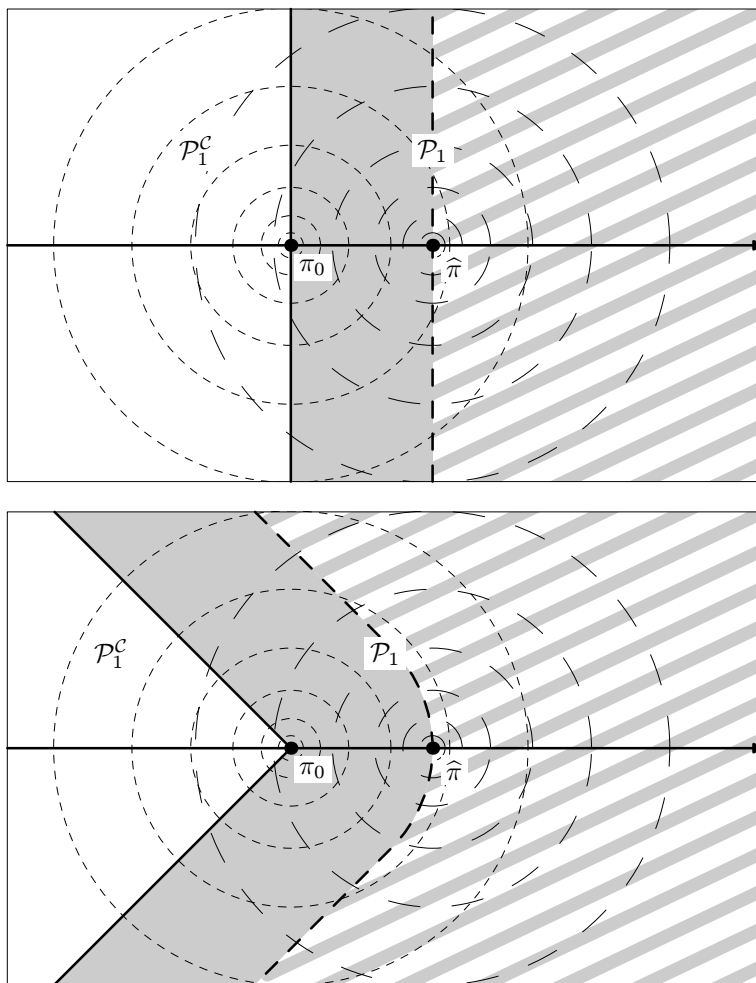
$$(28) \quad d(\hat{\mu}^{**}, \mathcal{P}_1^C) = \begin{cases} \hat{\mu}^{**} - \mu_0 & \text{pro } \hat{\mu}^{**} > \mu_0 \\ 0 & \text{jinak,} \end{cases}$$

odkud

$$(29) \quad \mathbf{P}_{\mu_0}(d(\hat{\mu}^{**}, \mathcal{P}_1^C) > \hat{\mu} - \mu_0) = \mathbf{P}_{\mu_0}(\hat{\mu}^{**} - \mu_0 > \hat{\mu} - \mu_0 > 0),$$

a tedy

$$(30) \quad p = 1 - \alpha_E$$



Obr. 3, 4 Závislost hladin α_F a α_E na tvaru oblasti: nezakřivené hranice (nahore), zakřivené hranice (dole).

Nově zavedená hladina α_E bere více v potaz tvar oblasti \mathcal{P}_1 než α_F . Jakým způsobem je zřejmé z obrázků 3 a 4, kde je chování obou hladin demonstrováno opět na

příkladu (tentokrát dvourozměrném) normálních rozdělení se stejným rozptylem a se střední hodnotou v $(\boldsymbol{\pi}_0, 0)$, resp. v $(\hat{\boldsymbol{\pi}}, 0)$. Nevybarvená plocha odpovídá oblasti \mathcal{P}_1^C (tj. odpovídá $1 - \alpha_F$). Z vybarvené oblasti \mathcal{P}_1 je pak šrafovaním vyznačena ta část, kde $d(\hat{\boldsymbol{\pi}}^{**}, \mathcal{P}_1^C) > d(\hat{\boldsymbol{\pi}}, \mathcal{P}_1^C)$, tedy ta, jež opovídá $1 - \alpha_E$. Zatímco pro nezakřivené hranice se obě hladiny shodují, v případě ryze konvexního tvaru oblasti \mathcal{P}_1^C platí:

