

## An Alternative Approach to Missing Information in the GUHA Method

TOMÁŠ HAVRÁNEK

In the present paper some new approaches to the treatment of missing data items in the GUHA procedures are presented. The first one is the dual approach to the usual conservative approach used in present GUHA procedures. Both these approaches are independent of the pollution process destroying the original data matrix with no missing data items. The second new approach is principally based on erasing objects with missing data items in just processed parts of data.

In the GUHA method, since [2], [1] a "pesimistic" way for the treatment of missing information has been applied, i.e. generalized quantifiers and/or underlying two-valued decision functions are extended to data containing missing information in the following conservative way:

Data are finite structures of the form  $M = \langle M, f_1, \dots, f_n \rangle$ , where  $M \neq \emptyset$ ,  $M \subseteq N$  and each  $f_i : M \rightarrow V$ ;  $n$  is called the type of data (elements of  $M$  are understood as names of observed objects and  $f_i(o)$  for an object  $o \in M$  as observed value of some investigated quantity). Data containing missing information are such that  $f_i : M \rightarrow V \cup \{ \times \}$ , where the value  $\times$  means "unspecified", "unobserved" etc. Now for given  $M$ ,  $V$  denote  $\mathcal{M}_{V,M}^\times$  the class of all structures (data) containing (possibly) missing information. By  $\mathcal{M}_{V,M}$  we denote structures with no missing values. Denote  $\mathcal{D} = \mathcal{P}_{fin}(N) - \{ \emptyset \}$ . Then put  $\mathcal{M}_V = \bigcup_{M \in \mathcal{D}} \mathcal{M}_{V,M}$ . Similarly for  $\mathcal{M}_V^\times$ . If  $V = \{0, 1\}$ , we write only  $\mathcal{M}_M$  etc. Let now  $M \in \mathcal{M}_{V,M}^\times - \mathcal{M}_{V,M}$ ; data  $N \in \mathcal{M}_{V,M}$  are called a *completion* of  $M$  if all values  $\times$  in  $M$  are in  $N$  replaced by some values from  $V$  and the rest is left unchanged.

The underlying philosophy is, that data with no missing information have been damaged by a "pollution" process and we have to make some statistical decisions without knowledge of the original data — "true" completion. No assumptions concerning the pollution process are done.

A decision function  $f: M_V \rightarrow \{0, 1\}$  (i.e. defined on data with complete information can be *conservatively* extended to incomplete data (i.e. onto  $\mathcal{M}_V^*$ ) as follows:

Put

$$f(M) = \begin{cases} 1 & \text{if for all completions } N \text{ of } M, f(N) = 1, \\ \times & \text{otherwise,} \\ 0 & \text{if for all completions } N \text{ of } M, f(N) = 0. \end{cases}$$

This is in accordance with the meaning of  $\times$  as "unknown". Usually an ordering of values  $0 < \times < 1$  is supposed. Then conservative extension can be obtained by using a "threshold" value 1. In [1] such an extension is called *secured* extension (some decision is accepted iff  $f(M) \geq 1$ ). Cf. also [8].

Logical connectives  $\vee$ ,  $\&$  and  $\neg$  can be extended using the "securedness" principle: Put the value 1 (or 0) if you are sure that the value must be 1 (or 0). Then we obtain for  $\vee$ ,  $\&$ ,  $\neg$  the following tables:

$\neg$	$0 \times 1$	$\&$	$0 \times 1$	$\vee$	$0 \times 1$
	$1 \times 0$	0	0 0 0	0	0 $\times$ 1
		$\times$	0 $\times$ $\times$	$\times$	$\times$ $\times$ 1
		1	0 $\times$ 1	1	1 1 1

The above conservative way can be criticised at least from the following two reasons (think decision 1 as acceptance of an alternative hypothesis):

(i) It is in contradiction with the philosophy of the GUHA methods (c. f. for example the introductory paper [4] — searching of all "interesting" in data and interestingness is equivalent with acceptance); it overestimates local reliability of results on the account of the global power of the method.

(ii) In practice, one can see that even relatively small "pollution" of data can lead to heavy consequences in omitting many results, which can be obtained e.g. if we exclude all objects with some missing data item from the processed data.

Hence some alternative ways for the treatment of missing information are to be suggested. The critical point of such suggestions is the possibility to implement new alternative decision functions into present computer GUHA procedures without destroying their architecture. The further suggested extensions will be discussed from this point of view, hence the present paper is rather technical. A preliminary version of this paper is obtained in the report [7].

