

# Aplikace matematiky

---

Zdeněk Renc

Automatická binarisace veličin

*Aplikace matematiky*, Vol. 15 (1970), No. 2, 117–124

Persistent URL: <http://dml.cz/dmlcz/103275>

## Terms of use:

© Institute of Mathematics AS CR, 1970

Institute of Mathematics of the Czech Academy of Sciences provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This document has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

AUTOMATICKÁ BINARISACE VELIČIN<sup>1)</sup>

ZDENĚK RENC

(Došlo dne 11. října 1968)

## I

Experimentální pracovník libovolného oboru (na př. lékařského biologického, sociologického atd.) se patrně dosti často ocitne v následující situaci. Shromáždil experimentální materiál skládající se z nějakých objektů (na př. pacientů, obyvatel nějaké oblasti, fyzikálních měření) a různých dat, která o těchto objektech zjistil. Pro svou další práci potřebuje zjišťovat různé zákonitosti a souvislosti mezi těmito daty. Abychom tuto situaci mohli lépe popsat, předpokládejme, že máme danou množinu  $M$  jakýchkoli  $m$  objektů a že jsou dále dány binární vlastnosti  $P_1, \dots, P_n, A_1, \dots, A_l$  a veličina  $V$ , při čemž pro každý objekt z naší množiny  $M$  je známo jednak zda splňuje po řadě vlastnosti  $P_1, \dots, A_l$ , jednak jakou hodnotu mu připisuje veličina  $V$ . Předpokládejme dále, že  $P_1, \dots, P_n$  jsou zadány tak, že každé  $x \in M$  splňuje právě jednu z nich a že každá z vlastností  $P_1, \dots, P_n$  je splněna alespoň jedním objektem množiny  $M$ . (Např.  $P_1, \dots, P_{n-1}$  mohou být různé vzájemně se vylučující druhy téže choroby,  $P_n$  znamená „býti zdrav“.) Máme zkoumat vztahy mezi  $P_1, \dots, P_n$  na jedné straně a  $A_1, \dots, A_l$  (což v našem příkladě mohou být nějaké další vlastnosti, které jsme u pacientů zjistili) a veličinou  $V$  (např. věkem pacienta) na straně druhé. Tyto vztahy je možno zjišťovat na příklad pomocí metody GUHA, popsané v pracích [2], [3], [4], na které čtenáře odkazují.

Abychom mohli metodu GUHA použít, musíme ovšem veličinu  $V$  binarizovat nebo, chceme-li, nahradit ji systémem binárních veličin  $V_1, \dots, V_k$  (to znamená, že otázku „jaký je věk pacienta“ potřebujeme nahradit systémem otázek „je věk pacienta v rozmezí  $L_1$  až  $L_2$  let“ atd.) Aby byla binarisace veličiny  $V$  rozumná, požadujeme, aby opět každý prvek  $x \in M$  splňoval právě jednu z vlastností  $V_1, \dots, V_k$  a aby každá z  $V_1, \dots, V_k$  byla splněna alespoň jedním objektem množiny  $M$  (to znamená, aby intervaly, na které jsme věkovou škálu rozdělili, se navzájem nepřekrývaly a aby dohromady pokryly celou tuto škálu).

Kromě této podmínky musíme mít ale na mysli ještě další požadavek, kterému můžeme říkat požadavek vhodnosti binarizace veličiny  $V$  vzhledem k  $P_1, \dots, P_n$ .

<sup>1)</sup> Tato práce byla referována na semináři aplikací matematické logiky na MFF KU a přednesena na XXIII. Fysiologických dnech, konaných v Železném Rudě.

Nechceme totiž binarizovat veličinu  $V$  jakkoliv, nýbrž tak, aby tato binarizace byla co nejvýhodnější v tom smyslu, že nám pomůže rozdělit obor hodnot veličiny  $V$  na částečné intervaly co nejlépe určující rozdělení prvků množiny  $M$  vzhledem k vlastnostem  $P_1, \dots, P_n$ . Tuto vhodnost binarizace veličiny  $V$  vzhledem k vlastnostem  $P_1, \dots, P_n$  budeme měřit způsobem, o kterém se čtenář může poučit v [1] a který nyní stručně popíšeme.