$$(31) \quad 1 - \alpha_F = \hat{\mathbf{P}}_{\hat{\boldsymbol{\pi}}^* | \hat{\boldsymbol{\pi}}}(\hat{\boldsymbol{\pi}}^* \in \mathcal{P}_1^C | \hat{\boldsymbol{\pi}} \in \mathcal{P}_1) \leq$$

$$(32) \quad \leq \mathbf{P}_{\hat{\boldsymbol{\pi}}^{**} | \boldsymbol{\pi}_0}(d(\hat{\boldsymbol{\pi}}^{**}, \mathcal{P}_1^C) > d(\hat{\boldsymbol{\pi}}, \mathcal{P}_1^C) | \boldsymbol{\pi}_0 \in \partial\mathcal{P}_1) = 1 - \alpha_E.$$

Tedy

$$(33) \quad \alpha_E \leq \alpha_F.$$

Nyní zbývá vyřešit otázku, jak spočíst α_E .

3.2. Intervaly spolehlivosti. Vzhledem ke složitosti prostoru \mathcal{P} není snadné vypočítat α_E přímo. Odložíme-li však na chvíli tento záměr, můžeme se přesunout do zdánlivě nesouvisících oblastí matematické statistiky – do problematiky intervalových odhadů, které nám v konečném důsledku pomohou nalézt vhodnou aproximaci α_E .

Představme si na okamžik, že stojíme před opačným problémem: kterak při zadané hladině určit interval spolehlivosti, tedy oblast, jež s danou spolehlivostí pokryje skutečnou hodnotu neznámého parametru.

Pro názornost pracujeme s jednoduchým modelem, kde data pocházejí z rozdělení určeného jednoparametrickým systémem hustot. Parametr θ odhadujeme reálnou statistikou $\hat{\theta}$ s rozdělením:

$$(34) \quad \hat{\theta} \sim f_\theta$$

3.2.1. Standardní asymptotický interval spolehlivosti. Využijeme-li normální aproximace

$$(35) \quad \frac{(\hat{\theta} - \theta)}{\hat{\sigma}} \stackrel{as}{\approx} N(0, 1),$$

kde $\stackrel{as}{\approx}$ označuje, že rozdělení statistiky $\sigma^{-1}(\hat{\theta} - \theta)$ je asymptoticky normální a kde $\hat{\sigma}$ je konzistentní odhad rozptylu, obdržíme často užívaný $100(1 - \alpha)$ procentní asymptotický interval spolehlivosti pro θ ve tvaru

$$(36) \quad \left(\hat{\theta} + \hat{\sigma}z^{(\frac{\alpha}{2})}; \hat{\theta} + \hat{\sigma}z^{(1 - \frac{\alpha}{2})} \right),$$

kde

$$(37) \quad z^{(\alpha)} = \Phi^{-1}(\alpha).$$

Tento interval ovšem může být v případě konečných výběrů značně nepřesný.

3.2.2. Percentilová metoda. Problému s normalitou $(\hat{\theta} - \theta)$ bychom se mohli pokusit vyhnout například odhadnutím percentilů pomocí bootstrapu, kdy odhadneme θ prostřednictvím statistiky $\hat{\theta}^*$ spočtené z B -tice bootstrapovaných populací. Z empirické distribuční funkce bootstrapovaného odhadu $\hat{\theta}^*$

$$(38) \quad \hat{F}_{\hat{\theta}}(t) = \hat{\mathbf{P}}_{\hat{\theta}}(\hat{\theta}^* \leq t) = B^{-1} \#\{b : \hat{\theta}_b^* \leq t\},$$

získáme pro pevné $\alpha \in (0; 0.5)$ meze $t_d^{(\frac{\alpha}{2})}$, $t_h^{(\frac{\alpha}{2})}$ intervalu spolehlivosti jakožto příslušné percentily. Přibližný $100(1 - \alpha)\%$ -ní interval spolehlivosti pro θ pak má tvar

$$(39) \quad (t_d^{(\frac{\alpha}{2})}; t_h^{(\frac{\alpha}{2})}) = \left(\widehat{F}_{\widehat{\theta}}^{-1}\left(\frac{\alpha}{2}\right); \widehat{F}_{\widehat{\theta}}^{-1}\left(1 - \frac{\alpha}{2}\right) \right).$$

Tento odhad – ač mnohdy lepší než standardní interval – není příliš důvěryhodný, jak co se týče vychýlení, tak i co do své délky (blíže viz např. [2]). Lze jej nicméně vhodně modifikovat do té míry, že se stane použitelným.

