

PŘESNÝ TEST NEZÁVISLOSTI V KONTINGENČNÍ TABULCE R x C

Miroslav Hartmann, Josef Bukač

Výpočetní středisko Lékařské fakulty UK
v Hradci Králové, Šimkova 870

1. Úvod

Tento postup byl vyvinut pro případy, kdy χ^2 -test dobré shody, používaný při testování nezávislosti v kontingenčních tabulkách, nevyhovuje pro velmi malé pozorované četnosti v některých polítkách. Metodu přesného výpočtu hladiny významnosti vyvinuli Freeman a Halton [1]. V této práci z roku 1951 se ještě nepředpokládá použití počítače a není ani zaměřena na jednoznačné stanovení algoritmu výpočtu. První program s přesným popisem postupu výpočtu publikoval March [2]. Při podrobném rozboru se ukazuje, že postup nevyužívá speciálních vztahů v kontingenčních tabulkách, které lze odvodit elementárními úvahami. Program je velmi jednoduchý a krátký, ale ne dostatečně účinný pro praktické použití. Některých speciálních vztahů pro tabulky typu 2x2 využívají Stucky a Vollmar [3]. Jimi sestrojený program pro tento typ tabulek je podstatně efektivnější, než Marchův.

Naše práce je zaměřena na obecné tabulky typu RxC. Vycházíme z Marchovy myšlenky, postup jsme podstatně upravili a tím urychlili několikanásobně vlastní výpočet.

2. STANOVENÍ PROBLÉMU

Kontingenční tabulkou typu RxC rozumíme matici celých nezáporných čísel

$$X = \begin{pmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,C} \\ \vdots & & & \\ x_{R,1} & x_{R,2} & \cdots & x_{R,C} \end{pmatrix} \quad (1)$$

Marginální součet i-tého řádku tabulky je definován vztahem

$$r_i = \sum_{j=1}^C x_{i,j} \quad (2)$$

obdobně marginální součet j-tého sloupce je definován vztahem

$$c_j = \sum_{i=1}^R x_{i,j} \quad (3)$$

Celkový počet pozorování v tabulce je dán výrazem

$$N = \sum_{i=1}^R \sum_{j=1}^C x_{i,j} = \sum_{i=1}^R r_i = \sum_{j=1}^C c_j \quad (4)$$

Pravděpodobnost výskytu tabulky X pro předem dané marginální součty r_i a c_j je při platnosti hypotézy nezávislosti dána výrazem

$$P(X) = \frac{\prod_{i=1}^R (r_i!) \prod_{j=1}^C (c_j!)}{N! \prod_{i=1}^R \prod_{j=1}^C (x_{i,j}!)} \quad (5)$$

Hladinu významnosti $P_s(X_0)$ testu nezávislosti v dané tabulce definujeme součtem

$$P_s(X_0) = \sum_{P(X) \leq P(X_0)} P(X) \quad (6)$$

Tuto definici vyjadříme slovně takto:

- a. zadaná tabulka X_0 určuje konkrétní hodnoty marginálních součtů r_i a c_j a hodnotu N .
- b. stanovíme pravděpodobnost $P(X_0)$ zadané tabulky.

- c. vytvoříme všechny možné tabulky \mathbf{X} , které vyhovují daným marginálním součtům a pro každou z nich vypočteme její pravděpodobnost $P(\mathbf{X})$.
d. hladina významnosti je pak dána jako součet všech pravděpodobností $P(\mathbf{X})$, pro které platí $P(\mathbf{X}) \leq P(\mathbf{X}_0)$.

Vzhledem k tomu, že hodnoty r_i , c_j , N jsou ve všech vytvářených tabulkách totálně, je vhodné přepsat vzorec (5) do tvaru