Some important questions concerning so called helpful quantifiers connected with the presentend approach is postponed into the forthcoming paper [5].

**1.0.** The “optimistic” extension of a decision function is defined as follows:  $f(\mathbf{M}) = 1$  if there is a completion  $\mathbf{N}$  of  $\mathbf{M}$  such that  $f(\mathbf{N}) = 1$ , otherwise  $f(\mathbf{M}) = 0$ .

**1.1.** Or, in three valued logic we use threshold value  $\times$ : set of designed values (interesting values) is defined as  $\{1, \times\}$ . It means that we seek for “potentially” true sentences. For the sake of simplicity, we follow only the way 1.0. without any loss of generality.

**1.2.** In the rest of paper we use freely technical means of observational functor calculi as described in [3] (or partially in [1]). We are concerned in generalized quantifiers defined by some decision functions. If  $q$  is such a quantifier, its defining function is denoted  $\text{Asf}_q$  and it is called the associated function of  $q$ . We shall use open formulas with one variable  $x$  (and hence we do not write the variable explicitly).

For example if  $q$  joins two open formulas (names of composite qualities)  $\varphi, \psi$ , the value of  $q(\varphi, \psi)$  on a structure of an appropriate type is defined by  $\|q(\varphi, \psi)\|_{\mathbf{M}} = \text{Asf}_q(\mathbf{M}_{\varphi, \psi})$ , where  $\mathbf{M}_{\varphi, \psi} = \langle \mathbf{M}, \|\varphi\|_{\mathbf{M}}, \|\psi\|_{\mathbf{M}} \rangle$  (note, that  $\|\varphi\|_{\mathbf{M}}$  is a function on  $\mathbf{M}$ , the domain of  $\mathbf{M}$ ). Moreover, we shall assume that  $V = \{0, 1\}$ .

**1.3.** For associational quantifiers (cf. 3.2.2 of [3]) we have to check its “associativity” for optimistic extension. Consider quantifiers of the type 2 (i.e. joining two open formulas). Their associated functions are defined on structures of the form  $\mathbf{M} = \langle \mathbf{M}, f_1, f_2 \rangle$ . For each pair  $\mathbf{u} = \langle u, v \rangle \in \{0, \times, 1\}^2$ , define  $\text{fr}(\mathbf{M}, \mathbf{u}) = \text{card} \{o \in \mathbf{M}; \langle f_1(o), f_2(o) \rangle = \mathbf{u}\}$ . Then we denote  $a(\mathbf{M}) = \text{fr}(\mathbf{M}, \langle 1, 1 \rangle)$ ,  $b(\mathbf{M}) = \text{fr}(\mathbf{M}, \langle 1, 0 \rangle)$ ,  $c(\mathbf{M}) = \text{fr}(\mathbf{M}, \langle 0, 1 \rangle)$ ,  $d(\mathbf{M}) = \text{fr}(\mathbf{M}, \langle 0, 0 \rangle)$ .

A quantifier is associational if  $\text{Asf}_q(\mathbf{M}) = 1$  and  $\mathbf{N} \geq^a \mathbf{M}$  implies  $\text{Asf}_q(\mathbf{N}) = 1$ , where  $\geq^a$  is the  $a$  - better relation, i.e.  $\mathbf{M} \leq^a \mathbf{N}$  if  $a(\mathbf{M}) \leq a(\mathbf{N})$ ,  $b(\mathbf{M}) \geq b(\mathbf{N})$ ,  $c(\mathbf{M}) \geq c(\mathbf{N})$ ,  $d(\mathbf{M}) \leq d(\mathbf{N})$ . Similarly for  $\geq^i$ , the  $i$  - better ordering:  $\mathbf{M} \leq^i \mathbf{N}$  if  $a(\mathbf{M}) \leq a(\mathbf{N})$ ,  $b(\mathbf{M}) \geq b(\mathbf{N})$ , we can define *implicational quantifiers*.

For incomplete data, define  $i(\mathbf{M}) = \text{fr}(\mathbf{M}, \langle 1, \times \rangle)$ ,  $o(\mathbf{M}) = \text{fr}(\mathbf{M}, \langle \times, 1 \rangle)$ ,  $n(\mathbf{M}) = \text{fr}(\mathbf{M}, \langle \times, \times \rangle)$ ,  $p(\mathbf{M}) = \text{fr}(\mathbf{M}, \langle \times, 0 \rangle)$ ,  $q(\mathbf{M}) = \text{fr}(\mathbf{M}, \langle 0, \times \rangle)$ .