3.2.3. Oprava vychýlení percentilové metody. Předpokládejme, že hustota rozdělení θ^* náleží do stejného systému jako hustota statistiky $\widehat{\theta}$, s tím rozdílem, že namísto hodnotou θ je určena jejím bodovým odhadem, tj. $\widehat{\theta}$. Tedy:

$$(40) \quad \widehat{\theta}^* \sim f_{\widehat{\theta}}.$$

Distribuční funkci $\widehat{\theta}^*$ lze tedy vyjádřit jako:

$$(41) \quad F_{\widehat{\theta}}(t) = \int_{-\infty}^t f_{\widehat{\theta}}(\widehat{\theta}^*) d\widehat{\theta}^* = \mathbf{P}_{\widehat{\theta}}(\widehat{\theta}^* \leq t)$$

Dále předpokládejme, že normality $(\widehat{\theta} - \theta)$ i stability rozptylu lze dosáhnout vhodnou transformací; tedy, že existuje taková monotónně rostoucí funkce $g : \Theta \rightarrow \Gamma$ a konstanty a a $\sigma > 0$, že pro veličiny

$$(42) \quad \gamma = g(\theta), \quad \widehat{\gamma} = g(\widehat{\theta}) \quad \text{a} \quad \widehat{\gamma}^* = g(\widehat{\theta}^*)$$

platí:

$$(43) \quad (\widehat{\gamma} - \gamma) \sim N(-z_0\sigma; \sigma^2) \quad \text{a} \quad (\widehat{\gamma}^* - \widehat{\gamma}) \sim N(-z_0\sigma; \sigma^2).$$

Rádi bychom stanovili interval spolehlivosti pro θ , tedy $(t_d^{(\frac{\alpha}{2})}, t_h^{(\frac{\alpha}{2})})$ tak, aby platilo

$$(44) \quad \mathbf{P}_{\theta}(t_d^{(\frac{\alpha}{2})} < \theta < t_h^{(\frac{\alpha}{2})}) = 1 - \alpha.$$

Pro zjednodušení značení nebudeme nadále zdůrazňovat zřejmou závislost mezi intervalu spolehlivosti $t_d^{(\frac{\alpha}{2})}$ a $t_h^{(\frac{\alpha}{2})}$ na hladině α . Napříště budeme tedy psát pouze t_d a t_h , nebude-li nutno jinak.

Díky monotonní transformace g lze pracovat v prostoru Γ :

$$(45) \quad \begin{aligned} 1 - \alpha &= \mathbf{P}_{\theta}(g(t_d) < \gamma < g(t_h)) = \\ &= \mathbf{P}_{\theta}\left(\frac{\widehat{\gamma} - g(t_d)}{\sigma} + z_0 > \frac{\widehat{\gamma} - \gamma}{\sigma} + z_0 > \frac{\widehat{\gamma} - g(t_h)}{\sigma} + z_0\right). \end{aligned}$$

Z normovaného normálního rozdělení prostředního výrazu plyne:

$$(46) \quad g(t_d) = \widehat{\gamma} + z_0\sigma - z^{(1-\frac{\alpha}{2})}\sigma \quad g(t_h) = \widehat{\gamma} + z_0\sigma + z^{(1-\frac{\alpha}{2})}\sigma.$$

Pro distribuční funkci veličiny $\widehat{\gamma}^*$

$$(47) \quad G(s) = \mathbf{P}_{\widehat{\gamma}}(\widehat{\gamma}^* \leq s)$$

díky monotonní funkce g platí:

$$(48) \quad G(g(t)) = \mathbf{P}_{\widehat{\gamma}}(\widehat{\gamma}^* \leq g(t)) = \mathbf{P}_{\widehat{\gamma}}(g(\widehat{\theta}^*) \leq g(t)) = \mathbf{P}_{\widehat{\theta}}(\widehat{\theta}^* \leq t) = F_{\widehat{\theta}}(t).$$