Představme si následující situaci. Máme dány binární vlastnosti  $P_1, \dots, P_n$  (s požadavky, o nichž jsme výše hovořili). Vybereme náhodně prvek  $x \in M$  a snažíme se co nejlépe „uhádnout“ tu vlastnost z  $P_1, \dots, P_n$ , kterou  $x$  splňuje a to nejprve v první a potom v druhé z následujících dvou situací

1. Není nám nic známo o vlastnostech  $V_1, \dots, V_k$ .

2. Je nám znám rozklad  $V$  na systém binárních vlastností  $V_1, \dots, V_k$ .

1. Necht'  $a_{.1}, \dots, a_{.n}$  jsou po řadě počty prvků z  $M$ , splňujících  $P_1, \dots, P_n$ ,  $a_{.max} = \text{Max } a_{.j}$ . Nejlepší odhad dostaneme, zvolíme-li tu vlastnost  $P_j$ , pro níž je  $a_{.max} = a_{.j}$ . Předpokládáme-li, že pravděpodobnosti jevu „objekt  $x \in M$  splňuje vlastnost  $P_j$ “ je rovna  $p_{.j} = a_{.j}/m$  (to znamená, že tato pravděpodobnost je rovna relativní frekvenci  $P_j$  v množině  $M$ ), můžeme říci, že nejlepší odhad dostaneme, zvolíme-li tu třídu  $P_j$ , pro níž je  $p_{.max} = p_{.j}$  (kde  $p_{.max} = \text{Max } p_{.j}$ ). Chyba tohoto odhadu je  $1 - p_{.max}$ .

2. Necht'  $a_{i1}, \dots, a_{in}$  (pro  $i = 1, \dots, k$ ) jsou po řadě počty prvků z  $M$ , splňujících  $V_i$  a současně  $P_1, \dots, P_n$ ,  $a_{imax} = \text{Max } a_{ij}$  pro  $i = 1, \dots, k$ . Necht'  $p_{i1}, \dots, p_{in}$  ( $i = 1, \dots, k$ ),  $p_{imax} = \text{Max } p_{ij}$  jsou odpovídající relativní frekvence (o nichž opět předpokládáme, že se rovnají pravděpodobnostem příslušných jevů). Nejlepší odhad dostaneme takto: je-li  $x \in V_1$ , splňuje tu vlastnost  $P_j$ , pro níž je  $p_{imax} = p_{1j}$ , ..., je-li  $x \in V_k$ , splňuje tu vlastnost  $p_j$ , pro níž je  $p_{kmax} = p_{kj}$ . Chyba tohoto odhadu je  $1 - \sum_{i=1}^k p_{imax}$ . Zřejmě chyba v situaci 2. nemůže být větší než v situaci 1. Označme

$$\lambda = \frac{1 - p_{.max} - (1 - \sum_{i=1}^k p_{imax})}{1 - p_{.max}} = \frac{\sum_{i=1}^k p_{imax} - p_{.max}}{1 - p_{.max}}$$

$\lambda$  je tedy relativní pokles chyby našeho odhadu při přechodu od situace 1. k situaci 2. Čím větší je hodnota čísla  $\lambda$ , tím lépe nám znalost rozkladu veličiny  $V$  na systém binárních vlastností  $V_1, \dots, V_k$  pomáhá určit rozdělení prvků množiny  $M$  podle vlastností  $P_1, \dots, P_n$ .