**1.4.** Let  $\mathbf{M} \in \mathbf{M}^\times$ . Let  $\mathbf{M}^b$  is any completion of  $\mathbf{M}$  such that  $a(\mathbf{M}^b) = a(\mathbf{M}) + i(\mathbf{M}) + o(\mathbf{M}) + n(\mathbf{M})$  and  $b(\mathbf{M}^b) = b(\mathbf{M})$ . Then  $\mathbf{M}^b$  is called the *best implicational completion* of  $\mathbf{M}$ .

**Theorem.** For each implicational quantifier  $\rightarrow$  and its optimistic extension  $\rightarrow^*$  we have  $\text{Asf}_{\rightarrow^*}(\mathbf{M}) = 1$  iff  $\text{Asf}_{\rightarrow}(\mathbf{M}^b) = 1$  for the best implicational completion of  $\mathbf{M}$ .

**Proof.** Clearly, for each other completion  $\mathbf{N}$  of  $\mathbf{M}$ ,  $a(\mathbf{N}) \leq a(\mathbf{M}^b)$  and  $b(\mathbf{N}) \geq b(\mathbf{M}^b)$ .

**1.5.** Our first aim is now to inspect observational inference rules (immediate

148 consequence rules) used in the implicational GUHA procedure as concerns their soundness for the optimistic extension of implicational quantifiers.

**1.6. Theorem.** The despecifying — dereducing inference rule (3.2.20 [3]) is sound for optimistic extension of any implicational quantifier.

**Proof.** Let  $\rightarrow^*$  be an optimistic extension of an implicational quantifier. The inference rule in question is defined as

$$\text{SpRd} = \left\{ \frac{\varphi \& \psi \rightarrow^* \delta}{\varphi \rightarrow^* \delta \vee \neg \psi \vee \chi} ; \begin{array}{l} \varphi, \psi \text{ elementary conjunctions,} \\ \delta, \chi \text{ elementary disjunctions and} \\ \varphi, \psi, \delta, \chi \text{ mutually disjoint.} \end{array} \right\}$$

Consider two steps (using transitivity of the rule SpRd):

(i) Let  $a, b, \dots$  concern the structure  $M_{\varphi, \delta}$ ,  $a', b', \dots$  the structure  $M_{\varphi, \delta \vee \chi}$  (both derived from a structure  $M$ ). Inspect the following Table 1:

Table 1.

frequencies	cards	frequencies
	$\varphi \quad \delta \quad \chi$	
$a$	$\left\{ \begin{array}{l} 1 \quad 1 \quad \text{everything} \\ 1 \quad 0 \quad 1 \end{array} \right\}$	$a'$
$i$	$\left\{ \begin{array}{l} 1 \quad \times \quad 1 \\ 1 \quad \times \quad 0 \\ 1 \quad \times \quad \times \end{array} \right\}$	$i'$
$o$	$\left\{ \begin{array}{l} 1 \quad 0 \quad \times \\ \times \quad 1 \quad 0 \\ \times \quad 1 \quad \times \\ \times \quad 1 \quad 1 \\ \times \quad 0 \quad 1 \\ \times \quad \times \quad 1 \end{array} \right\}$	$o'$
$n$	$\left\{ \begin{array}{l} \times \quad \times \quad 0 \\ \times \quad \times \quad \times \\ \times \quad \times \quad \times \\ \times \quad 0 \quad \times \end{array} \right\}$	$n'$

Then  $a(M_{\varphi, \delta \vee \chi}^b) = a' + i' + o' + n' \geq a + i + o + n = a(M_{\varphi, \delta}^b)$  and  $b(M_{\varphi, \delta \vee \chi}^b) \leq b(M_{\varphi, \delta}^b)$ .

