

С. М. Гусакова, С. О. Кузнецов

СХОДСТВО В ОБОБЩЕННОМ ДСМ-МЕТОДЕ И АЛГОРИТМЫ ЕГО ПОРОЖДЕНИЯ

Рассматривается порождение гипотез в «обобщенном» ДСМ-методе, заключающемся в порождении «условных» гипотез о причинно-следственных связях по положительным и отрицательным примерам. Условность гипотез состоит в том, что причины, о которых выдвигаются гипотезы, «срабатывают» (вызывают проявление свойства) в отсутствие некоторых подструктур, называемых «тормозами» этих причин. В отличие от «обычных» отрицательных причин, выражающих сходства отрицательных примеров, тормоза обладают некоторой «причинно-следственной структурой», неся в себе свою противоположность — условные положительные причины. В работе предложено алгебраическое описание сходства в обобщенном ДСМ-методе, а также алгоритмы порождения обобщенных гипотез (вместе с анализом вычислительной сложности).

1. ОБОБЩЕННЫЙ МЕТОД И СХОДСТВО

Основная идея ДСМ-метода [1, 2] — установление причинно-следственных отношений между подструктурами структурированных объектов и подмножеством свойств этих объектов на базе определения существенного сходства объектов — получает естественное расширение в обобщенном ДСМ-методе, принцип которого состоит в учете контекста (в структуре объекта) при установлении причинно-следственного отношения.

Такой учет контекста вызван реалиями решаемых с помощью ДСМ-метода задач из области химии, фармакологии, социологии.

В работе [3] было предложено формальное определение обобщенного ДСМ-предиката. Неформальный смысл его сводится к следующему: при наличии описания некоторых положительных и отрицательных примеров для некоторого класса объектов (характеризуемого определенным свойством) порождаются «условные» гипотезы. Условность их состоит в том, что причины «срабатывают» (вызывают проявление свойства) в отсутствие некоторых подструктур, называемых «тормозами» этих причин. В отличие от «обычных» (см. [1, 2]) отрицательных причин, выражающих сходства отрицательных примеров, тормоза обладают некоторой «каузальной (причинно-следственной) структурой», неся в себе свою противоположность — условные положительные причины. Таким образом бинарное отношение $V \Rightarrow_2 W$, читаемое как «подобъект V есть причина наличия множества свойств W » простого ДСМ-метода заменяется на тернарное отношение $T(V, x, W)$, читаемое как « V есть причина наличия множества свойств W при отсутствии «тормозов» x » в обобщенном ДСМ-методе.

Ограничения объема статьи и нежелание повторять неоднократно изложенный материал вынуждает авторов отослать читателя, незнакомого с ДСМ-методом и имеющего желание познакомиться с ним, к работам [1—4], содержащим все необходимые сведения.

Однако необходимый для прочтения настоящей работы минимум сведений, понятий и обозначений приводится ниже.

ДСМ-метод автоматического порождения гипотез используется для анализа и обработки данных, представленных в базе данных с неполной информацией. Класс

задач, решаемый ДСМ-методом характеризуется следующими условиями:

данные предметной области хорошо структурированы;

имеется множество объектов и множество свойств, на каждом из которых заданы определенные операции; на множестве объектов определено отношение «быть подобъектом объекта»;

на произведении множеств объектов и свойств заданы частично определенные отношения «объект обладает множеством свойств» (обозначается через \Rightarrow_1^*) и «подобъект является причиной множеств свойств» (обозначается через \Rightarrow_2^*);

имеется множество положительных и отрицательных примеров первого отношения, т. е. множество пар вида (x, y) , где x — объект, y — свойство, удовлетворяющих и неудовлетворяющих отношению.

Примеры получены эмпирически.

Наличие свойств у объектов или их отсутствие определяется «положительными» (+) и «отрицательными» (−) причинами, т. е. наличие некоторого явления обусловлено некоторым множеством (+)-причин, а отсутствие явления обусловлено или наличием (−)-причин или отсутствием соответствующих (+)-причин.

Как правило, задачи указанного класса возникают в плохо формализуемых областях знаний.

Исходными данными для ДСМ-системы являются матрицы частично определенных отношений \Rightarrow_1^* и \Rightarrow_2^* . Строки матриц соответствуют объектам в первом случае и подобъектам — во втором, столбцы — элементарным свойствам. На пересечении i -й строки и j -го столбца стоит +1, если $x_i \Rightarrow_j^*$ выполнено, −1,

если не выполнено и т, если не известно, имеет ли место $x_i \Rightarrow_j^*$ ($l = 1, 2$). Подстрока матрицы, соответствующая всем +1 (−1), называется положительным (отрицательным) примером.

