

# МЕТОД ПРЕДСКАЗАНИЯ В ЯЗЫКЕ ПЕРВОГО ПОРЯДКА

Демин<sup>1</sup> А.В., Витяев<sup>2</sup> Е.Е.

<sup>1</sup>Институт систем информатики имени А. П. Ершова СО РАН г. Новосибирск

<sup>2</sup>Институт математики СО РАН г. Новосибирск, e-mail: [vityaev@math.nsc.ru](mailto:vityaev@math.nsc.ru)

## Аннотация

В работе продолжается рассмотрение метода и программной системы «Discovery» обнаружений знаний в данных, реализующие разработанный ранее реляционный подход к обнаружению знаний. Рассматривается метод предсказания, использующий обнаруженные системой «Discovery» закономерности в языке первого порядка с вероятностными оценками. Предлагаемый метод предсказания нетривиален и аналогов не имеет.

## §1. Определение вида гипотез

Данная работа является продолжением работ [1, 2], в которых описывается метод обнаружения закономерностей и его применения. Предлагаемый метод предсказания, опирается на это описание метода и соответствующую ему терминологию, а также использует идеи метода предсказания, изложенного в [3].

Напомним основные понятия многосортной логики первого порядка [5], которые нам понадобятся впоследствии. *Сигнатурой* называется упорядоченная шестерка  $\Sigma = \langle S, P, F, \eta, \mu, \kappa \rangle$ , где  $S$  – множество сортов,  $P$  – множество предикатных символов,  $F$  – множество функциональных символов,  $\eta: P \cup F \rightarrow \omega$  – отображение местности (арности) символов,  $\mu: \bigcup_{a \in P \cup F} \{a\} \times \eta(a) \rightarrow S$  – отображение, сопоставляющее каждому аргументу символа его сорт,  $\kappa: F \rightarrow S$  – отображение, сопоставляющее каждому функциональному символу сорт его значения.

Упорядоченная пара  $A = \langle A, \pi \rangle$  называется *многосортной алгебраической системой* сигнатуры  $\Sigma = \langle S, P, F, \eta, \mu, \kappa \rangle$ , где

- $A = \{A_s \mid s \in S\}$  – непустое множество, называемое носителем или основным множеством алгебраической системы  $A$ ;
- $\pi$  – отображение множества  $P \cup F$  в множество отношений и операций на множестве  $A$ , называемое интерпретацией сигнатуры  $\Sigma$  на  $A$ ;

- если  $P \in \mathcal{P}$ , то  $\pi(P) \subseteq A_{\mu(P,1)} \times \dots \times A_{\mu(P,n)}$ , где  $n = \eta(P)$  – местность символа  $P$ , и  $\pi(P)$  называется *многочесным отношением* на  $A$ ;
- если  $f \in \mathcal{F}$ , то  $\pi(f)$  – отображение  $\pi(f): A_{\mu(f,1)} \times \dots \times A_{\mu(f,n)} \rightarrow A_{\kappa(f)}$ , где  $n = \eta(f)$  – местность символа  $f$ , и  $\pi(f)$  – называется *многочесной операцией* на  $A$ .

Для определения понятия *терм* будем считать, что с каждым сортом  $s$  связано счетное множество символов  $V_s$  – переменных сорта  $s$ . Каждая переменная имеет только один сорт.

*Термом* сорта  $s$  называется любая переменная или константа сорта  $s$ , а также любое конечное выражение вида  $f(t_1, \dots, t_n)$ , где  $f \in \mathcal{F}$ ,  $\eta(f) = n$ ,  $\kappa(f) = s$ ,  $t_i$  – терм сорта  $\mu(f, i)$  для каждого  $i = 1, \dots, n$ .

*Атомарной формулой* называется выражение вида  $P(t_1, \dots, t_n)$ , где  $P \in \mathcal{P}$ ,  $\eta(P) = n$ ,  $t_i$  – терм сорта  $\mu(P, i)$  для каждого  $i = 1, \dots, n$ .

Значения многосортных термов и формул первого порядка практически не отличаются от односортного случая, с единственным ограничением, что переменные сорта  $s$  должны интерпретироваться элементами в множестве  $A_s$ .

