

SIMILARITY IN THE GENERALIZED JSM METHOD AND ITS GENERATION ALGORITHMS

S. M. Gusakova and S. O. Kuznetsov

The article considers generation of hypotheses in the "generalized" JSM method, which consists in generation of "conditional" hypotheses about cause-and-effect relations for positive and negative examples. The conditionality of the hypotheses is that the causes about which the hypotheses are advanced "operate" (cause the manifestation of a property) in the absence of certain substructures called "brakes" of these causes. In contrast to "ordinary" negative causes expressing similarities of negative examples, the brakes possess some "cause-and-effect structure," bearing their own opposite: conditional positive causes. An algebraic description of similarity in the generalized JSM method is suggested, and also algorithms for generating generalized hypotheses (together with analysis of computational complexity).

1. THE GENERALIZED METHOD AND SIMILARITY

The basic idea of the JSM method [1,2], to establish cause-and-effect relations between substructures of structured objects and a subset of properties of these objects on the basis of determination of the objects' essential similarity, is naturally expanded in the generalized JSM method, the principle of which consists in taking into account the context (in the structure of the object) when establishing the cause-and-effect relations.

Such consideration of context is due to realities of the problems to be solved with the help of the JSM method from the field of chemistry, pharmacology, and sociology.

A formal definition of a generalized JSM predicate was suggested in [3]. Its informal meaning amounts to the following: in the presence of a description of some positive and negative examples for some class of objects (characterized by a particular property), "conditional" hypotheses are generated. Their conditionality is that the causes "operate" (cause manifestation of the property) in the absence of certain substructures called "brakes" of these causes. In contrast to "ordinary" (see [1,2]) negative causes expressing similarities of negative examples, the brakes possess some "causal (cause-and-effect) structure," bearing their own opposite: conditional positive causes. Thus, the binary relation $V \Rightarrow W$, which is read as "the subobject V is a cause of the presence of a set of properties W," of the simple JSM method is replaced by the ternary relation $T(V, \kappa, W)$, which is read as "V is a cause of the presence of a set of properties W in the absence of 'brakes' κ ," in the generalized JSM method.

Restrictions on the length of the article and the wish to avoid repeating the material set forth forces the authors to refer readers unfamiliar with the JSM method who want to get to know it to [1-4], which contain all of the necessary information.

However, the minimum information, concepts, and notations necessary for reading the present work are given below.

The JSM method of automatic hypothesis generation is used to analyze and process data represented in a data base with incomplete information. The class of problems to be solved by the JSM method is characterized by the following conditions:

- the subject area's data are well structured;
- there is a set of objects and a set of properties in each of which certain operations are assigned;
- in the set of objects, the relation "to be a subobject of an object" is determined;
- in the product of sets of objects and properties, the partially determined relations "an object possesses

a set of properties" (represented by $\overset{*}{\Rightarrow}_1$) and "a subobject is the cause of a set of properties" (represented by $\overset{*}{\Rightarrow}_2$) are assigned;

- there is a set of positive and negative examples of the first relation, i.e., a set of pairs of the type (x, y), where x is an object and y is a property, that satisfy and do not satisfy the relation;

The examples are obtained empirically.

The presence of properties of the objects or their absence is determined by "positive" (+) and "negative" (-) causes, i.e., the presence of some phenomenon is due to some set of (+) causes, and the absence of a phenomenon is due either to the presence of (-) causes or the absence of corresponding (+) causes.

As a rule, problems of this class arise in poorly formalized fields of knowledge.

The initial data for the JSM method are matrices of partially determined relations: $\overset{*}{\Rightarrow}_1$ and $\overset{*}{\Rightarrow}_2$. The rows of the matrices correspond to objects in the former case and subobjects in the latter; and the columns, to elementary properties. At the intersection of the i-th row and the j-th column, there is +1 if $x_i \overset{*}{\Rightarrow}_l y_j$ is fulfilled, -1 if it is not fulfilled, and τ if it is not known whether or not $x_i \overset{*}{\Rightarrow}_l y_j$ ($l = 1, 2$) takes place. The subrow of the matrix corresponding to all +1 (-1) is called a positive (negative) example.

With the help of logical combinatorial algorithms, the simple JSM method generates hypotheses of the type

$$J_{\langle v, n \rangle} (C \Rightarrow_1 A), J_{\langle v, n \rangle} (\bar{C} \Rightarrow_2 A), \\ J_{\langle \tau, n \rangle} (C \Rightarrow_1 A) \text{ and } J_{\langle \tau, n \rangle} (\bar{C} \Rightarrow_2 A),$$

where J is a one-place operator; and $v = (1, -1, 0)$ and τ are types of truth values denoting actual truth (+1), actual falsehood (-1), actual contradiction (0), and indeterminacy (τ).

With the help of significantly more complex algorithms, the generalized method, which reflects deeper ideas about the nature of cause-and-effect relations, generates hypotheses of the type

$$J_{\langle v, n \rangle} T(\bar{C}, \bar{x}, A) \text{ and } J_{\langle \tau, n \rangle} T(\bar{C}, \bar{x}, A),$$

where C and A are an object and a set of properties, respectively; \bar{C} is a subobject; and \bar{x} is a set of brakes.