(ii) Similarly for  $\varphi \& \psi \rightarrow^* \delta$  and  $\varphi \rightarrow^* \delta \vee \neg \psi$ : let  $a, b, \dots$  concern  $M_{\varphi \& \psi, \delta}$  and  $a', b', \dots$  concern  $M_{\varphi, \delta \vee \neg \psi}$ . Inspect the Table 2:

Table 2.

frequencies	cards	frequencies
	$\varphi \quad \psi \quad \delta \quad \neg\psi$	
$a$	$\left\{ \begin{array}{l} 1 \quad 1 \quad 1 \quad 0 \\ 1 \quad \text{ev.} \quad 1 \quad \text{ev.} \\ 1 \quad 0 \quad \times \quad 1 \\ 1 \quad 0 \quad 0 \quad 1 \end{array} \right\}$	$a'$
$i$	$\left\{ \begin{array}{l} 1 \quad 1 \quad \times \quad 0 \\ 1 \quad \times \quad \times \quad \times \\ 1 \quad \times \quad 0 \quad \times \end{array} \right\}$	$i'$
$o$	$\left\{ \begin{array}{l} \times \quad 1 \quad 1 \quad 0 \\ 1 \quad \times \quad 1 \quad \times \\ \times \quad \times \quad 1 \quad \times \rightarrow a' \\ \times \quad 1 \quad 1 \quad 0 \\ \times \quad 0 \quad 1 \quad 1 \\ \times \quad 0 \quad 0 \quad 1 \\ \times \quad \times \quad 0 \quad \times \\ \times \quad \times \quad \times \quad \times \\ \times \quad 1 \quad \times \quad 0 \\ 1 \quad \times \quad \times \quad \times \rightarrow i' \\ \times \quad \text{ev.} \quad \times \quad 0 \\ \times \quad \text{ev.} \quad 0 \quad \times \\ \times \quad \text{ev.} \quad \times \quad \times \end{array} \right\}$	$o'$
$n$		$n'$

Then  $a(M_{\varphi, \delta \vee \neg\psi}^b) = a' + i' + o' + n' \geq a + i + o + n = a(M_{\varphi \& \psi, \delta})$  and  $b(M_{\varphi, \delta \vee \neg\psi}^b) = b(M_{\varphi \& \psi, \delta})$ .

1.7. For a  $\{0, \times, 1\}$  - structure  $M$  of the type 2, define the *best symmetric completion* as any completion  $M^b$  such that:  $a(M^b) = a(M) + i(M) + o(M) + n_1$ ,  $d(M^b) = d(M) + j(M) + p(M) + n_2$ ,  $c(M^b) = c(M)$ ,  $b(M^b) = b(M)$  with  $n_1, n_2 \in \mathbb{N}$  such that  $a(M^b)d(M^b) = \max_{g+h=n(M), g, h \in \mathbb{N}} (a(M) + i(M) + o(M) + g)(d(M) + j(M) + p(M) + h)$ .

1.8. **Theorem.** Let  $\sim$  be one of the quantifiers  $\sim, \sim_a, \sim_a^2$  (i.e. in words, the simple associational, Fisher and chi-square quantifier). Then for each  $M$  and each best symmetric completion  $M^b$  of  $M$ , we have  $\text{Asf}_{\sim}(M) = 1$  iff  $\text{Asf}_{\sim}(M^b) = 1$ .

1.9. We first prove the following lemma: For each associational quantifier  $\sim$  we have the following: for each  $M$ ,  $\text{Asf}_{\sim}(M) = \text{Asf}_{\sim}(M^*)$ , where  $M^* = M(\langle 1, \times \rangle : \langle 1, 1 \rangle) (\langle \times, 0 \rangle : \langle 0, 0 \rangle) (\langle \times, 1 \rangle : \langle 1, 1 \rangle) (\langle 0, \times \rangle : \langle 0, 0 \rangle)$ , i.e.  $M^*$  is obtained from  $M$  replacing each card  $\langle 1, \times \rangle$  by  $\langle 1, 1 \rangle$  etc.

Proof. (i) for each completion of  $M^*$  we can construct an  $a$ -better completion of  $M$ ; clearly each completion of  $M^*$  is a completion of  $M$ , (ii) Let  $N$  is a completion of  $M$ , then changing completed as cards suggested above we obtain an completion of  $M^*$   $a$ -better than  $N$ .