Tedy

$$(49) \quad F_{\widehat{\theta}} = G \circ g, \quad \text{odkud} \quad F_{\widehat{\theta}}^{-1} = g^{-1} \circ G^{-1}.$$

Odtud z normality (43) obdržíme jednak opravu vychýlení z_0 :

$$(50) \quad F_{\hat{\theta}}(\hat{\theta}) = G(g(\hat{\theta})) = \mathbf{P}_{\hat{\gamma}}(\hat{\gamma}^* \leq \hat{\gamma}) = \hat{\mathbf{P}}_{\hat{\gamma}}\left(\frac{\hat{\gamma}^* - \hat{\gamma} + z_0\sigma}{\sigma} \leq z_0\right) = \Phi(z_0),$$

a tedy

$$(51) \quad z_0 = \Phi^{-1}(F_{\hat{\theta}}(\hat{\theta})),$$

navíc díky (43), tj. stejnému rozdělení rozdílů $(\hat{\gamma} - \gamma)$ a $(\hat{\gamma}^* - \hat{\gamma})$, vzhledem k (46) platí:

$$(52) \quad \begin{aligned} F_{\hat{\theta}}(t_h) &= G(g(t_h)) = \mathbf{P}_{\hat{\gamma}}(\hat{\gamma}^* \leq \hat{\gamma} + z_0\sigma + z^{(1-\frac{\alpha}{2})}\sigma) = \\ &= \mathbf{P}_{\gamma}(\hat{\gamma} \leq \gamma + z_0\sigma + z^{(1-\frac{\alpha}{2})}\sigma) = \mathbf{P}_{\gamma}\left(\frac{\hat{\gamma} - \gamma + z_0\sigma}{\sigma} \leq \frac{2z_0\sigma + z^{(1-\frac{\alpha}{2})}\sigma}{\sigma}\right) \\ &= \Phi(2z_0 + z^{(1-\frac{\alpha}{2})}). \end{aligned}$$

Dosadíme-li rovnici (52) do argumentu $F_{\hat{\theta}}^{-1}$, získáme horní a analogicky i dolní mez intervalu spolehlivosti:

$$(53) \quad t_d^{(\frac{\alpha}{2})} = F_{\hat{\theta}}^{-1}(\Phi(2z_0 - z^{(1-\frac{\alpha}{2})})) \quad t_h^{(\frac{\alpha}{2})} = F_{\hat{\theta}}^{-1}(\Phi(2z_0 + z^{(1-\frac{\alpha}{2})}))$$

Poznámka 3.1. Povšimněme si, že o zobrazení g stačí předpokládat pouze jeho existenci; nemusíme tedy znát jeho explicitní tvar.

Poznámka 3.2. Od rozdělení statistik $(\hat{\gamma} - \gamma)$ a $(\hat{\gamma}^* - \hat{\gamma})$ navíc nemusíme nutně žádat, aby bylo normální, postačilo by libovolné symetrické rozdělení, jehož kvantily jsme schopni spočítat.

Poznámka 3.3. $F_{\hat{\theta}}$ lze odhadnout např. pomocí empirické distribuční funkce $\hat{F}_{\hat{\theta}}$, viz (38)

Poznámka 3.4. Pro rozdělení $F_{\hat{\theta}}$, pro něž je $\hat{\theta}$ rovno mediánu, vyjde korekční člen $z_0 = 0$.

Poznámka 3.5. Požadavky na stabilitu rozptylu rozdílů $(\hat{\gamma} - \gamma)$ a $(\hat{\gamma}^* - \hat{\gamma})$ a jejich normalitu si mohou protirečit.