In various complicated versions, the logical combinatorial algorithms of the JSM method realize a simple idea of J. S. Mill (hence the name JSM method): the cause of similarity of objects' properties is similarity of their structures.

Therefore, the concept of similarity is one of the central concepts of JSM theory, and determination of essential similarity in specific applied problems and for specific data structures is one of the central problems of creating an applied JSM system.

JSM systems use the operation of similarity to distinguish the causes of subobjects, and local (relation) and global (family of sets) similarities to construct algorithms.

In JSM methods, the similarity operation is understood as an idempotent, commutative, and associative operation on pairs of objects. These properties make it possible to unambiguously express similarity of a set of objects in terms of pair similarities, regardless of the order of the objects' arrangement in the data base.

The operation of similarity is represented as Π . A strict definition of it can be found, for example, in [5].

Definition 1. Objects X_{i_1}, \dots, X_{i_n} are locally similar if $\prod_{j=1}^n X_{i_j} \neq \emptyset$.

Local similarity is an n-ary tolerance relation (see [6]).

Definition 2. Objects X_{i_1}, \dots, X_{i_k} are globally similar if

$$\left(\prod_{j=1}^k X_{i_j} = h \right) \& \forall X_{i_m} \left(\left(\prod_{j=1}^k X_{i_j} \right) \Pi X_{i_m} = h \right) \& \\ \& \left(\&_{r=1}^k \neg (m=r) \right).$$

Global similarity links all objects that contain a certain subobject (result of the similarity operation) and therefore is not a relation.

The structure of global similarity is described by a family of sets (M, G) , where M is the set of objects under consideration; $G = \{g_i\}$; and g_i is the subset of globally similar objects. We recommend [7,8] to readers who want to familiarize themselves with the theory of global and local similarity.

For further exposition, it will be convenient for us to use the operational interpretation of local and global similarity (see [5]), in accordance with which:

Definition 1'. h is local similarity of X_1, \dots, X_n from S (S is a set of objects), if

$$\prod_{i=1}^n X_i = h.$$

Definition 2'. $\langle h, \{X_1, \dots, X_k\} \rangle$ is global similarity, if $X_1, \dots, X_k \in S$, $\prod_{i=1}^k X_i = h$, and $\forall Y \in S \setminus \{X_1, \dots, X_k\}$, $Y \cap h \neq h$ takes place.

The logical combinatorial algorithms realizing the simple JSM method are based on finding the global similarity of objects possessing some property.

Algorithms for realizing the generalized JSM method seek triads in the relation $T(V, \kappa, W)$. The set of triads of this relation is the region of truth of a generalized JSM similarity predicate, which is represented as $M_{ag,n}^+$ (we are considering the case of a positive predicate, remembering that all of the reasoning can be transferred to the case of a negative predicate $M_{ag,n}^-$ with no difficulty).

On the strength of the capaciousness of the indicated predicate's formulation, we will give it here in a simplified form. For an exact formulation, we refer the reader to [9], which is devoted to study of this predicate.

We will represent the predicate $M_{ag,n}^+(V, \kappa, W)$ in the form of a conjunction of several parts, to which we will give mnemonic notations for the sake of convenience:

$$M_{ag,n}^+(V, \kappa, W) = EX \& ED \& CE \& B/.$$

EX describes a set of examples of the type $J_{(1, n)}(X_i \Rightarrow Y_i)$ which are the basis of a plausible deduction;

ED is an empirical dependence between V, W , and κ expressing the fact that " V is a cause of W in the absence of brakes $\kappa = \{V_1, \dots, V_r\}$."

This dependence has the form

$$\forall X \forall Y ((J_{(1, n)}(X \Rightarrow Y) \& \forall U (J_{(1, n)}(X \Rightarrow Y) \rightarrow U \subseteq Y) \& V \subset C X) \& \neg (V_1 \subset X \vee \dots \vee V_r \subset X) \rightarrow W \subseteq Y \& W \neq \emptyset).$$

CE is the condition of exhaustibility, which provides for consideration of all suitable examples from the data base:

$$\bigvee_{i=1}^k (X = Z_i).$$

$B/$ describes the set of examples that generate brakes and includes two parts, $B/_1$ and $B/_2$, which describe sets of positive and negative examples, respectively:

$$\begin{aligned} B/_1: & \bigwedge_{j=1}^r W_j ((\neg J_{(1, n)}(Z_j \Rightarrow W_j) \& W \subseteq W_j) \& V \subset Z_j, \\ B/_2: & (J_{(1, n)}(Z_j \Rightarrow W_j) \& \neg (W \subseteq W_j)) \& V \subset Z_j, \\ B/ : & (B/_1 \vee B/_2) \& (V \neq \bigcap_{p=1}^l Z_{j_p}) \& V \subset V_q. \end{aligned}$$

The rest of the fragments of the predicate $M_{ag,n}^+$ characterize the conditions of exhaustibility and imaginarity for the brakes.