**1.10. Proof of Theorem 1.1.** For quantifiers  $\sim$  and  $\sim_\alpha^2$  the proof is now trivial, due their symmetry w. r. t.  $|a - d|$  and  $|b - c|$ . We have to treat the case of  $\sim_\alpha$ ; namely prove the monotonicity of the statistic

$$D(a, r, k, m) = \sum_{i=a}^{\min(r, k)} \sigma(i, r, k, m),$$

where  $\sigma(i, r, k, m) = r! s! k! l! / m! i! b! c! d!$  with  $b = k - i$ ,  $c = r - i$ ,  $s = m - r$ ,  $l = m - k$ ,  $d = m + i - r - k$ . Consider now tables completed due to the Lemma 1.9.:

$$\begin{array}{cc|c} a + n' & b & k + n' \\ c & d + n - n' & l + n - n' \\ \hline r + n' & s + n - n' & m \end{array},$$

where  $n' \in [0, n]$ . We prove first an auxiliary inequality for the contingency table

$$\begin{array}{cc|c} a & b & k \\ c & d & l \\ \hline r & s & m \end{array}$$

with  $a \geq d$  and  $ad > bc$ . Consider inequality

$$(1) \quad sl/d \leq rk/a.$$

We can use the following equivalent forms of (1):

$$(2) \quad (a + d)(d + c)/d \leq (a + c)(a + b)/a;$$

$$(3) \quad d + bc/d \leq a + bc/a;$$

$$(4) \quad d(ad) + a(bc) \leq a(ad) + d(bc)$$

and, due  $d - a \leq 0$ , from  $(d - a)ad \leq bc(d - a)$  we have  $ad \geq bc$ , which is true. Due to the symmetry of the statistic we can suppose that in any case  $a + n' \geq d + n - n'$  and prove the monotonicity of  $H(n') = D(a + n', k + n', r + n', m)$  for  $n'$  increasing. I.e. we have to prove  $H(n') \leq H(n' + 1)$  for each  $n'$  such that  $a + n' \geq d + n - n'$  and  $n' + 1 \leq n$ . Both sums have the same number of members; we can compare  $\sigma(i + n', r + n', k + n', m)$  and  $\sigma(i + n' + 1, r + n' + 1,$

$k + n' + 1, m)$  for  $i = a, \dots, \min(k, r)$ . Denote  $i' = i + n'$ ,  $r' = r + n'$  etc. The desired inequality is then:  $\sigma(i', r', k', m) \leq \sigma(i' + 1, r' + 1, k' + 1, m)$  for  $i' = a' = a + n', \dots, \min(k', r')$ . We can reduce

$$r! (m - r')! (m - k')! k'! / m! i'! (k' - i')! (r' - i')! (m + i' - r' - k')! \leq (r' + 1)! (m - r' - 1)! (m - k' - 1)! (k' + 1)! / m! (i' + 1)! (k' + 1 - i' - 1)! (r' + 1 - i' - 1)! (m + i' + 1 - r' - 1 - k' - 1)! \text{ to}$$

$$(5) \quad (m - r')(m - k')(m + i' - r' - k') \leq (r' + 1)(k' + 1)(i' + 1)$$

for  $i' = a', \dots, \min(r', k')$ . First, note that

$r'k'/i' \leq (r' + 1)(k' + 1)/(i' + 1)$ ; the last inequality is equivalent to  $r'k'i' + r'i' + k'i' + i' \geq r'k'i' + r'k'$  which is equivalent to  $i'(r' + k') + i \geq r'k'$  and then  $a'((a' + b') + (a' + c')) + a' \geq (a' + b')(a' + c')$ . It is equivalent to simple inequality  $a'^2 + a' \geq bc'$ ; using  $a' \geq d'$  we obtain  $a'd' \geq b'c'$  which is true, and from which the previous inequality follows. Now we can consider instead of (5) the inequality

$$(6) \quad (m - r')(m - k')(m + i' - r' - k') \leq r'k'/i'.$$

For  $m - r' - k' \geq 0$  the ratio  $(m + i' - r' - k')/i'$  is decreasing function of  $i'$ , hence we can test (6) for  $i' = k'$  or  $i' = r'$ . Then we obtain  $(m - k') \leq r'$  (or  $(m - r') \leq k'$ ) which is true due the fact  $a' \geq d'$ . For  $m - r' - k' < 0$  we must test the inequality for  $i' = a'$  and use the inequality (1).

## 2. THE UNBIASED APPROACH

**2.1.** We can extend an quantifier  $q$  into a new quantifier  $q^*$  by the following way:  $\text{Asf}_{q^*}(M) = \text{Asf}_q(M^e)$ , where  $M^e$  is obtained from  $M$  by errasing all rows containing at least one value  $\times$ .

**2.2.** Particularly, in GUHA procedures with quantifiers of the type 2, we evaluate a sentence  $q(\varphi, \psi)$  on  $M_{\varphi, \psi}^e$  obtained from  $M_{\varphi, \psi}$  by erasing all cards  $\langle \times, \times \rangle, \langle 1, \times \rangle, \langle 0, \times \rangle, \langle \times, 1 \rangle, \langle \times, 0 \rangle$ . It means that objects with missing information are not omitted from  $M$  before computing, but only in evaluating particular sentence  $q^*(\varphi, \psi)$ .