Úkolem, který se budeme v této práci snažit vyřešit, bude (zhruba řečeno) nalézt pro dané vlastnosti  $P_1, \dots, P_n$  a danou veličinu  $V$  nejvhodnější binarizaci této veličiny  $V$  (ve smyslu naznačeném výše). V části II. zavedeme přesně pojmy, které jsme zde názorně vyložili a dokážeme ty jejich vlastnosti, které nám v části III. umožní popsat algoritmus pro práci samočinného počítače.<sup>2)</sup>

<sup>2)</sup> Na potřebu podobných úvah upozornil již dříve ing. M. CHYTL.

## II

**Definice.** Modelem nazveme  $(n + 2)$ -tici  $\mathfrak{M} = \langle M, P_1, \dots, P_n, V \rangle$  takovou, že

- (a)  $M$  je konečná množina o  $m$  prvcích
- (b)  $P_1, \dots, P_n$  jsou unární predikáty na  $M$  takové, že každé  $x \in M$  splňuje právě jeden z  $P_1, \dots, P_n$  a že každý z  $P_1, \dots, P_n$  je splněn alespoň jedním  $y \in M$
- (c)  $V$  je reálná funkce na  $M$ .

Nechť  $V$  nabývá na  $M$   $h$  ( $h \geq 3$ ) různých hodnot. Uspořádejme je podle velikosti a označme postupně  $1, \dots, h$ . Tomu odpovídá částečné uspořádání množiny  $M$  takové, že na prvních  $m_1$  prvcích nabývá  $V$  hodnoty 1, na dalších  $m_2$  prvcích hodnoty 2, atd. ( $m_1 + \dots + m_h = m$ ). Jestliže  $V$  nabývá na  $M$  právě  $h$  hodnot  $1, \dots, h$ , řekneme, že  $\mathfrak{M}$  je dobrý model. Předpokládáme v dalším, že model  $\mathfrak{M}$  je dobrý.

**Definice.** Řekneme, že  $D = \{h_1, \dots, h_k\}$  je dělení intervalu  $\langle 1, h \rangle$  stupně  $k$  ( $2 \leq k \leq h$ ), jestliže  $h_1, \dots, h_k$  je vybraná posloupnost z posloupnosti  $1, \dots, h$  taková, že  $1 \leq h_1 < h_2 < \dots < h_k = h$ .

**Definice.** Rozkladem stupně  $k$   $R^k = \{I_1, \dots, I_k\}$  funkce  $V$  na množině  $M$ , indukovaným dělením  $D = \{h_1, \dots, h_k\}$  rozumíme systém částečných intervalů  $I_1 = \langle 1, h_1 \rangle, I_2 = \langle h_1, h_2 \rangle, \dots, I_k = \langle h_{k-1}, h_k \rangle$  intervalu  $\langle 1, h \rangle$ .

**Definice.**  $V_1, \dots, V_k$  buďte unární predikáty, definované na množině  $M$  takto:  $V_i(x) \equiv V(x) \in I_i, i = 1, \dots, k$ .

K danému rozkladu veličiny  $V$  (neboli k danému dělení  $D$ ) sestavme tabulku:

|          | $P_1$         | $P_2$         | ...      | $P_n$         |              |
|----------|---------------|---------------|----------|---------------|--------------|
| $V_1$    | $a_{11}$      | $a_{12}$      | ...      | $a_{1n}$      | $a_{1\cdot}$ |
| $V_2$    | $a_{21}$      | $a_{22}$      | ...      | $a_{2n}$      | $a_{2\cdot}$ |
| $\vdots$ | $\vdots$      | $\vdots$      | $\vdots$ | $\vdots$      | $\vdots$     |
| $V_k$    | $a_{k1}$      | $a_{k2}$      | ...      | $a_{kn}$      | $a_{k\cdot}$ |
|          | $a_{\cdot 1}$ | $a_{\cdot 2}$ | ...      | $a_{\cdot n}$ | $m$          |