С помощью логико-комбинаторных алгоритмов простой ДСМ-метод порождает гипотезы вида:

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

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

где J — одноместный оператор, $v = \{1, -1, 0\}$ и τ — типы истинностных значений, обозначающих фактическую истину (+1), фактическую ложь (-1), фактическое противоречие (0), недоопределенность (τ).

Обобщенный метод, отражающий более глубокие представления о характере причинно-следственных зависимостей, порождает с помощью существенно более сложных алгоритмов гипотезы вида:

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

где C и A есть соответственно объект и множество свойств, \bar{C} — подобъект, \bar{x} — множество тормозов.

Логико-комбинаторные алгоритмы ДСМ-метода, с помощью которых порождаются гипотезы, реализуют в различных усложненных вариантах простую идею Д. С. Милля (отсюда и название ДСМ-метод) — причинной сходства свойств объектов является сходство их структур.

Поэтому понятие сходства является одним из центральных понятий ДСМ-теории, а определение существенного сходства в конкретных прикладных задачах и для конкретных структур данных — одной из центральных проблем создания прикладной ДСМ-системы.

В ДСМ-системах используется операция сходства для выделения причин подобъектов и локальное (отношение) и глобальное (семейство множеств) сходства для конструирования алгоритмов.

Операция сходства понимается в ДСМ-системах как идемпотентная, коммутативная и ассоциативная операция на парах объектов. Эти свойства позволяют однозначно выражать сходство множества объектов через попарные сходства независимо от порядка расположения объектов в базе данных.

Операция сходства обозначается через Π . Строгое ее определение можно найти, например, в [5].

Определение 1. Объекты X_1, \dots, X_n локально сходны, если $\prod_{j=1}^n X_j \neq \emptyset$.

Локальное сходство есть n -арное отношение толерантности (см. [6]).

Определение 2. Объекты X_1, \dots, X_k глобально сходны если

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

Глобальное сходство объединяет все объекты, содержащие некоторый подобъект (результат операции сходства) и потому не является отношением.

Структура глобального сходства описывается семейством множеств (M, G) , где M — рассматриваемое множество объектов, $G = \{g_i\}$, g_i — подмножество глобально сходных объектов. Для ознакомления с теорией глобального и локального сходства рекомендуем читателю [7, 8].

Для дальнейшего изложения нам будет удобно пользоваться операциональной трактовкой локального и глобального сходства (см. [5]), в соответствии с которой:

Определение 1'. h есть локальное сходство X_1, \dots, X_n из S (S — множество объектов), если $\prod_{i=1}^n X_i = h$.

Определение 2'. $\langle h, \{X_1, \dots, X_k\} \rangle$ есть глобаль-

ное сходство, если $X_1, \dots, X_k \in S$, $\prod_{i=1}^k X_i = h$ и $\forall V \in S \setminus \{X_1, \dots, X_k\}$, имеет место $V \Pi h \neq h$.

Логико-комбинаторные алгоритмы, реализующие простой ДСМ-метод, основаны на нахождении глобального сходства объектов, обладающих некоторым свойством.

Алгоритмы для реализации обобщенного ДСМ-метода ищет тройки, находящиеся в отношении $T(V, x, W)$. Множество троек этого отношения есть область истинности обобщенного ДСМ-предиката сходства, обозначаемого через $M_{ag,n}^+$ (мы рассматриваем случай положительного предиката, помня, что все рассуждения могут быть без труда перенесены на случай отрицательного предиката $M_{ag,n}^-$).

В силу объемности формулировки указанного предиката приведем его здесь в упрощенном виде, отсылая читателя за точной формулировкой к статье [9], которая посвящена изучению этого предиката.

Представим предикат $M_{ag,n}^+(V, x, W)$ в виде конъюнкции нескольких частей, которым для удобства дадим мнемонические обозначения

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

EX описывает множество примеров вида $J_{(1,n)}(X_i \Rightarrow Y_i)$, которые являются основанием правдоподобного вывода;

ED — эмпирическая зависимость между V, W и x , выражающая тот факт, что « V является причиной W при отсутствии тормозов $x = \{V_1, \dots, V_r\}$ ».

Эта зависимость имеет вид:

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

CE — условие исчерпываемости, обеспечивающее рассмотрение всех подходящих примеров из базы данных:

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

BI — описывает множество примеров, порождающих тормоза и включает в себя две части BI_1 и BI_2 , описывающие множества соответственно положительных и отрицательных примеров:

$$BI_1: \bigwedge_{j=1}^s W_j ((\neg J_{(1,n)}(Z_j \Rightarrow W_j) \& W \subseteq W_j) \& V \subset Z_j,$$

$$BI_2: (J_{(1,n)}(Z_j \Rightarrow W_j) \& \neg (W \subseteq W_j)) \& V \subset Z_j,$$

$$BI: (BI_1 \vee BI_2) \& \left(V \supseteq \bigcap_{p=1}^i Z_p \right) \& V \subset V \supseteq.$$

Остальные фрагменты предиката $M_{ag,n}^+$ характеризуют условия исчерпываемости и минимальности для тормозов.