3.2.4. *Opravená percentilová metoda bez konstantního rozptylu.* Upustíme-li od předpokladu stability rozptylu rozdílů (43) a požadujeme-li, aby pro vhodné konstanty a a z_0 (tj. $a > -1/\gamma$ pro $\gamma > 0$ a $a < -1/\gamma$ pro $\gamma < 0$) platilo:

$$(54) \quad \begin{aligned} &(\hat{\gamma} - \gamma) \sim N(-z_0\sigma_{\gamma}; \sigma_{\gamma}^2), & \text{kde } \sigma_{\gamma} &= 1 + a\gamma, \\ \text{a } &(\hat{\gamma}^* - \hat{\gamma}) \sim N(-z_0\sigma_{\hat{\gamma}}; \sigma_{\hat{\gamma}}^2), & \text{kde } \sigma_{\hat{\gamma}} &= 1 + a\hat{\gamma}, \end{aligned}$$

dospějeme pomocí stejných úvah jako v předchozím případě, tj. (44)–(46), k nerovnostem udávajícím interval spolehlivosti pro γ :

$$(55) \quad \begin{aligned} &\hat{\gamma} + z_0\sigma_{\gamma} - z^{(1-\frac{\alpha}{2})}\sigma_{\gamma} < \gamma < \hat{\gamma} + z_0\sigma_{\gamma} + z^{(1-\frac{\alpha}{2})}\sigma_{\gamma} \\ &\hat{\gamma} + (z_0 - z^{(1-\frac{\alpha}{2})})(1 + a\gamma) < \gamma < \hat{\gamma} + (z_0 + z^{(1-\frac{\alpha}{2})})(1 + a\gamma) \\ &\hat{\gamma} + z_0 - z^{(1-\frac{\alpha}{2})} + \gamma a(z_0 - z^{(1-\frac{\alpha}{2})}) < \gamma < \hat{\gamma} + z_0 + z^{(1-\frac{\alpha}{2})} + \gamma a(z_0 + z^{(1-\frac{\alpha}{2})}) \end{aligned}$$

$$(56) \quad \gamma_d \stackrel{df}{=} \frac{\hat{\gamma} + z_0 - z^{(1-\frac{\alpha}{2})}}{1 - a(z_0 - z^{(1-\frac{\alpha}{2})})} < \gamma < \frac{\hat{\gamma} + z_0 + z^{(1-\frac{\alpha}{2})}}{1 - a(z_0 + z^{(1-\frac{\alpha}{2})})} \stackrel{df}{=} \gamma_h.$$

Obě meze lze ještě upravit:

$$\begin{aligned}
 \gamma_d &= \frac{\hat{\gamma} + z_0 + z^{(\frac{\alpha}{2})}}{1 - a(z_0 + z^{(\frac{\alpha}{2})})} = \hat{\gamma} \left(1 + a \frac{z_0 + z^{(\frac{\alpha}{2})}}{1 - a(z_0 + z^{(\frac{\alpha}{2})})} \right) + \frac{z_0 + z^{(\frac{\alpha}{2})}}{1 - a(z_0 + z^{(\frac{\alpha}{2})})} \\
 &= \hat{\gamma} + \frac{z_0 + z^{(\frac{\alpha}{2})}}{1 - a(z_0 + z^{(\frac{\alpha}{2})})} (1 + a\hat{\gamma}) = \hat{\gamma} + \frac{z_0 + z^{(\frac{\alpha}{2})}}{1 - a(z_0 + z^{(\frac{\alpha}{2})})} \sigma_{\hat{\gamma}}, \\
 (57) \quad \gamma_h &= \hat{\gamma} + \frac{z_0 + z^{(1-\frac{\alpha}{2})}}{1 - a(z_0 + z^{(1-\frac{\alpha}{2})})} \sigma_{\hat{\gamma}}
 \end{aligned}$$

Je patrné, že obě meze lze vyjádřit jedním předpisem jako funkci hladiny α . Necht

$$(58) \quad \gamma^{(\alpha)} \stackrel{df}{=} \hat{\gamma} + \frac{z_0 + z^{(\alpha)}}{1 - a(z_0 + z^{(\alpha)})} \sigma_{\hat{\gamma}}, \quad \text{pak} \quad \gamma_d = \gamma^{(\frac{\alpha}{2})} \quad \text{a} \quad \gamma_h = \gamma^{(1-\frac{\alpha}{2})}.$$

Vezmeme-li v úvahu vztahy mezi distribučními funkcemi $F_{\hat{\theta}}$ a G , tj. (49), dospějeme analogicky jako v předchozím případě k vyjádření intervalu spolehlivosti pro θ , viz (52).