kde  $a_{ij}$  je počet prvků množiny  $M$ , které splňují současně  $V_i$  a  $P_j$  ( $i = 1, \dots, k, j = 1, \dots, n$ ), dále  $a_i = \sum_{j=1}^n a_{ij}$ ,  $a_{.j} = \sum_{i=1}^k a_{ij}$ . Položme  $a_{i\max} = \text{Max}_j a_{ij}$ ,  $a_{\max} = \text{Max}_j a_{.j}$ . Zřejmě platí  $m = \sum_{j=1}^n a_{.j} = \sum_{i=1}^k a_i = \sum_{i,j} a_{ij}$ . Označme

$$\lambda_{D,M} = \frac{\sum_{i=1}^k a_{i\max} - a_{\max}}{m - a_{\max}}$$

Označíme-li ještě  $p_{i\max} = a_{i\max}/m$ ,  $p_{\max} = a_{\max}/m$ , je

$$\lambda_{D,M} = \frac{m\left(\sum_{i=1}^k p_{i\max} - p_{\max}\right)}{m(1 - p_{\max})} = \frac{\sum_{i=1}^k p_{i\max} - p_{\max}}{1 - p_{\max}}$$

To znamená, že  $\lambda_{D,M}$  nezávisí na  $m$  (t.j. počtu prvků množiny  $M$ ), neboli na faktických frekvencích v  $M$ , nýbrž pouze na relativních frekvencích. Dále ho budeme tedy značit  $\lambda_D$ .

Snadno se ověří následující vlastnosti čísla  $\lambda_D$  (srov. [1]).

- (1)  $\lambda_D$  je definováno pro každé dělení  $D$ .
- (2)  $0 \leq \lambda_D \leq 1$  pro každé dělení  $D$ .
- (3)  $\lambda_D = 0$  právě když existuje  $n_0$  přirozené tak, že  $1 \leq n_0 \leq n$  a že pro všechna  $i = 1, \dots, k$  je  $p_{i n_0} = p_{i\max}$ .
- (4)  $\lambda_D = 1$  právě když pro každé  $i = 1, \dots, k$  existuje právě jedno  $j$  tak, že  $1 \leq j \leq n$  a že  $p_{ij} \neq 0$ .
- (5)  $\lambda_D$  nezávisí na pořadí  $V_1, \dots, V_k$ .
- (6) jsou-li  $V_1, \dots, V_k$  a  $P_1, \dots, P_n$  statisticky nezávislé, je  $\lambda_D = 0$ .

(3) říká, že  $\lambda_D = 0$ , právě když znalost dělení  $D$  intervalu  $\langle 1, h \rangle$  vůbec nepomáhá k lepšímu rozdělení prvků množiny  $M$  vzhledem k  $P_1, \dots, P_n$ ; naopak (4) znamená, že  $\lambda_D = 1$  právě když dělení  $D$  úplně určuje rozdělení prvků množiny  $M$  vzhledem k  $P_1, \dots, P_n$ . Tvrzení (5) nám umožňuje převést každý model  $\mathfrak{M}$  na dobrý model.

**Věta.** *Necht  $D$  a  $\bar{D}$  jsou dvě dělení intervalu  $\langle 1, h \rangle$  příslušného k veličině  $V$  z  $\mathfrak{M}$ , necht  $\bar{D}$  je zjemnění  $D$ . Potom  $\lambda_{\bar{D}} \geq \lambda_D$ .*

Důkaz. Označme  $k$  resp.  $\bar{k}$  stupeň dělení  $D$  resp.  $\bar{D}$ . Jistě je  $k \leq \bar{k}$ . Je-li  $k = \bar{k}$ , je  $D \equiv \bar{D}$  a tedy  $\lambda_D = \lambda_{\bar{D}}$ . Necht tedy  $k < \bar{k}$ . Zřejmě stačí důkaz převést pro případ  $\bar{k} = k + 1$ . Bez újmy na obecnosti můžeme předpokládat, že  $\bar{D}$  je tvořeno posloupností čísel  $h_0, h_1, \dots, h_k$ , kde  $1 \leq h_0 < h_1 < h_2 < \dots < h_k = h$ ,  $h_0$  je některé z čísel  $1, \dots, h - 1$ . Tabulka, odpovídající dělení  $\bar{D}$ , je