Суть предиката $M_{ag,n}^+$ в том, что подобъект, найденный как операциональное сходство объектов X_1, \dots, X_k , является причиной свойства W , если ни один из указанных объектов не содержит тормозов из $x = \{V_1, \dots, V_r\}$, $V \subset V_i$, $i = 1, \dots, r$. Тормоза, в свою очередь, находятся из объектов Z_1, \dots, Z_s как для отрицательных примеров, так и для положительных примеров, содержащих V , но не обладающих свойством W .

Причем, в случае, если тормозов для V из имеющихся в базе данных примеров извлечь нельзя, то V не является причиной свойства W в обобщенном смысле.

В алгоритмическом плане основные трудности порождает тот факт, что при последовательном выполнении действий по проверке условий предиката, найденный как кандидат в причины свойства W подобъект V , может быть забракован. Следующий шаг в этом случае будет проведен над множеством объектов X_1, \dots, X_n без тех объектов, которые содержат тормоза, и может привести к появлению нового кандидата в причины, а это, в свою очередь, повлечет поиск новых тормозов. Наличие в предикате $M_{ag,n}^+$ условия, по которому тормоза ищутся и среди положительных примеров, не обладающих свойством W , еще более усложняют алгоритм.

Исследование сходств, определяемых на множестве объектов предикатом обобщенного сходства, дает возможность выделить различные ситуации, возникающие в зависимости от структуры базы данных, и среди этих ситуаций найти те, для которых алгоритм, реализующий обобщенный ДСМ-метод, проще, чем в общем случае.

Поэтому мы и приступим к этому исследованию.

Анализируя предикат $M_{ag,n}^+$ обобщенного сходства, с точки зрения характера сходств, порождаемых им на множестве объектов (так называемые процедурные сходства), мы должны отметить следующие существенные моменты:

1) процедурные сходства на множестве объектов — это глобальные сходства в силу условия исчерпываемости;

2) процедурное сходство, определяющее причину, есть некоторая композиция простых сходств. Это следует из тернарности предиката и соотношений между σ и κ .

Введем следующие обозначения:

$W = \{w_1, \dots, w_a\}$ — множество элементарных свойств.

$\bar{W} = 2^W = \{\bar{w}_1, \dots, \bar{w}_a\}$ — множество всех подмножеств множества W .

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

$\bar{g}_k^i = \{\bar{h}_k^i, \bar{S}_k^i\}$ — глобальное сходство, причем \bar{S}_k^i состоит из объектов, взятых из подмножества положительных примеров вида $J_{\langle 1,0 \rangle} (X_n \Rightarrow \bar{W}_i)$, \bar{S}_k^i — из объектов подмножества отрицательных примеров также для свойства \bar{w}_i , а \bar{h}_k^i — из объектов подмножества как положительных, так и отрицательных примеров для этого же свойства.

Таким образом $\{\bar{h}_k^i\}$ — множество всех максимальных положительных «пересечений» (точнее результатов операции сходства), $\{\bar{h}_k^i\}$ — всех отрицательных, т. е.

$$\bar{h}_k^i = \prod_{X_n \in \bar{S}_k^i} X_n, \quad \bar{S}_k^i = \{X_n \mid (J_{\langle 1,0 \rangle} (X_n \Rightarrow Y) \& (\bar{W}_i \subseteq Y))\};$$

$$D^{i,j} = \{d_m^{i,j}\}, \quad m=1, \dots, \bar{m}; \quad i, j=1, \dots, \bar{a};$$

$$i \neq j; \quad \bar{W}_i \subseteq \bar{W}_j;$$

$$d_m^{i,j} = \{\bar{h}_m^{i,j}, \bar{S}_m^{i,j}\}; \quad \bar{h}_m^{i,j} = \prod_{X_n \in \bar{S}_m^{i,j}} X_n;$$

$$\bar{S}_m^{i,j} = \{X_n \mid ((J_{\langle 1,0 \rangle} (X_n \Rightarrow Y_1) \& (\bar{W}_i \subseteq Y_1)) \vee ((J_{\langle 1,0 \rangle} (X_n \Rightarrow Y_2) \& (\bar{W}_j \subseteq Y_2))),$$

т. е. $\bar{S}_m^{i,j}$ состоит из объектов, взятых из положительных примеров для свойств \bar{w}_i и \bar{w}_j , причем $\bar{w}_i \not\subseteq \bar{w}_j$.