The essence of the predicate $M_{ag,n}^+$ is that a subobject found as operational similarity of objects X_1, \dots, X_k is a cause of the property W , if none of the indicated objects contains brakes from $\kappa = \{V_1, \dots, V_r\}$, $V \subset V_i$, $i = 1, \dots, r$. The brakes, in turn, are found from objects Z_1, \dots, Z_n , for negative, as well as positive examples containing V ,

but not possessing the property W.

If brakes for V cannot be extracted from the examples available in the data base, then V is not a cause of the property W in the generalized sense.

At the algorithmic level, the main difficulties are due to the fact that when actions checking the conditions of the predicate are performed in succession a subobject V found as a candidate for a cause of the property W may be thrown out. The next step in this case will be conducted on a set of objects X_1, \dots, X_n without the objects that contain brakes, and this can lead to the appearance of a new candidate for a cause, which, in turn, entails a search for new brakes. The presence in the predicate $M_{x,n}^*$ of a condition according to which brakes are also sought among positive examples that do not possess the property W complicates the algorithm even more.

Investigation of the properties determined in a set of objects by a predicate of generalized similarity provides the opportunity to distinguish different situations that arise depending on the structure of the data base, and among these situations to find those for which the algorithm realizing the generalized JSM method is simpler than in the general case.

Therefore, we will start this investigation.

Analyzing the predicate $M_{x,n}^*$ of generalized similarity from the point of view of the nature of similarities generated by it in a set of objects (so-called procedural similarities), we must note the following significant points:

- 1) procedural similarities in a set of objects are global similarities, on the strength of the condition of exhaustibility;
- 2) procedural similarity determining a cause is some composite of simple similarities. This follows from the trinariness of the predicate and the relations between v and x .

We will introduce the following notations:

$W = \{w_1, \dots, w_k\}$ is a set of elementary properties.

$\bar{W} = 2^W = \{\bar{w}_1, \dots, \bar{w}_{\bar{a}}\}$ is the set of all subsets of the set W.

$$\bar{G}^l = \{\bar{g}_k^l, k=1, \dots, \bar{k}; l=1, \dots, \bar{a}; \sigma = (+, -, \pm)\};$$

$\bar{g}_k^l = \{\bar{h}_k^l, \bar{S}_k^l\}$ is global similarity, with \bar{S}_k^l consisting of objects taken from the subset of positive examples of the type $J_{(1,0)}(X \Rightarrow_1 \bar{W}_i)$, \bar{S}_k^l of objects from the subset of negative examples also for the property \bar{w}_i , and \bar{S}_k^l of objects from the subset of positive, as well as negative examples of this property.

Thus, $\{\bar{h}_i^l\}$ is the set of all maximum positive "intersections" (or rather, the results of the similarity operation); and $\{\bar{h}^l\}$, of all negative ones, i.e.,

$$\begin{aligned} \bar{h}_k^l &= \prod_{x_n \in \bar{S}_k^l} X_n, \quad \bar{S}_k^l = \{X_n | (J_{(1,0)}(X_n \Rightarrow_1 Y) \& \\ &\& (\bar{W}_i \subseteq Y))\}; \\ D^{l,j} &= \{d_m^{l,j}\}, \quad m=1, \dots, \bar{m}; \quad l, j=1, \dots, \bar{a}; \\ l \neq j; \quad \bar{W}_i &\subseteq \bar{W}_j; \\ d_m^{l,j} &= \{\bar{h}_m^{l,j}, \bar{S}_m^{l,j}\}; \quad \bar{h}_m^{l,j} = \prod_{x_n \in \bar{S}_m^{l,j}} X_n; \\ \bar{S}_m^{l,j} &= \{X_n | ((J_{(1,0)}(X_n \Rightarrow_1 Y_1) \& (\bar{W}_i \subseteq Y_1)) \vee \\ &\vee ((J_{(1,0)}(X_n \Rightarrow_1 Y_2) \& (\bar{W}_j \subseteq Y_2)))\}, \end{aligned}$$

i.e., $\bar{S}_m^{l,j}$ consists of objects taken from positive examples for properties \bar{w}_i and \bar{w}_j , with $\bar{w}_i \not\subseteq \bar{w}_j$.

The need to consider the property $D^{l,j}$ is determined by the part of the predicate B_l for seeking brakes.

We will define two ternary predicates:

$$\begin{aligned} \zeta_1(\bar{g}_k^l, \bar{g}_n^l, \bar{g}_m^l) &= \\ = 1 &\Rightarrow (\bar{h}_k^l \subseteq \bar{h}_n^l) \& (\bar{h}_n^l = \bar{h}_m^l) \& \forall l \neg (\bar{h}_l^l \subseteq \bar{h}_n^l); \\ \zeta_2(\bar{g}_k^l, \bar{g}_n^l, d_m^{l,j}) &= \end{aligned}$$

$$= 1 \Leftrightarrow (\bar{h}_k^i \subset \bar{h}_n^i) \& (\hat{h}_n^i = h_m^i) \& \forall i' \neg (\hat{h}_i' \subset \hat{h}_n^i).$$

It is clear that ζ_1 and ζ_2 are equal to 0 if even one of the terms of the conjunction is equal to zero (1 and 0 are the truth values "truth" and "falseness").