$$\begin{aligned}
 F_{\hat{\theta}}(t_d) &= G(g(t_d)) = G(\gamma_d) = \mathbf{P}_{\hat{\gamma}}(\hat{\gamma}^* \leq \gamma_d) = \mathbf{P}_{\hat{\gamma}}\left(\frac{\hat{\gamma}^* - \hat{\gamma}}{\sigma_{\hat{\gamma}}} + z_0 \leq \frac{\gamma_d - \hat{\gamma}}{\sigma_{\hat{\gamma}}} + z_0\right) \\
 (59) \quad &= \Phi\left(\frac{\gamma_d - \hat{\gamma}}{\sigma_{\hat{\gamma}}} + z_0\right) = \Phi\left(\frac{z_0 + z^{(\frac{\alpha}{2})}}{1 - a(z_0 + z^{(\frac{\alpha}{2})})} + z_0\right).
 \end{aligned}$$

Odtud

$$(60) \quad t_d^{(\frac{\alpha}{2})} = F_{\hat{\theta}}^{-1}\left(\Phi\left(\frac{z_0 + z^{(\frac{\alpha}{2})}}{1 - a(z_0 + z^{(\frac{\alpha}{2})})} + z_0\right)\right)$$

$$(61) \quad t_h^{(\frac{\alpha}{2})} = F_{\hat{\theta}}^{-1}\left(\Phi\left(\frac{z_0 + z^{(1-\frac{\alpha}{2})}}{1 - a(z_0 + z^{(1-\frac{\alpha}{2})})} + z_0\right)\right).$$

3.2.5. *Parametry z_0 a a .* Opravu vychýlení z_0 získáme úplně stejným způsobem jako v části 3.2.3, tj. viz (51). Tzv. „urychlovací konstantu“ a je možné (viz [1]) aproximovat výrazem:

$$(62) \quad a \doteq \frac{\text{SKEW}_{\theta=\hat{\theta}}(\hat{i}_{\theta})}{6} = \frac{\mathbf{E}(\hat{i}_{\theta} - \mathbf{E} \hat{i}_{\theta})^3}{6(\mathbf{E}(\hat{i}_{\theta} - \mathbf{E} \hat{i}_{\theta})^2)^{\frac{3}{2}}} \Big|_{\theta=\hat{\theta}}, \quad \text{kde} \quad \hat{i}_{\theta}(\hat{\theta}) = \frac{\partial}{\partial \theta} \log f_{\theta}(\hat{\theta}).$$

Poznámka 3.6. Analogicky bychom obdrželi i jednostranné $100(1 - \alpha)$ procentní intervaly spolehlivosti:

$$(63) \quad \left(-\infty; F_{\hat{\theta}}^{-1}\left(\Phi\left(\frac{z_0 + z^{(1-\alpha)}}{1 - a(z_0 + z^{(1-\alpha)})} + z_0\right)\right)\right)$$

$$(64) \quad \left(F_{\hat{\theta}}^{-1}\left(\Phi\left(\frac{z_0 + z^{(\alpha)}}{1 - a(z_0 + z^{(\alpha)})} + z_0\right)\right); +\infty\right)$$

Poznámka 3.7. Stanovení opraveného intervalu spolehlivosti, lze interpretovat jako opravu argumentu kvantilové funkce $F_{\hat{\theta}}^{-1}(\cdot)$, což je zřejmé ze srovnání výchozího intervalu spolehlivosti u neopravené percentilové metody, s výsledným. Např. jednostrannému intervalu

$$(65) \quad \left(\hat{F}_{\hat{\theta}}^{-1}(\alpha); +\infty\right)$$

odpovídá interval (64). Argumentem je u prvního z nich α , a ta se změní na

$$(66) \quad \Phi\left(\frac{z_0 + z^{(\alpha)}}{1 - a(z_0 + z^{(\alpha)})} + z_0\right), \quad \text{což označme jako } \Phi(z[\alpha]).$$