|          |               |               |          |               |                |
|----------|---------------|---------------|----------|---------------|----------------|
|          | $P_1$         | $P_2$         | ...      | $P_n$         |                |
| $V_0^1$  | $a_{01}^1$    | $a_{02}^1$    | ...      | $a_{0n}^1$    | $a_{0\cdot}^1$ |
| $V_0^2$  | $a_{01}^2$    | $a_{02}^2$    | ...      | $a_{0n}^2$    | $a_{0\cdot}^2$ |
| $V_2$    | $a_{21}$      | $a_{22}$      | ...      | $a_{2n}$      | $a_{2\cdot}$   |
| $\vdots$ | $\vdots$      | $\vdots$      | $\ddots$ | $\vdots$      | $\vdots$       |
| $V_k$    | $a_{k1}$      | $a_{k2}$      | ...      | $a_{kn}$      | $a_{k\cdot}$   |
|          | $a_{\cdot 1}$ | $a_{\cdot 2}$ | ...      | $a_{\cdot n}$ | $m$            |

kde  $a_{0j}^1 + a_{0j}^2 = a_{1j}$  pro  $j = 1, \dots, n$ ,  $a_{0\cdot}^1 + a_{0\cdot}^2 = a_{1\cdot}$ . Tedy též  $a_{0\cdot}^1 = \sum_{j=1}^n a_{0j}^1$ ,  
 $a_{0\cdot}^2 = \sum_{j=1}^n a_{0j}^2$ .

Položíme dále

$$a_{0\max}^1 = \text{Max}_j a_{0j}^1, \quad a_{0\max}^2 = \text{Max}_j a_{0j}^2,$$

tedy

$$\lambda_{\overline{D}} = \frac{a_{0\max}^1 + a_{0\max}^2 + \sum_{i=2}^k a_{i\max} - a_{\cdot\max}}{m - a_{\cdot\max}}$$

Jelikož

$$a_{1\max} = \text{Max}_j a_{1j} = \text{Max}_j (a_{0j}^1 + a_{0j}^2) \leq \text{Max}_j a_{0j}^1 + \text{Max}_j a_{0j}^2 = a_{0\max}^1 + a_{0\max}^2$$

dostáváme ihned  $\lambda_{\mathbf{D}} \leq \lambda_{\overline{D}}$ .

**Definice.** Necht'  $\mathbf{D}_1$  a  $\mathbf{D}_2$  jsou dvě dělení intervalu  $\langle 1, h \rangle$  příslušného k veličině  $V$  z  $\mathfrak{M}$ . Řekneme, že  $\mathbf{D}_1$  je vhodnější než  $\mathbf{D}_2$  vzhledem k  $P_1, \dots, P_n$ , jestliže  $\lambda_{\mathbf{D}_1} \geq \lambda_{\mathbf{D}_2}$ .

**Definice.** Maximálním dělením  $\mathbf{D}_{\max}$  intervalu  $\langle 1, h \rangle$  nazveme dělení intervalu  $\langle 1, h \rangle$ , tvořené posloupností čísel  $1, 2, \dots, h$ .

Toto dělení je stupně  $h$  a indukuje rozklad  $R^h$  veličiny  $V$  na intervaly  $I_1 = \langle 1, 1 \rangle$ ,

$I_2 = (1, 2), \dots, I_h = (h - 1, h)$  (tj. na „jednobodové intervaly“). Dělení  $D_{\max}$  je zjemněním každého dělení  $D$  intervalu  $\langle 1, h \rangle$  a podle právě dokázané věty k němu příslušné  $\lambda_{\max} \geq \lambda_D$  pro každé  $D$ .

**Definice.** Necht'  $D = \{h_1, \dots, h_k\}$  je dělení intervalu  $\langle 1, h \rangle$ , které indukuje rozklady  $R^k = \{I_1, \dots, I_k\}$ . Normou dělení  $D$  nazveme číslo  $v(D) = \text{Min}(\bar{I}_i)$ , kde  $\bar{I}_i$  je počet těch čísel  $z$  1, ...,  $h$ , která náleží do intervalu  $I_i$ .