We will define global similarity $\hat{p}^i = \{\hat{p}_k^i\}$, where $\hat{p}_k^i = (\hat{v}_k^i, \hat{\Sigma}_k^i)$ as follows:

$$\begin{aligned} \hat{v}_k^i &= \bar{h}_k^i, \\ \hat{\Sigma}_k^i &= \hat{S}_k^i \setminus \left[\left(\bigcup_{m,n} (\hat{S}_m^i \setminus \bar{S}_n^i) \mid \zeta_1 (\bar{g}_k^i, \bar{g}_n^i, \hat{g}_m^i) = 1 \right) \cup \right. \\ &\quad \left. \left(\bigcup_{m',n',j} (S_{m'}^i \setminus \hat{S}_{n'}^i) \mid \zeta_2 (\bar{g}_k^i, \hat{g}_{n'}^i, a_{m'}^i) = 1 \right) \right]; \\ \hat{p}_k^i &= \emptyset, \quad \text{if } (\forall n \forall m \zeta_1 (\bar{g}_k^i, \bar{g}_n^i, \hat{g}_m^i) = 0) \& \\ &\quad \& (\forall n' \forall m' \zeta_2 (\bar{g}_k^i, \hat{g}_{n'}^i, a_{m'}^i) = 0) \vee \exists r (\hat{\Sigma}_k^i = \hat{S}_r^i). \end{aligned}$$

Proposition 1. For any $\hat{p}_k^i = (\hat{v}_k^i, \hat{\Sigma}_k^i)$, $T(\hat{v}_k^i, \kappa, \bar{w}_i) = 1$ takes place, where

$$\begin{aligned} \kappa &= \{ \{ \bar{h}_n^i \mid \exists m \zeta_1 (\bar{g}_k^i, \bar{g}_n^i, \hat{g}_m^i) = 1 \} \cup \\ &\quad \cup \{ \bar{h}_{n'}^i \mid \exists m' \zeta_2 (\bar{g}_k^i, \hat{g}_{n'}^i, a_{m'}^i) = 1 \} \}. \end{aligned}$$

Proof.

1. From the definition of similarity \hat{G}^i , it follows that for each i and k the set \hat{S}_k^i constitutes exactly the set of examples described in the part *Ex* of the predicate $M_{ag,n}^+$. The condition of exhaustibility (CE) is provided by the globalness of the similarity \hat{g}_k^i . Consequently, $\hat{v}_k^i = \bar{h}_k^i$ satisfies the conditions of $M_{ag,n}^+$ since it is the maximum operational similarity of positive examples.

2. From the definition of \bar{g}_n^i , it follows that \bar{h}_n^i is the maximum operational similarity of negative examples. The condition of minimality \bar{h}_n^i with respect to embedding, which is imposed on the brakes, is provided by fulfillment of the predicate ζ_1 .

3. From the definition of d_m^i it follows that \hat{h}_n^i is the maximum operational similarity of positive examples that do not possess the property \bar{w}_i .

4. From the definition of $\hat{\Sigma}^i$ it follows that

$$\forall n \forall X \in \hat{\Sigma}_k^i \neg (\bar{h}_n^i \subset X) \& \forall n' \forall X' \in \hat{\Sigma}_k^i \neg (\hat{h}_{n'}^i \subset X').$$

Fulfillment of conditions (1-4) provides for $T(\hat{v}_k^i, \kappa, \bar{w}_i) = 1$ with κ described in the formulation of the proposition.

Now we will show that if $T(\hat{v}, \kappa', \bar{w}_i) = 1$, then $\exists k \hat{v} = \hat{h}_k^i$.

$$\begin{aligned} \exists n \forall m \exists n' \exists m' \kappa' &= \{ \{ \bar{h}_n^i \mid \zeta_1 (\bar{g}_k^i, \bar{g}_n^i, \hat{g}_m^i) = 1 \} \cup \\ &\quad \cup \{ \bar{h}_{n'}^i \mid \zeta_2 (\bar{g}_k^i, \hat{g}_{n'}^i, a_{m'}^i) = 1 \} \}. \end{aligned}$$

By the definition of the relation T , \hat{v} is the maximum operational similarity of some subset of positive examples, but all such similarities are included in the global similarity \hat{G}^i , consequently $\exists k \hat{v} = \hat{h}_k^i$.

Reasoning analogously, we find that

$$\forall v_q \in \kappa' (\exists n (v_q = \bar{h}_n^i)) \vee (\exists j \exists n' (v_q = \hat{h}_{n'}^i)).$$

Fulfillment of the predicates ζ_1 and ζ_2 follows from the definition of similarities \hat{G} and \hat{P} .

Strictly speaking, proposition 1 follows directly from the definition of the relation T , the predicate $M_{ag,n}^+$

and the similarities $\overset{\circ}{G}$ and $\overset{+}{P}$, which we can see directly from the proof.

Thus, the global similarity $\overset{\pm}{P}$ is generated by all those and only those objects the operational similarity of which is expressed by the cause $\overset{+}{v}$ of the property \bar{w} in the absence of brakes κ .