$$P(\mathbf{X}) = e(\mathbf{X}_0)/R(\mathbf{X}), \quad (7)$$

$$e(\mathbf{X}_0) = \prod_{i=1}^R (r_i!) \prod_{j=1}^C (c_j!) / (N!) \quad (8)$$

$$R(\mathbf{X}) = \prod_{i=1}^R \prod_{j=1}^C (x_{i,j}!). \quad (9)$$

3. GENEROVÁNÍ TABULEK

3.1 Marchova metoda

Jelikož jednotlivé hodnoty $x_{i,j}$ musí vyhovovat rovnici (2) a (3), je tabulka typu $R \times C$ dána $(R-1) \times (C-1)$ prvků, přičemž zbývající prvky jsou z této rovnice dopočítány. Připustnou tabulkou nazveme takovou tabulku, která splňuje rovnice (2), (3) a jejíž všechny prvky jsou nezáporné. První tabulku vytváří March tak, že pro $(R-1) \times (C-1)$ prvků vpravo dole, tj. s vynecháním prvního sloupce a prvního řádku, dosadí minimum z příslušných sloupcových a řádkových marginálních součtů. Po vypočtení prvků v prvním sloupci a řádku podle rovnic (2) a (3) zjistí, zda tabulka je připustná. Generování další tabulky vychází vždy z vygenerované tabulky předchozí. V této tabulce hledá postupně v prvcích $x_{2,2}, \dots, x_{2,C}, x_{3,2}, \dots, x_{3,C}, \dots, x_{R,C}$ první kladný prvek. Od tohoto prvního odečte jedničku a do všech předcházejících, které byly nulové, dosadí minimum z příslušných marginálních součtů. Opět vypočítá prvky v prvním sloupci a řádku a zjistí, zda vzniklá tabulka je připustná. Následuje opakováný návrat do vyhledávání kladného prvku. Celý postup je ukončen, jestliže není nalezen žádný další prvek.

V podstatě tento postup představuje $(R-1) \times (C-1)$ násobný cykl s pevnýmimezemi a ukazuje se, že obecně je při něm generován velký počet nepřipustných tabulek.

3.2 Nová metoda

Jde o způsob, podobný metodě severozápadního rohu, používané pro nalezení základního řešení lineárního dopravního problému. Z hlediska dopravního problému lze hodnoty marginálních součtů r_i , resp. c_j zadané tabulky \mathbf{X}_0 považovat za kapacity řádků, resp. sloupců tabulky. Toho využijeme následujícím způsobem:

Před započetím generování naplníme první řádek tabulky hodnotami c_j z tabulky \mathbf{X}_0 , tj. na prvním řádku jsou umístěny kapacity sloupců ($x_{1,j} = c_j$). Obdobně do prvního sloupce, počínaje druhým řádkem, umístíme hodnoty r_i z tabulky \mathbf{X}_0 , tj. kapacity řádků ($x_{i,1} = r_i$ pro $i=2,3,\dots,R$). Ostatní prvky tabulky nechť jsou nulové. Tato vzniklá tabulka splňuje rovnice (2) pro všechny řádky a výjimkou prvního a rovnice (3) pro všechny sloupce s výjimkou prvního.

Při generování jakéhokoliv prvku $x_{i,j}$ ($i > 1, j > 1$) v tabulce dbáme vždy této zásady: Jestliže prvek zvětšíme, resp. zmenšíme o určitou hodnotu, pak o tutéž hodnotu zmenšíme, resp. zvětšíme i prvky $x_{i,1}$ a $x_{1,j}$. Tím je tedy neustále zachována platnost rovnic (2) a (3) kromě prvního řádku a prvního sloupce.

První tabulku generujeme tak, že do $x_{R,C}$ dosadíme minimum z $x_{R,1}$ a $x_{1,C}$. Po dosazení okamžité hodnoty $x_{R,1}$ a $x_{1,C}$ upravíme. Obdobně do prvku $x_{R,j}$ ($j=C-1, C-2, \dots, 2$) dosadíme minimum z hodnot $x_{R,1}$ a $x_{1,j}$. Po tomto procesu již nebudeme považovat hodnotu $x_{R,1}$ za volnou kapacitu R-tého řádku, ale za hodnotu, která byla

v tomto prvku vygenerována. O tuto hodnotu $x_{R,1}$ snížime prvek $x_{1,1}$ /tj. kapacitu prvního sloupce/. Stejně jako v R-tém řádku postupujeme i v řádcích vyšších /tj. $i=R+1, R+2, \dots, 2$. Po vygenerování celého druhého řádku budeme také hodnoty, které zůstaly v prvním řádku, považovat za vygenerované. Tato první vygenerovaná tabulka je vždy přípustná /analogie metody severozápadního rohu u dopravního problému/ a navíc je totožná s první přípustnou tabulkou, generovanou Marchovou metodou.