Необходимость рассмотрения сходства $D^{i,j}$ определяется частью Bl_2 предиката для поиска тормозов.

Определим два тернарных предиката:

$$\zeta_1(\bar{g}_k^i, \bar{g}_n^i, \bar{g}_m^i) =$$

$$= 1 \Leftrightarrow (\bar{h}_k^i \subseteq \bar{h}_n^i) \& (\bar{h}_n^i = \bar{h}_m^i) \& \forall l \neg (\bar{h}_l^i \subseteq \bar{h}_n^i);$$

$$\zeta_2(\bar{g}_k^i, \bar{g}_n^j, d_m^{i,j}) =$$

$$= 1 \Leftrightarrow (\bar{h}_k^i \subseteq \bar{h}_n^j) \& (\bar{h}_n^j = \bar{h}_m^{i,j}) \& \forall l \neg (\bar{h}_l^i \subseteq \bar{h}_n^j).$$

Ясно, что ζ_1 и ζ_2 равны 0, если хотя бы один из членов конъюнкции равен 0 (1, 0 — истинностные значения «истина» и «ложь»).

Глобальное сходство $\bar{p}_k^i = \{\bar{v}_k^i, \bar{\Sigma}_k^i\}$ определим следующим образом:

$$\bar{v}_k^i = \bar{h}_k^i,$$

$$\bar{\Sigma}_k^i = \bar{S}_k^i \setminus \left[\left(\bigcup_{m,n} (\bar{S}_m^i \setminus \bar{S}_n^i) \mid \zeta_1(\bar{g}_k^i, \bar{g}_n^i, \bar{g}_m^i) = 1 \right) \cup \right.$$

$$\left. \left(\bigcup_{m,n',j} (\bar{S}_m^i \setminus \bar{S}_{n'}^j) \mid \zeta_2(\bar{g}_k^i, \bar{g}_n^j, d_m^{i,j}) = 1 \right) \right];$$

$$\bar{p}_k^i = \emptyset, \text{ если } (\forall n \forall m \zeta_1(\bar{g}_k^i, \bar{g}_n^i, \bar{g}_m^i) = 0) \&$$

$$\& (\forall n' \forall m' \zeta_2(\bar{g}_k^i, \bar{g}_n^j, d_m^{i,j}) = 0) \vee \exists r (\bar{\Sigma}_k^i = \bar{S}_r^i).$$

Утверждение 1. Для любого $\bar{p}_k^i = \{\bar{v}_k^i, \bar{\Sigma}_k^i\}$ имеет место $T(\bar{v}_k^i, \kappa, \bar{w}_i) = 1$, где

$$\kappa = \{ \{\bar{h}_n^i \mid \exists m \zeta_1(\bar{g}_k^i, \bar{g}_n^i, \bar{g}_m^i) = 1 \} \cup$$

$$\cup \{ \bar{h}_n^i \mid \exists m' \zeta_2(\bar{g}_k^i, \bar{g}_n^j, d_m^{i,j}) = 1 \} \}.$$

Доказательство.

1. Из определения сходства \bar{G}^i следует, что для каждого i, k множество \bar{S}_k^i представляет из себя как раз множество примеров, описываемых в части Ex предиката $M_{ag,n}^+$. Условие исчерпываемости (CE)

обеспечивается глобальностью сходства \bar{g}_k^i . Следовательно, $\bar{v}_k^i = \bar{h}_k^i$ удовлетворяет условиям $M_{ag,n}^+$, так как является максимальным операциональным сходством положительных примеров.

2. Из определения \bar{g}_n^i следует, что \bar{h}_n^i является максимальным операциональным сходством отрицательных примеров.

Условие минимальности \bar{h}_n^i по вложению, накладываемое на тормоза, обеспечивается выполнением предиката ζ_1 .

3. Из определения $d_m^{i,j}$ следует, что \bar{h}_n^j является максимальным операциональным сходством положительных примеров, не обладающих свойством \bar{w}_i .

4. Из определения $\bar{\Sigma}^i$ следует, что

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

Выполнение условий (1—4) обеспечивает $T(\bar{v}_k^+, \kappa, \bar{w}_i) = 1$ при κ , описанном в формулировке утверждения.

Покажем теперь, что, если $T(v', \kappa', \bar{w}_i) = 1$, то $\exists k v' = \bar{h}_k^+$.

$$\exists n \forall m \exists n' \exists m' \kappa' = \{ \{ \bar{h}_n^+ \mid \zeta_1(\bar{g}_k^+, \bar{g}_n^+, \bar{g}_m^+) = 1 \} \cup \cup \{ \bar{h}_n^+ \mid \zeta_2(\bar{g}_k^+, \bar{g}_n^+, \bar{g}_m^+) = 1 \} \}.$$

По определению отношения T , v есть максимальное операциональное сходство некоторого подмножества положительных примеров, но все такие сходства входят в глобальное сходство \bar{G}^+ , следовательно, $\exists k v = \bar{h}_k^+$.