**2.3. Theorem.** Let  $\rightarrow$  be an implicational quantifier. Then the inference rule SpRd is sound for its unbiased extension  $\rightarrow^*$ .

**Proof.** Consider the same steps as in 1.6. First: Clearly, by Table 1,  $a' \geq a$ . Similarly  $b' \leq b$ :

$$\begin{array}{c}
 \varphi \quad \delta \quad \chi \\
 b \left\{ \begin{array}{ccc} 1 & 0 & 0 \\ 1 & 0 & \times \\ 1 & 0 & 1 \end{array} \right\} \quad b'
 \end{array}$$

Second: by Table 2, we have  $a' \geq a$  again. Similarly  $b' = b$ :

$$\begin{array}{c}
 \varphi \quad \psi \quad \delta \quad \neg\psi \\
 b \left\{ \begin{array}{cccc} 1 & 1 & 0 & 0 \end{array} \right\} \quad b'
 \end{array}$$

**2.4.** Note that this way of extension of implicational and associational quantifiers is enabled by the new definition of associativity and implicativity used in [3] in the opposition to the old one used in [1], [2], that demanded for  $M \leq^a N$  to have  $\text{card}(M) = \text{card}(N)$ .

### 3. DISCUSSION

**3.1.** We shall summarize here some (more or less evident) properties of the three mentioned ways of extending quantifiers for the three valued  $\{0, \times, 1\}$  – calculus.

Consider quantifiers defined on  $\{0, 1\}$  structures of the type 2. Denote, for such an quantifier  $q$ , by  $q_s$ ,  $q_e$  and  $q_o$  its secured, unbiased and optimistic extension.

**3.2. Fact.** Let  $q$  be an associational quantifier. Then  $\text{Asf}_{q_s} \leq \text{Asf}_{q_e} \leq \text{Asf}_{q_o}$ .

*Proof.* Let  $M \in \mathcal{M}^\times$  and  $a, b, \dots$  be corresponding frequencies. Denote by  $a_c, b_c, \dots$  the frequencies obtained by completion. For the secured case consider frequencies of completions with  $a_c = a$ ,  $d_c = d$ ,  $b_c \geq b$ ,  $c_c \geq c$ . If  $q_s$  is secured and  $\text{Asf}_{q_s}(M) = 1$  then for all the above completions  $\text{Asf}_q(M^c) = 1$ . By associativity  $\text{Asf}_{q_e}(M) = 1$ . The case of  $q_o$  is analogous.

**3.3. Note.** If  $q$  is associational and saturable (or executive), then for each  $M \in \mathcal{M}^\times - \mathcal{M}$ , at least one of the inequalities in 3.2 is strict (for models with domains greater then a threshold value  $m_0$ ).

**3.4.** Suppose now that data (models, structures) are obtained by some random experiments, i.e. we have a probability space  $\langle \Sigma, \mathcal{R}, P \rangle$  and  $Q : \mathcal{P}_{\text{fin}}(N) \times \Sigma \rightarrow \mathcal{M}$  ( $Q$  partialized on  $M$ ,  $M \in \mathcal{P}_{\text{fin}}(N)$ , is measurable). Moreover, assume that for each  $M \in \mathcal{P}_{\text{fin}}(N)$  and each  $M \in \mathcal{M}_M$  we have  $P_M(\{\omega \in \Sigma; Q(\omega, M) = M\}) > 0$  (e.g. data are obtained by independent multinomial sampling with no zero probabilities).

**3.5.** Denote, for each  $M$ , by  $\mathcal{C} \subseteq \mathcal{M}_M \times \mathcal{M}_M^\times$  the completion relation, i.e.  $\langle M, N \rangle \in \mathcal{C}$  iff  $M$  is a completion of  $N$ .



Consider now a (deterministic) pollution mechanism **pol** as a mapping from  $\mathcal{M}_M$  into  $\mathcal{M}_M^*$ ,  $\mathbf{pol} \subseteq \mathcal{C}$ .