Při generování dalších tabulek vycházíme podobně jako March vždy z předcházející vygenerované tabulky. V prvcích $x_{2,2}, \dots, x_{2,C}, x_{3,2}, \dots, x_{3,C}, \dots, x_{R,2}, \dots, x_{R,C}$ hledáme první kladný prvek. Navíc před započetím prohledávání i-tého řádku přičteme hodnotu $x_{1,1}$ k hodnotě $x_{1,i}$. Od prvního kladného prvku, řekněme $x_{1,j}$, odečteme jedničku a navíc tuto jedničku přičteme k prvkům $x_{1,1}$ a $x_{1,j}$. V pořadí opačném k pořadí prohledávání znovu generujeme již popsáným způsobem /jako u první tabulky/ hodnoty prvků, předcházejících prvku $x_{1,j}$.

Ze způsobu generování plyne, že přípustnost vygenerované tabulky je závislá pouze na hodnotě prvku $x_{1,1}$. Tabulka je přípustná právě tehdy, když $x_{1,1} \geq 0$.

Vzhledem k tomu, že hodnota $x_{1,1}$ během generování tabulky představuje volnou kapacitu prvního sloupce, potom jestliže po vygenerování i-tého řádku ($i \geq 2$) je hodnota $x_{1,1}$ větší než hodnota $x_{1,i}$ a tedy po odečtení by hodnota $x_{1,1}$ byla záporná, pak až je ve vyšších řádcích tímto postupem generováno cokoliv, výsledná tabulka nebude nikdy přípustná. V takovém případě není tedy třeba zbytek tabulky generovat a hodnotu $x_{1,1}$ od hodnoty $x_{1,1}$ neodečteme. Můžeme vyhledat první kladný prvek, počínaje prvkem $x_{1,2}$ a od něj odečíst jedničku. Jestliže však odečteme jedničku od prvku $x_{1,2}$ /je-li to možné/, zvýšíme tím hodnotu $x_{1,1}$ a po odečtení této hodnoty od $x_{1,1}$ /kapacity 1. sloupce/ získáme ještě menší hodnotu $x_{1,1}$ a tudíž ani tato tabulka nebude přípustná. Analogicky můžeme uvážit:

Byla-li po vygenerování i-tého řádku přečerpána kapacita 1. sloupce a byly-li vyčerpány kapacity 2. až j-tého sloupce, pak jakékoli snížení hodnot v 2. až (j+1). prvku i-tého řádku nevede k přípustné tabulce. V důsledku toho můžeme před vyhledáváním prvního kladného prvku vrátit hodnoty prvků $x_{1,2}$ až $x_{1,j+1}$ do jím příslušných kapacit a tyto prvky vynulovat. Po odečtení jedničky v prvním kladném prvku a úpravách patřičných kapacit opět začneme generovat předcházející prvky dříve popsáným způsobem.

Efektivnost těchto úvah vyplývá z následujícího tvrzení: Jestliže po vygenerování i-tého řádku nebyla přečerpána kapacita 1. sloupce, pak jistě existuje alespoň jedna tabulka, která je přípustná a prvky v řádcích i-tém až R-tém má totožné s dosud vygenerovanými. Toto tvrzení plyne z analogie s dopravním problémem.

Vytváření všech tabulek je ukončeno tehdy, když všechny prvky $x_{i,j}$ ($i > 1, j > 1$) jsou nulové a nelze tedy nikde provést odečtení jedničky.

3.3 Postupy užitečné při programování nové metody

3.3.1 Odlišné generování u prvku $x_{2,2}$

Máme vygenerovanou přípustnou tabulku a nechť $x_{2,2} > 0$. Při generování další tabulky bychom tedy odečetli jedničku od prvku $x_{2,2}$. Musíme pak přičist jedničku k prvkům $x_{1,2}$ a $x_{2,1}$ a odečist jedničku od prvku $x_{1,1}$. Opakování použití tohoto postupu neneruší přípustnost tabulky v tom případě, kdy $x_{1,1} \geq x_{2,2}$. Toho lze využít k naprogramování těchto změn v jednoduchém cyklu. Hodnoty v prvcích o indexech (1,1), (1,2), (2,1) a (2,2) budou dány výrazy

$$x_{1,1} = K, x_{1,2} + K, x_{2,1} + K, x_{2,2} = K \quad (10)$$

kde $K = 0, 1, \dots, \min(x_{1,1}, x_{2,2})$.