Рассуждая аналогично, получаем, что

$$\forall v \in \kappa' (\exists n (v_n = \bar{h}_n^+) \vee (\exists j \exists n' (v_n = \bar{h}_n^+))).$$

Выполнение предиката ζ_1 и ζ_2 следует из определения сходств \bar{G} и \bar{P} .

Собственно говоря, утверждение 1 непосредственно следует из определения отношения T , предиката $M_{ag,n}^+$ и сходств \bar{G} и \bar{P} , что и видно непосредственно из доказательства.

Таким образом, глобальное сходство \bar{P}^+ порождается всеми теми и только теми объектами, операциональное сходство которых выражается причиной \bar{v}^+ свойства \bar{w}^+ при отсутствии тормозов κ .

Глобальное сходство \bar{P}^+ , как видно из его определения, существенно зависит от свойства \bar{w}^+ , поэтому структура множества свойств в базе данных определяет во многом сложность алгоритмической процедуры нахождения \bar{P}^+ . Так очевидно, что для случая: $W = \{w\}$, $\bar{W} = W$, процедура эта значительно упрощается. В самом деле, прежде всего упрощается предикат сходства $M_{ag,n}^+$, так как из него убирается добавка Bl_2 .

Сходство $\bar{P}^+ = \{ \bar{p}_j^+ \}$ принимает для этого случая вид:

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

\bar{P}^+ образуется с помощью сходств $\bar{G}^+ = \{ \bar{g}_j^+ \}$, $\bar{G}^- = \{ \bar{g}_n^- \}$ и $\bar{G} = \{ \bar{g}_m \}$.

Если сходство P найдено, то найдены все тройки $\{v, \kappa, w\}$, удовлетворяющие предикату T .

2. АЛГОРИТМЫ ПОСТРОЕНИЯ ОБОБЩЕННЫХ ГИПОТЕЗ

В конструктивности определения сходства \bar{P}^+ заложен, по существу, алгоритм 1 его нахождения, заключающийся в последовательности следующих действий:

1. Нахождение сходств \bar{G} и \bar{G} , осуществляемое с помощью, так называемого алгоритма максимальных пересечений, реализованного в версии ДСМ-системы для простого метода.

2. Проверка предиката ζ_1 .

3. Конструирование сходства \bar{P}^+ .

4. Проверка условия $\forall r (\bar{S}_k^+ \neq \bar{S}_r^+)$.

Продemonстрируем работу алгоритма для случая одного свойства на примере:

Множество объектов: $\{x_1, \dots, x_{11}\}; \Omega^+ = \{x_1, x_2, x_3, x_4, x_5\}, \Omega^- = \{x_6, x_7, \dots, x_{11}\}; x_1 = vaa', x_2 = vbdd', x_3 = vbd'd''m, x_4 = vd'cm, x_5 = vd''c'm, x_6 = vae, x_7 = vae', x_8 = vbdj, x_9 = vbdj', x_{10} = vd'gg', x_{11} = vd'gg''.$

Находим сходство \bar{G} из положительных примеров:

$$\bar{g}_1^+ = \{v; x_1, x_2, x_3, x_4, x_5\};$$

$$\bar{g}_2^+ = \{vbd'; x_2, x_3\};$$

$$\bar{g}_3^+ = \{vd'; x_2, x_3, x_4\};$$

$$\bar{g}_4^+ = \{vm; x_3, x_4, x_5\};$$

$$\bar{g}_5^+ = \{vd''m; x_2, x_3\}.$$

Сходство \bar{G} находим из отрицательных примеров.

$$\bar{g}_1^- = \{v; x_6, x_7, x_8, x_9, x_{10}, x_{11}\};$$

$$\bar{g}_2^- = \{va; x_6, x_7\};$$

$$\bar{g}_3^- = \{vbd; x_2, x_3\};$$

$$\bar{g}_4^- = \{vd'g; x_{10}, x_{11}\}.$$

Сходство \bar{G}^+ представлено следующим образом:

$$\bar{g}_1^+ = \{v; x_1, \dots, x_{11}\};$$

$$\bar{g}_2^+ = \{va; x_1, x_6, x_7\};$$

$$\bar{g}_3^+ = \{vb; x_2, x_3, x_8, x_9\};$$

$$\bar{g}_4^+ = \{vbd'; x_2, x_3\};$$

$$\bar{g}_5^+ = \{vbd; x_2, x_6, x_9\};$$

$$\bar{g}_6^+ = \{vd'; x_2, x_3, x_4, x_{10}, x_{11}\};$$

$$\bar{g}_7^+ = \{vd''m; x_2, x_3\};$$

$$\bar{g}_8^+ = \{vd'g; x_{10}, x_{11}\}.$$

Множество троек $\{ \bar{g}_i^+, \bar{g}_j^-, \bar{g}_m^+ \}$, для которых предикат ζ истинен, обозначим через Th_{ζ} .