Будем предполагать, что исходные данные представлены в виде таблицы значений  $D$ , строки которой соответствуют объектам, а колонки – признакам объектов. Т.е.  $D = \{D(1), \dots, D(N)\}$ , где  $D(i)$  –  $i$ -я строка таблицы (объект с номером  $i$ ),  $D(i) = \{D(i, 1), \dots, D(i, m)\}$ ,  $D(i, j)$  – значение таблицы на пересечении  $j$ -ой колонки и  $i$ -ой строки (значение  $j$ -го признака объекта  $D(i)$ ),  $D(i, j) \in \text{Re}$ , где  $\text{Re}$  – множество действительных чисел.

Зафиксируем фрагмент многосортного языка логики первого порядка сигнатуры  $\Sigma^* = \langle \mathcal{S}, \mathcal{P}, \mathcal{F}, \eta, \mu, \kappa \rangle$ , где

- $\mathcal{S} = \{s_{obj}, s_{Re}\}$ ,  $s_{obj}$  – сорт объектов таблицы  $D$ ,  $s_{Re}$  – сорт действительных чисел;
- $\mathcal{P}$  – множество предикатных символов, таких, что если  $P \in \mathcal{P}$ , то  $\mu(P, i) = s_{Re}$ , где  $i = 1, \dots, \eta(P)$ , т.е. все аргументы символов из  $\mathcal{P}$  имеют сорт  $s_{Re}$ ;
- $\mathcal{F} = \mathcal{F}_1 \cup \mathcal{F}_2 \cup \mathcal{F}_3$  – множество функциональных символов, где  $\mathcal{F}_1$  – множество функциональных символов  $f \in \mathcal{F}_1$ ,  $\mu(f, i) \in \{s_{obj}, s_{Re}\}$ ,  $i = 1, \dots, \eta(f)$ ,  $\kappa(f) = s_{obj}$ , т.е. аргументы символов из  $\mathcal{F}_1$  имеют сорт  $s_{obj}$  или  $s_{Re}$ , а их значения – сорт  $s_{obj}$ ;

$\mathbb{F}_2$  – множество функциональных символов  $f \in \mathbb{F}_2$ ,  $\eta(f) = 1$ ,  $\mu(f, I) = s_{obj}$  и  $\kappa(f) = s_{Re}$ , т.е. все функциональные символы из  $\mathbb{F}_2$  одноместные, их аргументы сорта  $s_{obj}$ , а значения – сорта  $s_{Re}$ .

$\mathbb{F}_3$  – множество функциональных символов  $f \in \mathbb{F}_3$ ,  $\mu(f, i) = s_{Re}$ , где  $i = 1, \dots, \eta(f)$ , и  $\kappa(f) = s_{Re}$ , т.е. все аргументы символов из  $\mathbb{F}_3$  и их значения имеют сорт  $s_{Re}$ ;

Введем многосортную алгебраическую систему  $\mathcal{A}^* = \langle \{D, Re\}, \pi \rangle$  сигнатуры  $\Sigma^*$ , где  $D$  – таблица данных, являющаяся носителем сорта  $s_{obj}$ , соответственно строки таблицы являются объектами сорта  $s_{obj}$ ,  $Re$  – множество действительных чисел, являющееся носителем сорта  $s_{Re}$ .

Поясним назначение множеств функциональных символов  $\mathbb{F}_1$ ,  $\mathbb{F}_2$  и  $\mathbb{F}_3$ .

Множество  $\mathbb{F}_1$  позволяет описывать функции, организующие доступ к различным строкам таблицы. Это могут быть, к примеру, функции, которые возвращают строки, смещенные на определенное число позиций в таблице относительно заданных строк, или функции, возвращающие строки по их номеру, или функции, организующие поиск строки в таблице.

Множество  $\mathbb{F}_2$  позволяет описывать функции, организующие доступ к различным признакам объектов. Это могут быть, к примеру, функции, возвращающие значение определенного признака заданной строки, или функции, возвращающие номер строки.

Множество  $\mathbb{F}_3$  позволяет описывать функции, осуществляющие различные преобразования над признаками объектов или значениями констант.

Для дальнейшего описания будем считать, что у нас есть счетное множество символов  $V_{obj}$  – переменных сорта  $s_{obj}$  и  $V_{Re}$  – переменных сорта  $s_{Re}$ .

Введем понятия *шаблона термов* и *шаблона предикатов*.