Po tomto cyklu přičteme hodnotu $x_{1,1} - K$ a $x_{2,2} - K$ k hodnotám $x_{1,1}$ a $x_{2,2}$.

* $x_{2,1}$ - Prvek $x_{2,2}$ vynulujeme a teprve nyní začneme hledat kladný prvek.

3.3.2 Uspořádání marginálních součtů

Ze vzorců (5) a (7) vyplývá, že jsou invariantní k záměně řádků a sloupců zadáné tabulky. Tím tedy ani výpočet hladiny významnosti (6) nezávisí na uspořádání marginálních součtů, pěsto, že přípustné tabulky mohou být generovány v jiném pořadí. Vzhledem k algoritmu generování tabulek a vzhledem k odstavci 3.3.1 se jeví nejvhodnější uspořádání marginálních součtů sestupně podle velikosti. To je výhodné zejména proto, že algoritmus zajíšťuje rychlejší změny v prvcích o nižších řádkových a v řádku o nižších sloupcových indexech.

3.4 Zhodnocení nové metody

Nová metoda zaručuje generování všech přípustných tabulek, navíc jsou generovány tyto přípustné tabulky ve stejném pořadí, v jakém jsou generovány Marchovou metodou. Proti Marchové metodě zde vzniká podstatně menší počet nepřípustných tabulek, navíc nepřípustnost generované tabulky můžeme zjistit dříve, než March a proces generování zastavit.

Není třeba vypočítávat prvky v 1. řádku, resp. 1. sloupci pomocí rovnic (2), resp. (3) a během jejich vypočtu kontrolovat přípustnost, jako to dělá March. Zde se stačí omezit na prvek $x_{1,1}$ a to v libovolné fázi výpočtu.

Sestupné uspořádání marginálních součtů snižuje podstatně dobu výpočtu a při praktických výpočtech se ukázalo, že vede i ke snížení počtu nepřípustných tabulek. Toto uspořádání by mohlo urychlit i Marchovu metodu.

4. VÝPOČET HLAĐINY VÝZNAMNOSTI

4.1 Výpočet pravděpodobnosti tabulky

Pravděpodobnost vygenerované tabulky se počítá podle vzorců (7), (8) a (9). Při praktickém využití je vhodné používat logaritmů, což zabrání přetečení. V Marchové metodě je výraz $\ln(R(X))$ počítán po každé vygenerované přípustné tabulce celým. Ukažuje se, že je účelné zavést tabulku S součtu logaritmů faktoriálů, jejíž využití urychlí výpočet $\ln(R(X))$. Tabulka S má také rozměr $R \times C$ a je svázána s tabulkou X tak, že prvek $s_{i,j}$ obsahuje součet logaritmů faktoriálů všech prvků, lexicograficky následujících za prvkem $x_{i,j}$. Např. prvek $s_{3,C}$ obsahuje součet logaritmů faktoriálů všech prvků tabulky X od čtvrtého řádku počínaje.

Výpočet prvků tabulky S se provádí současně s generováním tabulky X . Jestliže před generováním další tabulky X byla odečtena jednička na prvku $x_{i,j}$, potom tabulka S bude pře počítána jen v prvcích lexicograficky předcházejících prvku $s_{i,j}$.

4.2 Výpočet hladiny významnosti pomocí doplňku

Pravděpodobnosti $P(X)$ přípustných tabulek, které byly vypočteny v průběhu generování, jsou vždy porovnány s pravděpodobností $P(X_0)$ zadané tabulky X_0 a jestliže platí $P(X) \leq P(X_0)$, jsou přičteny do proměnné, která po skončení generování dává $P_s(X_0)$.