As we can see from the definition, the global similarity $\overset{\pm}{P}$ depends significantly on the property \bar{w} , therefore the structure of the set of properties in the data base largely determines the complexity of the algorithmic procedure for finding $\overset{\pm}{P}$. So, it is obvious that for the case $W = \{w\}$, $\bar{W} = W$ this procedure is significantly simplified. In fact, the predicate of similarity $M^*_{\kappa, n}$ is primarily simplified, since the addition B/l_2 is removed from it.

For this case, the similarity $\overset{\pm}{P} = \{p_j\}$ takes the form

$$\overset{\pm}{P}_j = \{ \overset{+}{v} = \overset{+}{h}_j; \overset{\pm}{S}_j = \overset{\pm}{S}_j \setminus \left(\bigcup_{m,n} (\overset{\pm}{S}_m \setminus \bar{S}_n) \mid \zeta_1, (\overset{+}{g}_j, \bar{g}_n, \overset{\pm}{g}_m) = 1 \right) \}.$$

$\overset{\pm}{P}$ is formed with the help of similarities $\overset{+}{G} = \{g_j\}$, $\bar{G} = \{\bar{g}_n\}$ and $\overset{\pm}{G} = \{g_m\}$.

If the similarity P is found, then all of the triads $\{v, \kappa, w\}$ satisfying the predicate T are found.

2. ALGORITHMS FOR CONSTRUCTING GENERALIZED HYPOTHESES

In essence, the way the definition of similarity $\overset{\pm}{P}$ is constructed incorporates algorithm 1 for finding it, which consists in a sequence of the following actions:

1. Finding similarities $\overset{+}{G}$ and \bar{G} , which is done with the help of the so-called algorithm of maximum intersections realized in the version of the JSM system for the simple method.

2. Verification of the predicate ζ_1 .

3. Construction of the similarity $\overset{\pm}{P}$.

4. Verification of the condition $\forall r (\overset{\pm}{S}_r \neq \bar{S}_r)$.

We will demonstrate the algorithm's operation for the case of one property on the following example:

A set of objects: $\{x_1, \dots, x_{11}\}$; $\Omega^+ = \{x_1, x_2, x_3, x_4, x_5\}$, $\Omega^- = \{x_6, x_7, \dots, x_{11}\}$; $x_1 = vad'$, $x_2 = vbd'd'$, $x_3 = vdb'd''m$, $x_4 = vd'cm$, $x_5 = vd''c'm$, $x_6 = vae$, $x_7 = vac'$, $x_8 = vbd'f$, $x_9 = vbd'f$, $x_{10} = vd'gg'$, $x_{11} = vd'gg''$.

We find the similarity $\overset{+}{G}$ from positive examples:

$$\begin{aligned} \overset{+}{g}_1 &= \{v; x_1, x_2, x_3, x_4, x_5\}; \\ \overset{+}{g}_2 &= \{vbd'd'; x_1, x_2\}; \\ \overset{+}{g}_3 &= \{vd'; x_2, x_3, x_4\}; \\ \overset{+}{g}_4 &= \{vm; x_3, x_4, x_5\}; \\ \overset{+}{g}_5 &= \{vd''m; x_3, x_4\}. \end{aligned}$$

We find the similarity \bar{G} from negative examples:

$$\begin{aligned} \bar{g}_1 &= \{v; x_6, x_7, x_8, x_9, x_{10}, x_{11}\}; \\ \bar{g}_2 &= \{va; x_6, x_7\}; \\ \bar{g}_3 &= \{vbd; x_6, x_8\}; \\ \bar{g}_4 &= \{vd'g; x_{10}, x_{11}\}. \end{aligned}$$

The similarity $\overset{\pm}{G}$ is represented in the following way:

$$\begin{aligned}
\bar{g}_1^{\pm} &= \{v; x_1, \dots, x_{11}\}; \\
\bar{g}_2^{\pm} &= \{va; x_1, x_8, x_7\}; \\
\bar{g}_3^{\pm} &= \{vb; x_2, x_3, x_9, x_8\}; \\
\bar{g}_4^{\pm} &= \{vbd'; x_1, x_2\}; \\
\bar{g}_5^{\pm} &= \{vbd; x_2, x_8, x_9\}; \\
\bar{g}_6^{\pm} &= \{vd'; x_1, x_2, x_4, x_{10}, x_{11}\}; \\
\bar{g}_7^{\pm} &= \{vd''m; x_9, x_8\}; \\
\bar{g}_8^{\pm} &= \{vd'g; x_{10}, x_{11}\}.
\end{aligned}$$

We will let Th_{ζ} represent the set of triads $\{\bar{g}_i^{\pm}, \bar{g}_j^{\pm}, \bar{g}_m^{\pm}\}$ for which ζ^{\pm} is true. In the given example, it is the following:

$$\begin{aligned}
Th_{\zeta^{\pm}} &= \{(\bar{g}_1^{\pm}, \bar{g}_2^{\pm}, \bar{g}_3^{\pm}), (\bar{g}_1^{\pm}, \bar{g}_3^{\pm}, \bar{g}_4^{\pm}), (\bar{g}_1^{\pm}, \bar{g}_4^{\pm}, \bar{g}_5^{\pm}), \\
&(\bar{g}_2^{\pm}, \bar{g}_4^{\pm}, \bar{g}_5^{\pm})\}.
\end{aligned}$$