For each such a mapping **pol**, we can define for any  $\mathcal{N} \subseteq \mathcal{M}_M^*$ ,  $P(\mathcal{N})$  as  $P(\mathcal{N}) = P_M(\mathbf{pol}^{-1}(\mathcal{N}))$ . Write now  $q(\mathbf{M})$  instead of  $\text{Asf}_q(\mathbf{M})$ . Consider now sets

$$\mathcal{N}_s = \{N \in \mathcal{M}_M^*; q_s(N) = 1\} \quad \text{and} \quad \mathcal{N}_o = \{N \in \mathcal{M}_M^*; q_o(N) = 1\}$$

and denote  $p_s = P(\mathcal{N}_s)$  and  $p_o = P(\mathcal{N}_o)$ . Pedantically, we should write  $p_s(\mathbf{pol}, M)$  etc. Put  $p = P_M\{\mathbf{M} \in \mathcal{M}_M; q(\mathbf{M}) = 1\}$ .

**3.6. Theorem.** Under assumptions of 3.3–3.5, we have for  $\text{card}(M) \leq m_o$ ,  $p_s = p = p_o = 0$  and, for  $\text{card}(M) > m_o$ ,  $p_s \leq p \leq p_o$ .

*Proof.* The first part is clear; for  $\text{card}(M) \leq m_o$  there is no completion with  $q(\mathbf{M}) = 1$ . For the second part, we have to prove that for each **pol** the following is true:

$$\mathbf{pol}^{-1}(\mathcal{N}_s) \subseteq \{\mathbf{M} \in \mathcal{M}_M; q(\mathbf{M}) = 1\} \subseteq \mathbf{pol}^{-1}(\mathcal{N}_o).$$

Clearly: if  $N \in \mathcal{N}_s$  then for each  $\mathbf{M}$  such that  $\langle \mathbf{M}, N \rangle \in \mathcal{C}$ ,  $q(\mathbf{M}) = 1$ , hence for each  $\mathbf{M} \in \mathbf{pol}^{-1}(\mathcal{N}_s)$ ,  $q(\mathbf{M}) = 1$ . The second inclusion can be proved similarly.

**3.7.** There are pollutions mappings, for which  $p_s < p < p_o$ ; moreover, there is an pollution mapping giving  $p_o = 1$  and  $p_s = 0$ . For example, a necessary condition for a mapping to satisfy  $p_s < p < p_o$  is: for each  $M$ ,  $\text{card}(M) > m_o$ , there are  $\mathbf{M}_1, \mathbf{M}_2 \in \mathcal{M}_M$  such that **pol**( $\mathbf{M}_1$ ) has a completion  $\mathbf{M}'_1$  with  $q(\mathbf{M}'_1) = 0$  and  $q(\mathbf{M}_1) = 1$  and **pol**( $\mathbf{M}_2$ ) has a completion  $\mathbf{M}'_2$  such that  $q(\mathbf{M}'_2) = 1$  and  $q(\mathbf{M}_2) = 0$ .

For the constant mapping  $\mathbf{pol}_x : \mathcal{M}_M \rightarrow X$ , where  $X$  is a model containing only  $x$ -values, clearly  $p_s = 0$  and  $p_o = 1$ .

**3.8.** We can consider a probabilistic model of the pollution process: Consider for each  $M$  a pair of (generally dependent) random variables  $\langle Q, \mathbf{pol}^* \rangle : \Sigma \rightarrow \mathcal{C}$ . Then for each  $\mathbf{M}$  we obtain a conditional distribution on

$$\mathbf{pol}^*(\mathbf{M}) = \{N; \exists \omega \in \Sigma, \langle Q, \mathbf{pol}^* \rangle(\omega) = \langle \mathbf{M}, N \rangle\}.$$

It can be extended onto  $\mathcal{M}_M^*$  putting  $P(N|\mathbf{M}) = 0$  for each  $\langle \mathbf{M}, N \rangle \notin \mathcal{C}$ . Assume that for each  $\langle \mathbf{M}, N \rangle \in \mathcal{C}$ ,  $P(\langle \mathbf{M}, N \rangle) > 0$ .

Now we can clearly speak about (for a given domain  $M$ )  $p_o = P(\mathcal{N}_o)$  and  $p_s = P(\mathcal{N}_s)$ .

**3.9. Theorem.** Under assumptions of 3.3, 3.4 and 3.8, we have  $p_o < p < p_s$  (for  $\text{card}(M) > m_o$ , else  $p_o = p = p_s = 0$ ).