Tímto způsobem počítá hladinu významnosti March. Rovněž počítá doplněk k hladině významnosti podle vzorce

$$P_s^D(X_0) = \sum_{P(X) > P(X_0)} P(X) \quad (11)$$

Součet $P_s(X_0) + P_s^D(X_0)$ by se měl, až na zaokrouhlovací chybu rovnat jedné. Tchoto vztahu používá pro kontrolu přesnosti výpočtu. Praktické výpočty na ODKE 1204 ukázaly, že přesnost je dostatečně vysoká a tato kontrola je zbytečná. Lze tedy bez

podstatné chyby počítat hladinu významnosti i ze vztahu

$$P_s(X_0) = 1 - P_s^D(X_0) \quad (12)$$

V nové metodě nepočítáme $P_s(X)$ podle vztahu (6) jako March, ale podle vztahu (12). Omezili jsme se tedy jen na ty tabulky X , pro které $R(X) > R(X_0)$. ze vztahu (7) vyplývá, že u těchto tabulek platí $R(X) < R(X_0)$.

5. PRAVDĚPODOBNOSTNÍ OMEZENÍ PŘI GENEROVÁNÍ

5.1 Pravděpodobnostní omezení

Jestliže se při výpočtu hladiny významnosti řídíme vztahem (12), omezíme se jen na tabulky X , pro které $R(X) < R(X_0)$. $R(X_0)$ je tedy horní mezí hodnot $R(X)$ pro všechny tabulky X , které je nutné generovat. Ostatní tabulky, které mají $R(X) \geq R(X_0)$, nejsou k výpočtu hladiny významnosti potřebné. Je tedy možno je z generování vypustit.

Vzhledem k tomu, že $\ln(R(X))$ je v nové metodě načítán průběžně v tabulce S , je možné po vygenerování jakéhokoliv prvku $x_{i,j}$ zjistit, zda $s_{i,j} + \ln(x_{i,j}!)$ je větší než $\ln(R(X_0))$. Jestliže ano, potom již není nutné takovou tabulku dále generovat, neboť pro ni bude jistě platit $R(X) > R(X_0)$, ať do zbývajících prvků vygenerujeme cokoliv.

5.2 Funkce MSF

Pro potřeby dalšího výkladu udělejme tuto úvahu. Máme-li rozdělit K jednotek do N poliček tak, aby součin faktoriálů / resp. součet logaritmů faktoriálů / počtu jednotek v těchto poličkách byl minimální, lze jednoduše ukázat, že to bude v tom případě, když budou rozděleny co nejrovnoměrněji / tj. počty v libovolných dvou poličkách se budou lišit nejvýše o jednotku/. K tomu si definujme funkci MSF (K, N), jejíž hodnotou bude minimální součet logaritmů faktoriálů při K jednotkách, rozdělených do N poliček.

5.3 Řádková a sloupcová pravděpodobnostní omezení

Podle úvahy v předchozím odstavci 5.2 můžeme pro jednotlivé řádky tabulek určit předem dolní mez soutě logaritmů faktoriálů prvků v těchto řádcích. Tato dolní mez je pro i -tý řádek dána hodnotou funkce MSF (r_i, c). Můžeme tedy pro i -tý řádek určit horní mez RB_i soutě logaritmů faktoriálů prvků v tomto řádku a v řádcích s indexem vyšším pomocí vztahu

$$RB_i = \ln(R(X_0)) - \sum_{k=1}^{i-1} MSF(r_k, c) \quad (13)$$

Hodnoty RB_i je možno vypočítat před započetím generování, neboť závisí pouze na $R(X_0)$, r_i a c . Omezení stanovené v bodě 5.1 je možné ještě zosavit tímto způsobem:

Jestliže po vygenerování prvku $x_{i,j}$ zjistíme, že

$$s_{i,j} + \ln(x_{i,j}!) > RB_i \quad (14)$$

není nutné tabulku dále generovat. Ať je do zbytku tabulky dogenerováno cokoliv, bude vždy $R(X) > R(X_0)$ pro každou přípustnou tabulku X , která má v prvku $x_{i,j}$ a v prvcích lexikograficky následujících hodnoty dosud vygenerované.

Obdobně jako pro řádky je možno uplatňovat během výpočtu i omezení pro sloupce, neboť zbývající sloupcové kapacity je nutno určitým způsobem do dosud nevygenerovaných prvků rozdělit. Je zřejmé, že pokud

$$x_{1,j} + \ln(x_{1,j}!) > \ln(R(X_0)) - \sum_{k=1}^{j-1} \text{MSF}(x_{1,k}, i) - \sum_{k=j}^c \text{MSF}(x_{1,k}, i-1), \quad (15)$$

je zbytečné tabulku degenerovávat.