Now we can construct the similarities \bar{p}_i^{\pm} :

$$\begin{aligned}
\bar{p}_1^{\pm} &= \{v; \bar{\Sigma}_1 \{ \{x_1, x_2, x_3, x_4, x_8\} \setminus \{ \{x_1, x_8, x_7 \setminus x_8, x_7\} \cup \\
&\cup \{x_2, x_3, x_9 \setminus x_2, x_9\} \cup \{x_{10}, x_{11} \setminus x_{10}, x_{11}\} \} \} = \\
&= \{x_2, x_4, x_8\}.
\end{aligned}$$

But since $\bar{\Sigma}_1 = \bar{\Sigma}_4$, then $\bar{p}_1^{\pm} = \emptyset$.

$\bar{p}_2^{\pm} = \emptyset$, since $\forall j \neg (\bar{h}_2 \subset \bar{h}_j)$ $j=1, \dots, 4$.

$\bar{p}_3^{\pm} = \{vd'; \bar{\Sigma}_3 \{ \{x_2, x_3, x_4\} \setminus \{x_{10}, x_{11} \setminus x_{10}, x_{11}\} \} = \{x_2, x_3, x_4\}$.

$\forall l \bar{\Sigma}_1 \neq \bar{\Sigma}_l$ ($l=1, \dots, 5$), consequently $\bar{p}_4^{\pm} \neq \emptyset$.

$\bar{p}_5^{\pm} = \emptyset$ and $\bar{p}_6^{\pm} = \emptyset$ for the same reason that $\bar{p}_1^{\pm} = \emptyset$.

As a result of the algorithm's operation, we find one triad: $\{vd', \kappa = vd'g, w\}$ that satisfies the predicate $M_{\alpha, \beta}^*$ and, consequently, belongs to the relation T.

An increase in the set of elementary properties even to two elements entails a significant increase in the diversity of situations and, consequently, complication of the algorithm.

We will figure an upper estimate of the number of computer operations necessary to generate generalized hypotheses with the help of the indicated procedure for an arbitrary set of positive examples Ω^+ , set of negative examples Ω^- , and set of elementary structures U (subsets of which are positive and negative examples). In [5], the ZO algorithm was suggested for seeking intersections (similarities), the time requirement of which is linear in relation to the number of intersections generated. Supposing that similarities (intersections) are sought with the help of the ZO algorithm, we will determine that finding the set \bar{G}^+ of all positive similarities requires spending

$$O(|\bar{G}^+| \cdot |U| \cdot |\Omega^+|)$$

computer operations; finding the set \bar{G}^- of all negative similarities requires spending

$$O(|\bar{G}^-| \cdot |U| \cdot |\Omega^-|)$$

computer operations; and finding the set of similarities \bar{G}^{\pm} of all examples (positive and negative) requires spending

$$\begin{aligned}
&O(|\bar{G}^{\pm}| \cdot |U| \cdot (|\Omega^+| + |\Omega^-|)) = \\
&= O(|\bar{G}^+| \cdot |\bar{G}^-| \cdot |U| \cdot (|\Omega^+| + |\Omega^-|))
\end{aligned}$$

computer operations. In the worst case, generating each set $\tilde{\Sigma}_j$ requires

$$\begin{aligned} O(|\tilde{G}_j| \cdot |\tilde{G}| \cdot \text{pol}_1(|U|, |\Omega^+|, |\Omega^-|)) = \\ = (|\tilde{G}|^2 \cdot |\tilde{G}| \cdot \text{pol}_1(|U|, |\Omega^+|, |\Omega^-|)) \end{aligned}$$

computer operations, where $\text{pol}_1(\cdot)$ is some small-degree polynomial depending on $|U|$, $|\Omega^+|$, and $|\Omega^-|$ and determined by a computer model.

We can suggest more efficient algorithms for constructing generalized hypotheses. The speed of these algorithms depends significantly on the type of similarity operation (even if one considers the time required to perform such operations to be the same). First we will give a description of such an algorithm for constructing generalized positive hypotheses in the case of representation of data by sets and the presence of just one property. Algorithm 2

1. All (+) intersections are constructed with the help of the ZO algorithm [5].

2. For each (+) intersection V:

2.1. Negative examples containing it are found.

2.2. All intersections of negative examples that are minimal with respect to embedding and contain V (i.e., brakes) are sought: Let N be the union of all negative examples containing V. Then the arbitrary minimal intersection of negative examples containing V has the form $X^j = \bigcap X_{t_j}$, where $X_{t_j} \in X_t = \{X | X \in \Omega^-, (V \cup \{q_i\}) \subseteq X\}$, where $q_i \in N$. Naturally, $\bigcup_j X_{t_j}$ is determined, if $j \geq 2$.