**Definice.** Necht'  $D = \{h_1, \dots, h_k\}$  je dělení intervalu  $\langle 1, h \rangle$ , které indukuje rozklad  $R^k = \{I_1, \dots, I_k\}$ . Normou dělení  $D$  vzhledem k množině  $M$  nazveme číslo  $v_M(D) = \text{Min}_i(\bar{I}_i^M)$ , kde  $\bar{I}_i^M$  je počet prvků množiny  $M$ , pro něž hodnota veličiny  $V$  náleží do intervalu  $I_i$ .

**Definice.** Řekneme, že dělení  $D$  intervalu  $\langle 1, h \rangle$  je prosté, jestliže pro každé dělení  $D$  intervalu  $\langle 1, h \rangle$  takové, že  $D$  je zjemněním  $D$ , platí  $\lambda_D < \lambda_{\bar{D}}$ .

Naši úlohu můžeme nyní formulovat takto. Je-li  $\mathfrak{M} = \langle M, P_1, \dots, P_n, V \rangle$  dobrý model,  $k, d$  daná přirozená čísla ( $2 \leq k \leq h, 1 \leq d \leq h/2$ ), hledáme prosté dělení intervalu  $\langle 1, h \rangle$  stupně nejvýše  $k$ , jehož norma  $v$  je větší nebo rovna  $d$ , nejvýhodnější vzhledem k  $P_1, \dots, P_n$ . Je-li takovýchto dělení více, uspořádáme je podle normy  $v_M$  vzhledem k množině  $M$ .

### III

V této části popíšeme algoritmus řešení dané úlohy se zřetelem k jeho realizaci na samočinném počítači.

(A) Přípravná část.

1) Zjistit počet  $h$  hodnot, které nabývá veličina  $V$  na množině  $M$ . Je-li  $h = 1$ , práce končí. Je-li  $h = 2$ , práce končí ( $V$  je již binarizována).

2) Provéřit, zda některý z predikátů  $P_1, \dots, P_n$  není na  $M$  pravdivý nebo lživý. Pokud ano, vyřadit.

3) Ověřit, zda  $d \leq h/2$ . Je-li  $d > h/2$ , práce končí.

4) Uspořádat funkční hodnoty veličiny  $V$  podle velikosti, přiřadit jim čísla 1, ...,  $h$  a sestavit tabulku viz str. 123 (t.j. tabulku příslušnou k dělení  $D_{\max}$ ), kde  $a_{ij}$  ( $i = 1, \dots, h, j = 1, \dots, n$ ) je počet objektů  $M$ , pro něž je veličina  $V$  nabývá hodnoty  $i$  a které splňují predikát  $P_j$ ,  $a_{i..}$ ,  $a_{.j}$  příslušné součty,  $a_{i\max}$ ,  $a_{\max}$  příslušná maxima.

5) Spočítat

$$\lambda_{\max} = \frac{\sum_{i=1}^h a_{i\max} - a_{\max}}{m - a_{\max}}$$

|          |               |          |               |              |
|----------|---------------|----------|---------------|--------------|
|          | $P_1$         | ...      | $P_n$         |              |
| 1        | $a_{11}$      | ...      | $a_{1n}$      | $a_{1\cdot}$ |
| 2        | $a_{21}$      | ...      | $a_{2n}$      | $a_{2\cdot}$ |
| $\vdots$ | $\vdots$      | $\ddots$ | $\vdots$      | $\vdots$     |
| $h$      | $a_{h1}$      | ...      | $a_{hn}$      | $a_{h\cdot}$ |
|          | $a_{\cdot 1}$ | ...      | $a_{\cdot n}$ | $m$          |

(B) Vlastní práce.