3.3. Aproximace. Poznámku 3.7 můžeme dále rozvést, uvědomíme-li si, že α v intervalu (65) je zároveň doplňkem jeho požadované spolehlivosti do jedné. Odtud $1 - \Phi(z[\alpha])$ odpovídá spolehlivostní hladině intervalu (64). Jak jsme již zmínili na počátku části 3.2 je náš původní problém opačný: k dané oblasti hledáme její spolehlivost. Vrátime-li se k analogii s příkladem 3.1, má tato oblast (ve zmíněném příkladu alternativa) tvar $(\pi_0; +\infty)$ a odpovídá tedy situaci (64). Proto:

$$(67) \quad \alpha_E \doteq 1 - \Phi(z[\alpha]) = \Phi(-z[\alpha]).$$

Zde narazíme na problém, jak odhadnout α . Nabízí se využít předchozí odhad hladiny spolehlivosti α_F , který umíme snadno spočítat. Ve shodě s výše zmíněným pozorováním o α by mělo platit:

$$(68) \quad \alpha \doteq 1 - \alpha_F.$$

Upravme dále:

$$(69) \quad \begin{aligned} z[\alpha] \doteq z[1 - \alpha_F] &= \frac{z_0 + z^{(1-\alpha_F)}}{1 - a(z_0 + z^{(1-\alpha_F)})} + z_0 \frac{z_0 - z^{(\alpha_F)}}{1 - a(z_0 - z^{(\alpha_F)})} + z_0 = \\ &= -\frac{z^{(\alpha_F)} - z_0}{1 - a(z^{(\alpha_F)} - z_0)} + z_0. \end{aligned}$$

Odtud:

$$(70) \quad \alpha_E \doteq \Phi\left(\frac{z^{(\alpha_F)} - z_0}{1 - a(z^{(\alpha_F)} - z_0)} - z_0\right),$$

kde Φ je distribuční funkce $N(0,1)$,

$$(71) \quad z^{(\alpha_F)} = \Phi^{-1}(\alpha_F)$$

$$(72) \quad a \quad z_0 = \Phi^{-1}\left(\mathbf{P}_{\pi_0}(\hat{\pi}^{**} \in \mathcal{P}_1)\right),$$

„Urychlovací člen“ a spočteme v následující části vhodnou úpravou výrazu (62) (viz 76).

3.4. Výpočet. (1) Pro rozsah B_1 spočteme algoritmem popsáním v 2.3 hladinu α_F pro danou větev. Odtud získáme $z^{(\alpha_F)}$ dle (71). Generování boot-populace ze sloupců matice \mathbf{X} by zároveň šlo popsat jako vygenerování

$$(73) \quad \mathbf{Y}_1^*, \dots, \mathbf{Y}_{B_1}^* \stackrel{iid}{\sim} \text{Multi}_n(n, \mathbf{p}), \quad \text{kde } \mathbf{p} = \frac{1}{n}(1, \dots, 1)'$$

Vektory $\mathbf{Y}^*(b)$ jsou analogií vektorů pozorovaných četností Υ (viz 8), v nichž se neuvažují nulové složky odpovídající „nepozorovaným“ vektorům ξ z prostoru \mathcal{X} . Mají tedy $n = 73$ složek. Stejným způsobem je svázán i vektor \mathbf{p} se svým protějškem π . Opět platí, že z vektoru pravděpodobností $\mathbf{p}_b^* = \mathbf{Y}_b^*/n$, kde $b = 1, \dots, B_1$ lze vypěstovat stromy $\hat{\Psi}_b^*$, podobně jako v (13).

(2) Zapamatujme si ty stromy, které se v dané větvi neshodují s původním stromem $\hat{\Psi}$. Označme příslušné pravděpodobnostní vektory, které je reprezentují, jako \mathbf{p}^j , $j = 1, \dots, J$. Víme o nich, že nemohou náležet do \mathcal{P}_1 , proto je můžeme využít k nalezení hraničních vektorů \mathbf{p}_0^j , analogických k π_0 . Nejprve zkusíme, zda bod ležící v polovině úsečky mezi body \mathbf{p} a \mathbf{p}^j patří do \mathcal{P}_1 . Pokud ano, nahradíme tímto bodem bod \mathbf{p} a pokračujeme s půlením této nové úsečky. Pokud ne, budeme

novým bodem nahrazovat bod \mathbf{p}^j . Pokračujeme-li takto u každého $j = 1, \dots, J$ až do požadované přesnosti (počet kroků označuji ve výsledcích parametrem L), obdržíme J -tici hraničních bodů \mathbf{p}_0^j .