Оно в данном примере следующее:

$$Th_{\zeta} = \{ (\bar{g}_1^+, \bar{g}_2^-, \bar{g}_3^+), (\bar{g}_1^+, \bar{g}_3^-, \bar{g}_5^+), (\bar{g}_1^+, \bar{g}_4^-, \bar{g}_8^+), (\bar{g}_2^+, \bar{g}_6^-, \bar{g}_3^+) \}.$$

Теперь можно построить сходства \bar{P}^+ :

$$\bar{P}_1^+ = \{v; \bar{\Sigma}_1^+ \setminus \setminus \{ (x_1, x_2, x_3, x_4, x_5) \setminus \setminus \{ (x_3, x_4, x_7 \setminus x_8, x_7) \cup \{ (x_2, x_3, x_9 \setminus x_6, x_9) \cup \{ (x_{10}, x_{11} \setminus x_{10}, x_{11}) \} \} \} \} = \{x_2, x_4, x_5\}.$$

Но так как $\bar{\Sigma}_1 = \bar{S}_1$, то $\bar{P}_1 = \emptyset$.

$\bar{P}_2 = \emptyset$, так как $\forall j \neg (\bar{h}_2 \subset \bar{h}_j)$ $j=1, \dots, 4$.

$\bar{P}_3 = \{vd'; \bar{\Sigma}_3 \{x_1, x_2, x_3\} \setminus \{x_{10}, x_{11} \setminus x_{10}, x_{11}\}\} = \{x_1, x_2, x_3\}$.

$\forall i \bar{\Sigma}_4 \neq \bar{S}_i$ ($i=1, \dots, 5$), следовательно, $\bar{P}_4 \neq \emptyset$.

$\bar{P}_4 = \emptyset$ и $\bar{P}_5 = \emptyset$ по той же причине, что и $\bar{P}_2 = \emptyset$.

В результате работы алгоритма мы найдем одну тройку: $\{vd', \kappa = vd'g, w\}$, удовлетворяющую предикату $M_{\alpha g, n}^+$ и, следовательно, принадлежащую отношению T .

Увеличение множества элементарных свойств хотя бы до двух элементов влечет за собой существенное увеличение разнообразия ситуаций, и следовательно, усложнение алгоритма.

Подсчитаем оценку сверху количества машинных операций, необходимых для порождения обобщенных гипотез с помощью указанной процедуры, для произвольных множества положительных примеров Ω^+ , множества отрицательных примеров Ω^- и множества элементарных структур U (подмножествами которого являются положительные и отрицательные примеры). В работе [5] предложен алгоритм ЗО поиска пересечений (сходств), временная трудоемкость которого линейна по отношению к числу порождаемых пересечений. Предположив, что сходства (пересечения) ищутся с помощью алгоритма ЗО, получим, что для нахождения

множества \bar{G} всех положительных сходств необходимо затратить

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

машинных операций; для нахождения множества \bar{G} всех отрицательных сходств необходимо затратить

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

машинных операций; для нахождения множества сходств \bar{G} всех примеров (положительных и отрицательных) необходимо затратить

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

машинных операций. Для порождения каждого множества \bar{S}_j в худшем случае необходимо

$$O(|\bar{G}| \cdot |\bar{G}| \cdot \text{pol}_1(|U|, |\Omega^+|, |\Omega^-|))$$

машинных операций. Всего по всем j (т. е. для нахождения всех положительных обобщенных гипотез) необходимо

$$O(|\bar{G}| \cdot |\bar{G}| \cdot \text{pol}_1(|U|, |\Omega^+|, |\Omega^-|)) = \\ = (|\bar{G}|^2 \cdot |\bar{G}|^2 \cdot \text{pol}_1(|U|, |\Omega^+|, |\Omega^-|))$$

машинных операций, где $\text{pol}_1(\cdot)$ — некоторый полином небольшой степени, зависящий от $|U|$, $|\Omega^+|$, $|\Omega^-|$ и определяемый вычислительной моделью.

Можно предложить более эффективные алгоритмы построения обобщенных гипотез. Быстродействие этих алгоритмов существенным образом зависит от типа операции сходства (даже если считать время выполнения таких операций одинаковым). Вначале дадим описание такого алгоритма для построения обобщенных поло-

жительных гипотез в случае представления данных множествами и наличия лишь одного свойства.