1)  $\bar{k} = 2$ . Generovat (v lexikografickém uspořádání) všechna dělení intervalu  $\langle 1, h \rangle$  stupně 2 normy alespoň  $d$ . Pro každé dělení  $D$  zjistit  $\lambda_D$ . Je-li větší nebo rovno dosud maximálnímu, zjistit ještě  $v(D)$ ,  $v_M(D)$ , jinak generovat další. Je-li  $\lambda_D = \lambda_{\max}$ , vytisknout (spolu se zjištěnými údaji) a práce končí. Je-li  $\lambda_D < \lambda_{\max}$  přejít k dalšímu dělení. Pamatovat průběžně dělení s největším  $\lambda$ . Po skončení práce pro  $\bar{k} = 2$  toto nejlepší dělení spolu se zjištěnými údaji vytisknout, je-li jich více, uspořádat nejprve podle normy  $v$ , potom podle normy  $v_M$ . První z tištěných dělení ponechat v paměti.

2)  $\bar{k} = \bar{k}_0 + 1$ . Ověřit  $\bar{k} \leq k$ . Je-li  $\bar{k} > k$ , práce končí. Ověřit  $d \leq h/(\bar{k}_0 + 1)$ . Neplatí-li, práce končí. Generovat postupně (v lexikografickém uspořádání) dělení intervalu  $\langle 1, h \rangle$  stupně  $\bar{k}$  normy alespoň  $d$ . Pro každé zjistit  $\lambda_D$ , je-li větší nebo rovno dosud maximálnímu a je-li prosté, zjistit ještě  $v(D)$ ,  $v_M(D)$ . Pamatovat průběžně prosté dělení s největším  $\lambda$  (narazí-li se na dělení  $D$  takové, že  $\lambda_D = \lambda_{\max}$  tisknout a konec) a po skončení práce pro  $\bar{k} = \bar{k}_0 + 1$  tisknout jako v 1).

3)  $\bar{k} = k$ . Ověřit  $d \leq h/k$ . Neplatí-li, práce končí. Jinak generovat, pamatovat a tisknout jako ve 2). Konec.

Literatura

- [1] L. A. Goodman, W. H. Kruskal: Measures of Association for Cross Classifications, JASA, 49 (1954).
- [2] P. Hájek, I. Havel, M. Chytil: GUHA-metoda systematického vyhledávání hypotéz, Kybernetika 2, 1966, 1, 31–47.
- [3] P. Hájek, I. Havel, M. Chytil: GUHA-metoda systematického vyhledávání hypotéz II, Kybernetika 3, 1967, 5, 430–437.
- [4] P. Hájek: Problém obecného pojetí metody GUHA, Kybernetika — v tisku.



## Summary

### AUTOMATIC BINARIZATION OF QUANTITIES

ZDENĚK RENC

When processing automatically biological (and other) experimental data the need often occurs to replace a certain measured quantity (e.g. pressure, electric potential, age etc.) by one or more properties (e.g. pressure heavy or light, temperature less than 30 °C, between 30 °C–60 °C, over 60 °C, etc.). As a property may be understood as a binary quantity, the problem of the replacement of a quantity by properties can be named the binarization problem. (This problem occurs, e.g., when applying the GUHA method of automatic hypotheses determination, see [2], [3], [4].) We ask which of possible binarizations is the best one (e.g. whether the binarization of temperature given above is better than the following one: less than 40 °C, between 40 °C–70 °C, over 70 °C). On the basis of the paper [1], it is possible to compare numerically which of two binarizations is more suitable with respect to one or more properties. In the present paper a method (called ABQ) is described for the automatic determination of the best binarization of a quantity with respect to given properties on the basis of given experimental material. When the research worker determines the maximal number of intervals that the quantity is to be divided into and the properties with respect to which the suitability is to be measured (and, of course, supplies the experimental data), the computer gives all the optimal binarizations. The computer counts also supplementary information concerning the binarizations obtained.

*Adresa autora:* RNDr. Zdeněk Renc, Matematicko-fyzikální fakulta KU, Sokolovská 83, Praha 8.