Tato omezení však již závisí na volných kapacitách a na indexech i, j a není proto možné je vypočítat předem.

Průběžná kontrola podle nerovnosti (15) by časově příliš ztížila proces generování a proto jsou předchozí úvahy prakticky realizovány tak, že nerovnost (14) je kontrolována při generování každého prvku. Teprve když je splněna, jsou prováděny kontroly jak pro sloupcové omezení (15), tak i pro zastřené řádkové omezení, dáné nerovností

$$x_{1,j} + \ln(x_{1,j}!) > \ln - \text{MSF}(x_{1,1}, j-1). \quad (16)$$

Pokud je daný prvek $x_{1,j}$ splněna alespoň jedna z nerovnosti (15) a (16) a $x_{1,j} > 0$, můžeme od $x_{1,j}$ jedničku a provedeme příslušné úpravy kapacit $x_{1,1} \dots x_{1,j-1}$. Odík kontrolované splnění nerovnosti (15) a (16). Jestliže již $x_{1,j} = 0$ a dosud alespoň jedna z nerovnosti (15) a (16) platí, pokračujeme v obdobném testování první lexicografické následujícího. Tento zpětný proces je ukončen, jestliže obě nerovnosti nejsou splněny /pak se pokračuje v generování/, případně tehdy, když jsme v tomto procesu dostali k poslednímu prvku tabulky $X_{R,C}$ a tento je již nulový, přičemž alespoň jedna z nerovnosti (15) a (16) platí /pak je generování tabulek u konce/. Při přechodu na nízší řádek je v tomto procesu samozřejmě nutné upravit i kapacitu 1. sloupcu /k hodnotě $x_{1,1}$ přičist hodnotu $x_{1,1} /$.

Toto strategie generování vyžaduje složitější program, ale nezatěžuje zbytečnými časovými ztrátami výpočty pro tabulky s nízkou hladinou významnosti.

5.4. Úprava generování u prvku $x_{2,2}$

Počet uvedený v bodě 3.3.1 můžeme při využívání pravděpodobnostních omezení v generování upnout následujícím způsobem:

Množina těch K , které vystupují v (10) můžeme rozdělit na dvě podmnožiny, pro které se $R(X)$ chovají monotoně.

$$M_1 = \{K; K \leq L\},$$

$$M_2 = \{K; K > L\}.$$

kde

$$L = \frac{x_{1,1} + x_{2,2} - (x_{1,2} + 1) \cdot (x_{2,1} + 1)}{x_{1,1} + x_{1,2} + x_{2,1} + x_{2,2} + 2} \quad (17)$$

V množině M_1 se pro rostoucí K $R(X)$ snižuje, v M_2 naopak zvyšuje. Je možné, že jedna z množin je prázdná.

Výpočet pak bude muset provádět raději ve dvou cyklech, odpovídajících množinám M_1 a M_2 . V M_1 postupně proklesající K , v M_2 pro rostoucí. Nemusíme pak zbytečně počítat v tabulkách: $\leq R(X) > R(X)$.

V případě, že nebyla v obou cyklech nalezena žádná tabulka X , pro kterou $R(X) \leq R(X)$, rozšiřujeme přímo v generování dalších tabulek, ale počínaje prvkem $x_{2,2}$ a využívame testování vzhledem k pravděpodobnostním omezením podle bodu 5.3.

5.5 Zhodnocení metody generování při pravděpodobnostních omezeních

Metoda popsána v kapitole 5 má proti metodě generování všech přípustných tabulek z kapitoly 3 tu přednost, že z procesu generování vynechává tabulky pravděpodobnostně nevhodné, tj. tabulky, pro které $P(X) \leq P(X_0)$. Z toho plynou, že efektivnost této metody se projeví podstatně předeším v těch případech, kdy $P(X_0)$ bude co největší vzhledem k pravděpodobnostem ostatních přípustných tabulek. Algoritmus je ovšem navržen tak, aby nezdržoval výpočet hladiny významnosti u tabulek, pro které $P(X_0)$ je malé.