Алгоритм 2

1. Строятся все (+)-пересечения с помощью алгоритма ЗО [5].

2. Для каждого (+)-пересечения V :

2.1. Находятся отрицательные примеры его содержащие.

2.2. Ищутся все минимальные по вложению пересечения отрицательных примеров, содержащие V (т. е. тормоза): Пусть N есть объединение всех отрицательных примеров, содержащих V . Тогда произвольное минимальное пересечение отрицательных примеров, содержащее V , имеет вид $X^i = \bigcap_j X_{ij}$,

где $X_{ij} \in X_i = \{X \mid X \in \Omega^-, (V \cup \{q_i\}) \subset X\}$, где $q_i \in N$. Естественно, что $\bigcup_j X_{ij}$ определено, если $j > 2$.

Чтобы исключить порождение одинаковых пересечений и, тем самым, ускорить процесс их порождения, при построении очередного множества $\bigcup_j X_{ij}$, соот-

ветствующую «добавку» к v можно выбирать не из N , а из N за вычетом всех предыдущих «добавок» и порожденных пересечений, или формально, $N_1 = N$, $N_{k+1} = N_k \setminus (\{q_{k+1}\} \cup \bigcap_j X_{kj})$.

2.3. Из множества (+)-примеров, содержащих V , удаляются все такие, которые содержат какие-либо тормоза, оставшиеся (+)-примеры вновь пересекаются. Если результат совпадает с V , то V вместе с найденными минимальными пересечениями составляет обобщенную гипотезу, если нет, то (+)-обобщенной гипотезы относительно пересечения V не существует, переходим к следующему (+)-пересечению.

Временная сложность построения одной обобщенной положительной гипотезы для фиксированного (+)-пересечения V на шагах 2.1—2.3 такого алгоритма есть

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

Таким образом, сложность построения всех положительных обобщенных гипотез есть

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

Рассмотрим действие алгоритма для множеств положительных и отрицательных примеров, приведенных выше. Множество N в этом случае есть $N = \{a, b, d, d', e, f, f', g, g', g'', v\}$. Пусть на втором этапе действия алгоритма нами выбрано (+)-пересечение v . Тогда минимальные пересечения (—)-примеров, содержащие v , находятся как

$$X^a = \bigcap_j X_j^a, \text{ где } X_j^a \in X_a = \{X \in \Omega^- \mid \{v, a\} \subset X\} = \{v, a\},$$

$$N_2 = N \setminus \{a\} = \{b, d, d', e, e'f, f', g, g', g''\};$$

$$X^b = \bigcap_j X_j^b, \text{ где } X_j^b \in X_b = \{X \in \Omega^- \mid \{v, b\} \subset X\} = \{v, b, d\},$$

$$N_3 = N_2 \setminus \{b, d\} = \{d', e, e'f, f', g, g', g''\}.$$

Пересечение, соответствующее $\{v, d\}$, не порождается, так как на этом шаге $d \notin N$.

$$X^{d'} = \bigcap_j X_j^{d'}, \text{ где } X_j^{d'} \in X_{d'} = \{X \in \Omega^- \mid \{v, d'\} \subset X\} =$$

$$= \{v, d', g\}, N_4 = N_3 \setminus \{d', g\} = \{e, e'f, f', g', g''\};$$

$$X^e = \bigcap_j X_j^e, \text{ где } X_j^e \in X_e = \{X \in \Omega^- \mid \{v, e\} \subset X\},$$

$$X^e \text{ не определено, } N_5 = N_4 \setminus \{e\} = \{e'f, f', g', g''\};$$

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

$$X^{e'} \text{ не определено, } N_6 = N_5 \setminus \{e'\} = \{f, f', g', g''\}$$

$X^f = \bigcap_j X_j^f$, где $X_j^f \in X_f = \{X \in \Omega \mid \{v, f\} \subset X\}$,

X^f не определено, $N_7 = N_6 \setminus \{f\} = \{f', g', g''\}$;

$X^{f'} = \bigcap_j X_j^{f'}$, где $X_j^{f'} \in X_{f'} = \{X \in \Omega \mid \{v, f'\} \subset X\}$,

$X^{f'}$ не определено, $N_8 = N_7 \setminus \{f'\} = \{g', g''\}$.

Минимальное пересечение, соответствующее $\{v, g\}$, не порождается, так как $g \notin N_8$,

$X^{g'} = \bigcap_j X_j^{g'}$, где $X_j^{g'} \in X_{g'} = \{X \in \Omega \mid \{v, g'\} \subset X\}$,

$X^{g'}$ не определено, $N_9 = N_8 \setminus \{g'\} = \{g''\}$;

$X^{g''} = \bigcap_j X_j^{g''}$, где $X_j^{g''} \in X_{g''} = \{X \in \Omega \mid \{v, g''\} \subset X\}$,

$X^{g''}$ не определено, $N_{10} = N_9 \setminus \{g''\} = \emptyset$.