In order to rule out generation of the same intersections and thus to accelerate the process of generating them, in constructing the next set $\bigcup_j X_{t_j}$, the corresponding "addition" to v can be chosen not from N, but from N minus all of the preceding "additions" and generated intersections, or formally, $N_{k+1} = N_k \setminus (\{q_{k+1}\} \cup \bigcap_j X_{t_j})$.

2.3. From the set of (+) examples containing V, all of those that contain any brakes are removed, and the remaining (+) examples are reintersected. If the result coincides with V, then V, together with the minimal intersections found, is the generalized hypothesis; if not, then no (+) generalized hypothesis exists in relation to the intersection V, and we move on to the next (+) intersection.

The time complexity for constructing one generalized positive hypothesis for a fixed (+) intersection V in steps 2.1-2.3 of this algorithm is

$$O(N \cdot |\Omega^-|).$$

Thus, the complexity of constructing all positive generalized hypotheses is

$$O(|\tilde{G}| \cdot N \cdot |\Omega^-|).$$

We will consider the algorithm's action for the sets of positive and negative examples given above. In this case, the set N is $N = \{a, b, d, d', e, f, f', g, g', g'', v\}$. Suppose that in the second stage of the algorithm's action we chose the (+) intersection v. Then the minimal intersections of (-) examples containing v are found as

$$X^a = \bigcap X_j^a, \text{ where } X_j^a \in X_a = \{X \in \Omega^- | \{v, a\} \subset X\} = \{v, a\},$$

$$N_1 = N \setminus \{a\} = \{b, d, d', e, e'f, f', g, g', g''\};$$

$$X^b = \bigcap X_j^b, \text{ where } X_j^b \in X_b = \{X \in \Omega^- | \{v, b\} \subset X\} = \{v, b, d\},$$

$$N_2 = N_1 \setminus \{b, d\} = \{d', e, e'f, f', g, g', g''\}.$$

The intersection corresponding to $\{v, d\}$ is not generated, since in this step $d \notin N$.

$$X^{d'} = \bigcap X_j^{d'}, \text{ where } X_j^{d'} \in X_{d'} = \{X \in \Omega^- | \{v, d'\} \subset X\} =$$

$$= \{v, d', g\}, N_3 = N_2 \setminus \{d', g\} = \{e, e'f, f', g', g''\};$$

$$X^e = \bigcap X_j^e, \text{ where } X_j^e \in X_e = \{X \in \Omega^- | \{v, e\} \subset X\},$$

$$X^e \text{ is not determined, } N_4 = N_3 \setminus \{e\} = \{e'f, f', g', g''\};$$

$$X^{e'} = \bigcap_j X_j^{e'}, \text{ where } X_j^{e'} \in X_{e'} = \{X \in \Omega - \{v, e'\} \subset X\},$$

$$X^{e'} \text{ is not determined. } N_6 = N_5 \setminus \{e'\} = \{f, f', g', g''\}.$$

$$X^f = \bigcap_j X_j^f, \text{ where } X_j^f \in X_f = \{X \in \Omega - \{v, f\} \subset X\},$$

$$X^f \text{ is not determined. } N_7 = N_6 \setminus \{f\} = \{f', g', g''\};$$

$$X^{f'} = \bigcap_j X_j^{f'}, \text{ where } X_j^{f'} \in X_{f'} = \{X \in \Omega - \{v, f'\} \subset X\},$$

$$X^{f'} \text{ is not determined. } N_8 = N_7 \setminus \{f'\} = \{g', g''\}.$$

The minimal intersection, which corresponds to $\{v, g\}$, is not generated, since $g \notin N_8$,

$$X^{g'} = \bigcap_j X_j^{g'}, \text{ where } X_j^{g'} \in X_{g'} = \{X \in \Omega - \{v, g'\} \subset X\},$$

$$X^{g'} \text{ is not determined. } N_9 = N_8 \setminus \{g'\} = \{g''\};$$

$$X^{g''} = \bigcap_j X_j^{g''}, \text{ where } X_j^{g''} \in X_{g''} = \{X \in \Omega - \{v, g''\} \subset X\},$$

$$X^{g''} \text{ is not determined. } N_{10} = N_9 \setminus \{g''\} = \emptyset.$$

Actions relating to step 2.3 remain to be carried out. From the set of (+) examples containing the intersection v , i.e., from $\{x_1, x_2, x_3, x_4, x_5\}$, we will remove all examples containing the minimal intersections (-): examples containing v , i.e., va , vbd , and $vd'gg'$. Such (+) examples will be $x_1 = vaa'$ and $x_2 = vbdd'$. We will intersect the remaining (+) examples, i.e., x_3, x_4 , and x_5 , and get $vm \neq v$. Hence, there will be no generalized (+) hypothesis corresponding to the similarity of (+) examples v .

We will carry out the computations of steps 2.1-2.3 for the (+) intersection vd' . The (-) examples containing it will be $x_{10} = vd'gg'$ and $x_{11} = vd'gg''$. Hence, the sole minimal intersection of (-) examples containing vd' will be $vd'g'$. Since none of the (+) examples contain it, in step 2.3 none of the (+) examples will be thrown out, and the triad $\{vd', k = vd'g, w\}$ will be accepted as a generalized (+) hypothesis.