Vezmeme-li v úvahu výsledky kapitoly 3 a této kapitoly, je zřejmé, že požadovanou hladinu významnosti $P_s(X_0)$ získáme v podstatně kratším čase, než March.

6. PROCEDURA CONP

6.1 Popis procedury

Nová metoda, jejíž popis byl uveden v předcházejících kapitolách 3 až 5, byla zpracována do algolské procedury. Hlavíčka procedury má tento tvar:

```
PROCEDURE CONP (MATRIX, R, C, PT, PS);
```

```
VALUE MATRIX, R, C;
```

```
INTEGER R, C;
```

```
REAL PT, PS;
```

```
INTEGER ARRAY MATRIX;
```

```
END CONP;
```

Parametry mají tento význam:

MATRIX Dvourozměrné celočíselné pole s rozsahem $[1:R, 1:C]$, které obsahuje zadanou kontingenční tabulku X_0 .

R Počet řádků tabulky X_0 .

C Počet sloupců tabulky X_0 .

PT Výstupní parametr, při návratu z procedury obsahuje pravděpodobnost $P(X_0)$.

PS Výstupní parametr, při návratu z procedury obsahuje hladinu významnosti $P_s(X_0)$.

Procedura pracuje s nelokálními proměnnými typu **ARRAY**. Je to jednorozměrné pole $FL [0:\alpha]$, jehož prvky obsahují logaritmy faktoriálů, tedy $FL [i] = \ln(i!)$ a dvourozměrné pole $MSFL [0:\beta, 1:\gamma]$, kde jednotlivé prvky tohoto pole obsahují minimální součet logaritmů faktoriálů, tedy $MSFL [K,N] = MSF (K,N)$, viz fce **MSF**, definovaná v bodě 5.2. Horní meze α, β, γ v deklaracích těchto polí je nutno volit s ohledem na rozměry tabulky a potřhy pozorování v tabulce. Hodnota α musí být větší nebo rovna celkovému počtu pozorování v tabulce X_0 , hodnota β větší nebo rovna největšímu z řádkových a sloupcových marginálních součtů, hodnota γ větší nebo rovna maximu z hodnot R a C. Pokud tyto podmínky nejsou dodrženy, může dojít k chybě při práci procedury /SUBSCRIPT/. Hodnoty do polí **FL** a **MSFL** je nutno pochopitelně dosadit před vyvoláním procedury **CONP**. Prvky nelokálních polí **FL** a **MSFL** je možno samozřejmě nahradit funkčními procedurami, které poskytnou požadované hodnoty, to ale podstatně prodlouží dobu výpočtu.

6.2 Doba výpočtu

Odhad doby výpočtu nelze v jednotlivých případech stanovit. Lze jen uvést určité fakta, ze kterých je možné usuzovat na složitost a tedy i délku výpočtu příkladu. Délka výpočtu závisí obecně na počtu přípustných tabulek, u kterých $P(X) > P(X_0)$. obtížnost je zde záleží především neznalosti celkového počtu přípustných tabulek pro dané marginální součty a neznalosti rozvrstvení jejich pravděpodobnosti. Vzorec pro výpočet počtu přípustných tabulek nám není znám, ze zkušenos-

tí však víme, že tento počet závisí především na velikosti tabulky /jejích rozměrech/ a dále na velikosti a rozdělení marginálních součtů v řádcích a sloupcích. Čím rovnoměrnější jsou marginální součty v řádcích, resp. sloupcích rozděleny, tím větší počet přípustných tabulek můžeme očekávat. Také se vzrůstajícím počtem pozorování v tabulce počet přípustných tabulek narůstá.

Z praktického hlediska, jak ukázaly výpočty zkušebních příkladů, je vhodné zadávat tabulku χ^2 -výdy tak, aby $R \geq C$. V tomto případě dochází k lepšímu využití pravděpodobnostních omezení RB, v jednotlivých řádcích, proces generování pravděpodobnostně nevhodných tabulek ($P(X) \leq P(X_0)$) je zastavován dříve.