(3) Nyní už nic nebrání nagenarovat pro $j = 1, \dots, J$:

$$(74) \quad \mathbf{Y}_{j1}^{**}, \dots, \mathbf{Y}_{jB_2}^{**} \stackrel{iid}{\sim} \text{Multi}_n(n, \mathbf{p}_0^j)$$

a zjistit počty stromů vypěstovaných z $\mathbf{Y}_{j1}^{**}, \dots, \mathbf{Y}_{jB_2}^{**}$ souhlasících v dané větvi s $\widehat{\Psi}$, označme je W_j . Průměrem z pozorovaných relativních četností W_j/B_2 odhadneme $\mathbf{P}_{\pi_0}(\widehat{\pi}^{**} \in \mathcal{P}_1)$, odtud dle (72)

$$(75) \quad z_0 = \Phi^{-1} \left(\frac{\sum_j W_j}{J \cdot B_2} \right).$$

(4) Poslední veličinou, která nám zbývá je urychlovací člen a . Závisí na směru vektoru spojujícího vektor \mathbf{p} a $\partial\mathcal{P}_1$. Položme tedy $\mathbf{u}_j = \mathbf{p}^j - \mathbf{p}$ pro $j = 1, \dots, J$

$$(76) \quad a(\mathbf{u}) = \frac{\sum_{k=1}^n u_k^3}{6 \left(\sum_{k=1}^n u_k^2 \right)^{\frac{3}{2}}}$$

Odvození tohoto vztahu lze nalézt v [1], jedná se vlastně o rozpis vztahu (62) pro případ exponenciálního systému hustot. Výrazy \mathbf{u}_j odpovídají empirickým influenčním funkcím pro odhad vektoru $\boldsymbol{\pi}$, tedy střední hodnoty rozdělení $\mathcal{L}(\widehat{\boldsymbol{\pi}} | \boldsymbol{\pi})$ – viz (18). V našem případě vezměme $a = J^{-1} \sum_j a(\mathbf{u}_j)$

(5) Dle vztahu (70) můžeme nyní dopočíst α_E .

3.5. Výsledky. V tabulce jsou větve číslovány vzestupně podle pořadí, v němž se slučovaly. Obě metody se výrazně liší v odhadu hladiny spolehlivosti především u větvi č. 7, 8 a 9, tj. u větvi, které, soudě dle nízké hodnoty α_F , mají při bootstrapu větší tendenci se přeskupovat.

č. větve	α_F			B_1	α_E		
	B	B	B		B_2	B	B
	2000	5000	10000		2000	2000	10000
	2000	5000	10000		40	200	50
1	0.796	0.811	0.891		0.727	0.754	0.751
2	1.000	1.000	1.000		1.000	1.000	1.000
3	0.958	0.958	0.954		0.912	0.904	0.878
4	0.848	0.842	0.841		0.787	0.765	0.768
5	0.666	0.673	0.670		0.737	0.724	0.723
6	0.752	0.736	0.739		0.745	0.739	0.716
7	0.297	0.318	0.308		0.734	0.745	0.752
8	0.357	0.347	0.357		0.740	0.750	0.742
9	0.387	0.386	0.386		0.725	0.713	0.726
10	0.672	0.641	0.653		0.718	0.696	0.705
11	0.630	0.618	0.620		0.712	0.714	0.706

Tab. 1 Srovnání hladin, $L = 20$

LITERATURA

- [1] Efron, B.: Better bootstrap confidence intervals. *JASA*, 82, 171–185, 1987.
- [2] Efron, B.: The Jackknife, the Bootstrap, and Other Resampling Plans. CBMS, 38, SIAM-NSF, 1982.
- [3] Efron, B., Halloran, E., Holmes, S.: Bootstrap confidence levels for phylogenetic trees. *Proc. Natl. Acad. Sci. USA*, 93, 13429–13434, 1996.
- [4] Escalante, A. A., Grebert, H. M., Chaiyaroj, S. C., Magris, M., Biswas, S., Nahlen, B. L., Lal A. A.: Polymorphism in the gene encoding the apical membrane antigen-1 (AMA-1) of *Plasmodium falciparum*. X. Asembo Bay Cohort Project. *Molecular & Biochemical Parasitology*, 113, 279–287, 2001.
- [5] Robert, C. P.: *The Bayesian Choice*. Springer-Verlag, New York, 1996.

KPMS MFF UK PRAHA, SOKOLOVSKÁ 83, 186 75 PRAHA 8 - KARLÍN
E-MAIL: betinec@karlin.mff.cuni.cz