As we already said above, algorithm 2 is intended for data represented by sets. This circumstance makes it possible to accomplish step 2.2 very efficiently. When the data are represented by objects from an arbitrary semilattice with similarity operation Π , the situation is worse: we cannot form an intersection of (-) examples that is minimal with respect to embedding and contains the given intersection of (+) examples v by "building up" v with atoms of the lattice of (-) intersections (for information about the lattice of intersections, see, for example, [5]) - we simply cannot know these atoms. Thus, for an arbitrary similarity operation, instead of step 2.2 of algorithm 2 we must seek the minimal intersections "from the top down:" by taking the intersections of a growing number of (-) examples containing v . As soon as the intersection of (-) examples becomes equal to v in the next step, we fall back on the preceding intersection of (-) examples: it will be minimal among those containing v . In this case, the upper estimate of the number of computer operations necessary to generate all of the minimal intersections containing v will be $O(|\bar{G}| \text{pol}_2(|U|, |\Omega^+|, |\Omega^-|))$; consequently, the algorithm's overall complexity will be $O(|G| \cdot |\bar{G}| \text{pol}_2(|U\Omega|, |\Omega^+|, |\Omega^-|))$, where $\text{pol}_2(\cdot)$ is some small-degree polynomial depending on $|U|$, $|\Omega^+|$, and $|\Omega^-|$ and determined by a computer model.

REFERENCES

1. V. K. Finn, "Plausible deductions and plausible reasoning," in: Results of Science and Technology. Probability Theory, Mathematical Statistics, and Theoretical Cybernetics Series [in Russian], vol. 28, pp. 3-34, VINITI, Moscow, 1988.
2. S. O. Kuznetsov, "The JSM method as an automatic teaching system," in: Results of Science and Technology, Information Science Series [in Russian], vol. 15, pp. 17-54, VINITI, Moscow, 1991.
3. V. K. Finn, "The generalized JSM method of automatic hypothesis generation," Semiotika i Informatika, no. 29, pp. 93-123, 1989.
4. S. M. Gusakova and V. K. Finn, "Similarity and plausible deduction," Izv. AN SSSR, Ser. Tekhn. Kibernetika, vol. 5, pp. 42-63, 1987.
5. S. O. Kuznetsov, "Quick algorithm for constructing all intersections of objects from a finite semilattice," Nauchno-Tekhnicheskaya Informatsiya, Ser. 2, no. 1, pp. 17-20, 1993.

6. S. M. Gusakova, Canonical representations of similarities," *Nauchno-Tekhnicheskaya Informatsiya*, Ser. 2, no. 9, pp. 19-21, 1987.
7. S. M. Gusakova and V. K. Finn, "New means of formalizing local and global similarities," *Nauchno-Tekhnicheskaya Informatsiya*, Ser. 2, no. 10, pp. 14-22, 1987.
8. S. M. Gusakova, "Formalization of the concept of similarity and its application in intelligent information systems," Author's abstract of candidate's dissertation, VINITI, Moscow, 1988.
9. V. K. Finn and M. A. Mikheenkova, "Some problems of the generalized JSM method of automatic hypothesis generation," *Semiotika i Informatika*, no. 33, pp. 136-163, 1983.

5 May 1995

THE ALLERTON PRESS JOURNAL PROGRAM

AUTOMATIC DOCUMENTATION & MATHEMATICAL LINGUISTICS

Selected major articles from

NAUCHNO-TEKHNICHESKAYA INFORMATSIYA

Seriya 2. Informatsionnye Protsessy i Sistemy

Editor: R. S. Gilyarevskii

Executive Secretary: N. P. Zhukova

Editorial Board:

G. T. Artamonov	G. S. Pospelov
G. G. Belonogov	S. A. Rozhkov
I. A. Boloshin	A. Ya. Rodionov
I. A. Bol'shakov	E. P. Semenyuk
O. E. Buryi-Shmar'yan	V. R. Serov
A. V. Butrimenko	S. S. Tereshchenko
O. I. Voverene	A. D. Ursul
B. M. Gerasimov	V. A. Uspenskii
V. A. Gubanov	Yu. Yu. Ukhin
N. N. Ermoshenko	A. I. Chernyi
O. V. Kedrovskii	O. B. Shatberashvili
V. P. Leonov	A. V. Shileiko
Yu. N. Marchuk	Yu. A. Shreider
V. F. Medvedev	M. G. Yaroshevskii
V. K. Popov	

© Vsesoyuznyi Institut Nauchnoi i Tekhnicheskoi Informatsii, 1995

© 1996 by Allerton Press, Inc.

Published six times a year

All rights reserved. This publication or parts thereof may not be reproduced in any form without permission of the publisher.

ALLERTON PRESS, INC.
150 Fifth Avenue New York, N.Y. 10011

ISSN 0005-1055

AUTOMATIC DOCUMENTATION AND MATHEMATICAL LINGUISTICS

**(Nauchno-Tekhnicheskaya
Informatsiya, Seriya 2)**

Vol. 29, No. 3

ALLERTON PRESS, INC.