7. APLIKACE METODY

7.1 Mohoucí χ^2 -testu dobré shody

Je známo, že není vhodné používat χ^2 -testu, jestliže napozorované četnosti jsou nízké. Při srovnání s přesným testem se ukazuje, že testové kriterium χ^2 není vhodné užívat ani jako přibližnou míru závislosti v tabulce. Proti jeho použití svědčí dvě skutečnosti. Jednak to, že výsledky získané χ^2 -testem se mohou značně lišit od výsledků získaných přesným testem a to tím více, čím menší jsou napozorované četnosti. Druhá skutečnost je ta, že výsledek přesného testu závisí jen na pravděpodobnosti a marginálních součtech zadané tabulky, zatímco χ^2 je závislý ne na pravděpodobnosti $P(X_0)$, ale konkrétních napozorovaných četnostech. Pro dvě různé tabulky se stejnou pravděpodobností a se stejnými marginálními součty dává přesný test tutéž hladinu významnosti, zatímco χ^2 -test může dát hodnoty $P_g(X_0)$ značně odlišné. To se projevuje podstatně u tabulek s nízkými četnostmi.

S ohledem na to, že doba výpočtu hladiny významnosti $P_g(X_0)$ při použití přesného testu rychle roste s rostoucími rozměry tabulky a s rostoucími četnostmi pozorování, nelze běžně počítat s jeho hromadným použitím jako náhradou za běžně používaný χ^2 -test. Přesného testu lze však použít v těch případech, kdy podmínky pro aplikaci χ^2 -testu nejsou dobré splněny.

7.2 Dopravní problém

Generování všech přípustných tabulek, které je součástí metody popsane v kapitole 3, je z hlediska matematického programování výčtem všech přípustných řešení celočíselného dopravního problému. Toho lze využít při minimalizaci nákladů v dopravních problémech, jejichž účelová funkce není lineární. Je-li účelová funkce dopravního problému taková, že její předpis je shodný pro všechna přípustná řešení a příspěvek každého prvku matice přepravovaných množství je nezávislý na hodnotách jiných prvků této matice a přitom je nezáporný, můžeme při minimalizaci použít budebných úvah, jako v bodě 4.1 a kapitole 5 této práce. Navíc omezení, analogické pravděpodobnostnímu omezení uvedenému v bodě 5.1 lze upravovat během výpočtu podle dosud zjištěného minima účelové funkce.

8. Závěry

Pro ilustraci uvádíme dva příklady, na nichž bude demonstrováno srovnání uvedených metod výpočtu hladiny významnosti. Jednotlivé metody označme takto:

1. Výpočet pomocí χ^2 -testu.
2. Výpočet Marchovou metodou.
3. Jako 2., ale pro transponovanou tabulku.
4. Výpočet novou metodou generování tabulek bez pravděpodobnostních omezení (generovány všechny přípustné tabulky).
5. Jako 4., ale pro transponovanou tabulku.
6. Výpočet novou metodou generování tabulek s pravděpodobnostními omezeními.
7. Jako 6., ale pro transponovanou tabulku.

Příklad 1.

Tabulka typu 3×3 **metoda** **hladina významnosti** **počet generovaných přípust. tabulek** **doba výpočtu s**

2	3	1	1	0.53831	-	-
0	1	2	2	0.71218	278	37
6	5	5	3	0.71218	278	39
			4	0.71218	278	8
			5	0.71218	278	9
			6	0.71218	13	1
			7	0.71218	13	2

Příklad 2.

Tabulka typu 3×4 **metoda** **hladina významnosti** **počet generovaných přípust. tabulek** **doba výpočtu s**

3	7	4	8	1	0.46631	-	-
5	2	3	3	2	?	?	3600
2	1	0	3	3	?	?	3600
				4	0.50790	31892	762
				5	0.50790	31892	768
				6	0.50790	818	70
				7	0.50790	818	53

Výpočet metodou 2 a 3 byl zastaven po jedné hodině bez získání výsledků.

LITERATURA:

1. Freeman, G. H., Halton, J. H.: Note on an exact treatment of contingency, goodness of fit, and other problems of significance. Biometrika, 38(1951), 141-149.
2. March, D. L.: Exact probabilities for $R \times C$ contingency tables. Algorithm 434. Communications of the Association for Computing Machinery 15(1972) 11, 991-992.
3. Stucky, W., Vollmar, J.: Ein Verfahren zur exakten Auswertung von $2 \times C$ - Häufigkeitstafeln. Biometrische Zeitschrift 17(1975) 3, 147-162.