*Proof.* Consider the case of  $\text{card}(M) > m_o$ . Then  $p > 0$ . For each  $\mathcal{N} \subseteq \mathcal{M}_M^*$  denote  $\mathcal{N}^M = \{N \in \mathcal{N}; \langle \mathbf{M}, N \rangle \in \mathcal{C}\}$ . Consider  $P(\mathcal{N}_o) = \sum_{\mathbf{M} \in \mathcal{M}_M} P(\mathcal{N}_o|\mathbf{M}) P(\mathbf{M}) =$

154  $= \sum_{M \in \mathcal{M}_M} P(\mathcal{N}_s^M/M) P(M)$ . Note that  $p = P_M(\{M \in \mathcal{M}_M; q(M) = 1\}) = \sum_{M \in \mathcal{M}_M} q(M) P(M)$  and that if  $q(M) = 1$  then  $P(\mathcal{N}_s^M/M) = 1$ ; clearly  $p < p_o$ . Similarly for  $p_s < p$  (if  $q(M) = 1$  then  $P(\mathcal{N}_s^M/M) = 0$ ). The strict inequality follows from the fact that for each  $M$ ,  $P(\langle M, X \rangle) > 0$ .

**3.10.** In the present general framework the properties of  $q_e$  can be hardly discussed. But under some particular conditions,  $q_e$  has a great advantage, namely  $p_e = p$ .

Consider the special case in which for each  $M$ ,  $\langle Q, \text{pol}^* \rangle$  is a sequence of independent and identically and independently of  $M$  distributed random variables (d-homogeneity and d-independence condition of [3]). Moreover, let  $Q$  and  $\text{pol}^*$  be mutually independent.

**3.11.** Under assumptions of 3.3, 3.4 and 3.10 we have the following: if  $p$  is independent of  $M$ , then  $p_e = P(\mathcal{N}_e) = p$  for each  $M$ .

**Proof.** Note that in the present context, all probabilities depend only on cardinality  $m$  of  $M$ , not on  $M$ . We can write  $P_m(\mathcal{N}_e) = \sum_{k=0}^m P_k(\{M \in \mathcal{M}_m; q(M) = 1\}) P_m(\{N; N \text{ has exactly } k \text{ rows without } \times\}) = p \sum_{k=0}^m P_m(\{N; N \text{ has exactly } k \text{ rows without } \times\}) = p$ .

**3.12.** For the usual quantifiers, the condition that  $p$  is independent of  $m$  is satisfied not quite exactly (usually  $p_m \rightarrow p$ ). This fact is due mainly to the non continuity of underlying distributions. The condition could be "easily" satisfied by randomization (under for example independency hypothesis and considering exact distributions).

**3.13.** The considerations of the present chapter give clearly some guidelines showing to us when use different ways of extension of quantifiers in the GUHA procedures.

(Received August 14, 1979.)

#### REFERENCES

- [1] P. Hájek: Automatic listing of important observational statements I—III, *Kybernetika* 9 (1973), 187—205, 251—271 and 10 (1974), 95—124.
- [2] P. Hájek, K. Bendová, Z. Renc: The GUHA method and the three valued logic. *Kybernetika* 7 (1971), 13—21.
- [3] P. Hájek, T. Havránek: Mechanizing hypothesis formation — mathematical foundations for a general theory. Springer-Verlag, Berlin—Heidelberg—New York, 1978.
- [4] P. Hájek, T. Havránek: The GUHA method — its aims and techniques. *Int. J. Man-Machine Studies* 10 (1978), 3—22.
- [5] P. Hájek, T. Havránek: A theory of helpful quantifiers and its application to approaches to missing information in the GUHA method. (In preparation for *Kybernetika*.)
- [6] P. Hájek, I. Havel, T. Havránek, M. Chytil, Z. Renc, J. Rauch: *Metoda GUHA*. Dům techniky, České Budějovice, sec. ed. 1977.
- [7] T. Havránek: Some alternative ways of treatment of missing information in GUHA pro-

cedures I—III, GUHA res. rep. 124, 125, 126, Czech. Acad. Sci. Cent. of Biomathematics, Prague 1977. 155

- [8] W. Lipski: Information systems with incomplete information. Proceedings of the third international symposium on automata theory, languages and programs, S. Michaelson, R. Milner (eds.). University Press, Edinburgh 1976, 120—130.

*RNDr. Tomáš Havránek, CSc., Matematické středisko biologických ústavů ČSAV (Center of Biomathematics — Czechoslovak Academy of Sciences), Vítězná 1083, 142 20 Prague 4, Czechoslovakia.*