Осталось произвести действия, относящиеся к шагу 2.3. Из множества (+)-примеров, содержащих пересечение v , т. е. из $\{x_1, x_2, x_3, x_4, x_5\}$, удалим все примеры, содержащие минимальные пересечения (-)-примеров, содержащие v , т. е. $va, vbd, vd'gg'$. Такими (+)-примерами будут $x_1 = vaa'$, $x_2 = vbd'd'$. Пересечем оставшиеся (+)-примеры, т. е. x_3, x_4, x_5 , и получим $vm \neq v$. Значит, обобщенной (+)-гипотезы, соответствующей сходству (+)-примеров v , не будет.

Продолжаем вычисления шагов 2.1—2.3 для (+)-пересечения vd' . (-)-примерами, его содержащими, будут $x_{10} = vd'gg'$ и $x_{11} = vd'gg''$. Значит единственное минимальное пересечение (-)-примеров, содержащее vd' , будет $vd'g'$. Так как ни один (+)-пример его не содержит, то на шаге 2.3 ни один (+)-пример не будет отброшен, и тройка $\{vd', k = vd'g, w\}$ будет признана обобщенной (+)-гипотезой.

Как уже говорилось выше, алгоритм 2 предназначен для данных, представленных множествами. Последнее обстоятельство позволяет весьма эффективно осуществлять шаг 2.2. В случае, когда данные представлены объектами из произвольной полурешетки с операцией сходства Π , дело обстоит хуже: мы не можем образовать минимальное по вложению пересечение (-)-примеров, содержащее данное пересечение (+)-примеров v путем «наращивания» v атомами решетки (-)-пересечений (о решетке пересечений см., например, [5]) — мы можем просто не знать этих атомов. Таким образом, для произвольной операции сходства, вместо шага 2.2 алгоритма 2 мы должны искать минимальные пересечения «сверху — вниз»: беря пересечения нара-

стающего количества (-)-примеров, содержащих v . Как только на очередном шаге пересечение (-)-примеров становится равным v , мы отходим к предыдущему пересечению (-)-примеров: оно и будет минимальным среди содержащих v . При этом верхняя оценка числа машинных операций, необходимых для порождения всех минимальных пересечений, содержащих v , будет $O(|\bar{G}| \text{pol}_2(|U|, |\Omega^+|, |\Omega^-|))$, следовательно, общая сложность алгоритма будет $O(|\bar{G}| \cdot |\bar{G}| \text{pol}_2(|U\Omega|, |\Omega^+|, |\Omega^-|))$, где $\text{pol}_2(\cdot)$ — некоторый полином небольшой степени, зависящий от $|U|$, $|\Omega^+|$, $|\Omega^-|$ и определяемый вычислительной моделью.

СПИСОК ЛИТЕРАТУРЫ

1. Финн В. К. Правдоподобные выводы и правдоподобные рассуждения // Итоги науки и техники. Сер. Теория вероятностей. Математическая статистика. Теоретическая кибернетика. Т. 28.— М.: ВИНТИ, 1988.— С. 3—84.
2. Кузнецов С. О. ДСМ-метод как система автоматического обучения // Итоги науки и техники. Сер. Информатика. Т. 15.— М.: ВИНТИ, 1991.— С. 17—54.
3. Финн В. К. Об обобщенном ДСМ-методе автоматического порождения гипотез // Семантика и информатика.— 1989.— Вып. 29.— С. 93—123.
4. Гусакова С. М., Финн В. К. Сходство и правдоподобный вывод // Изв. АН СССР. Сер. Техн. кибернетика.— 1987.— Т. 5.— С. 42—63.
5. Кузнецов С. О. Быстрый алгоритм построения всех пересечений объектов из конечной полурешетки // НТИ. Сер. 2.— 1993.— № 1.— С. 17—20.
6. Гусакова С. М. Канонические представления сходств // НТИ. Сер. 2.— 1987.— № 9.— С. 19—21.
7. Гусакова С. М., Финн В. К. О новых средствах формализации локального и глобального сходств // НТИ. Сер. 2.— 1987.— № 10.— С. 14—22.
8. Гусакова С. М. Формализация понятия сходства и его применение в интеллектуальных информационных системах: Автореф. дис... к. ф.-м. н.— М.: ВИНТИ, 1988.— 21 с.
9. Финн В. К., Михеев М. А. Некоторые проблемы обобщенного ДСМ-метода автоматического порождения гипотез // Семантика и информатика.— 1983.— Вып. 33.— С. 136—163.

Материал поступил в редакцию 05.05.95.