*Шаблон термов* – это пара  $Tf = \langle f, \{\Psi_1, \dots, \Psi_n\} \rangle$ , где  $f \in \mathbb{F}_3$ ,  $n = \eta(f)$ ,  $\Psi_i = \{t_1, \dots, t_{m_i}\}$ ,  $i = 1, \dots, n$ ,  $t_k$ ,  $k = 1, \dots, m_i$  – терм сорта  $s_{Re}$ , такой, что любая переменная, входящая в терм  $t_k$  имеет сорт  $s_{obj}$ . Шаблон терма  $Tf = \langle f, \{\Psi_1, \dots, \Psi_n\} \rangle$  определяет множество термов, состоящее из всех термов вида  $f(t_1, \dots, t_n)$ , где  $t_i \in \Psi_i$ ,  $i = 1, \dots, n$ . Будем обозначать через  $[Tf]$  – множество термов, определяемое шаблоном  $Tf$ . Мы можем оп-

ределить множество термов  $\Psi_i$  при помощи других шаблонов термов:

$\Psi_i = [Tf_1] \cup \dots \cup [Tf_k]$ , где  $Tf_j$ ,  $j = 1, \dots, k$  – некоторые шаблоны термов.

Приведем пример шаблона термов.

Пусть  $k$ -я колонка таблицы данных содержит значения некоторого временного ряда (к примеру, цена закрытия акции).

Пусть  $r$  – двухместный функциональный символ,  $r \in \mathbb{F}_1$ ,  $\mu(r, 1) = s_{obj}$  – первый аргумент  $r$  сорта  $s_{obj}$ ,  $\mu(r, 2) = s_{Re}$  – второй аргумент  $r$  сорта  $s_{Re}$ . Пусть интерпретацией символа  $r$  на  $\mathcal{A}^*$  является отображение  $\pi(r)(D(i), k) = D(i+k)$ , которое для заданной  $i$ -ой строки таблицы возвращает строку с номером  $(i+k)$  (предполагаем, что  $r$  определена только для целочисленных  $k$ ).

Пусть  $h_k$  – одноместный функциональный символ,  $h_k \in \mathbb{F}_2$ , интерпретацией которого на  $\mathcal{A}^*$  является отображение  $\pi(h_k)(D(i)) = D(i, k)$ , которое для заданного объекта ( $i$ -ой строки таблицы), возвращает значения его  $k$ -го признака (значение таблицы на пересечении  $k$ -ой колонки и  $i$ -ой строки).

Рассмотрим шаблон термов

$\langle t_1 - t_2, \{\Psi_1, \Psi_2\} \rangle$ , где  $\Psi_1 = \{h_k(r(i, 0))\}$ ,  $\Psi_2 = \{h_k(r(i, -1)), h_k(r(i, -2)), h_k(r(i, -3))\}$ .

Легко видеть, что данный шаблон определяет три терма, интерпретации которых на  $\mathcal{A}^*$  представляют собой функции, задающие временные лаги от 1 до 3:

- 1)  $D(i, k) - D(i-1, k)$ ,
- 2)  $D(i, k) - D(i-2, k)$ ,
- 3)  $D(i, k) - D(i-3, k)$ .

*Шаблон предикатов* – это пара  $Tr = \langle P, \{\Theta_1, \dots, \Theta_n\} \rangle$ , где  $P \in \mathbb{P}$ ,  $n = \nu(P)$ ,  $\Theta_i = \{Tf_1, \dots, Tf_{m_i}\}$ ,  $i = 1, \dots, n$ ,  $Tf_k$ ,  $k = 1, \dots, m_i$  – шаблоны термов. Шаблон предиката  $Tr = \langle P, \{\Theta_1, \dots, \Theta_n\} \rangle$  определяет множество атомарных формул, состоящее из всех атомарных формул вида  $P(t_1, \dots, t_n)$ , где  $t_i \in [Tf_1] \cup \dots \cup [Tf_{m_i}]$ ,  $i = 1, \dots, n$ ,  $Tf_k \in \Theta_i$ ,  $k = 1, \dots, m_i$ . Будем обозначать через  $[Tr]$  – множество атомарных формул, определяемых шаблоном  $Tr$ .

Приведем пример шаблона предикатов.

Предположим, что  $k$ -я колонка таблицы данных представляет временной ряд со значениями дневной цены закрытия какой-нибудь ценной бумаги. Предположим, что мы

хотим определить множество предикатов, сравнивающих друг с другом цены закрытия последних пяти дней. Это можно сделать, определив следующий шаблон предикатов:

$$Tp = \langle t_1 < t_2, \{\Theta_1, \Theta_2\} \rangle, \text{ где } \Theta_1 = Tf, \Theta_2 = Tf,$$

$$Tf = \langle x, \{\Psi\} \rangle, \Psi = \{h_k(r(i, 0)), h_k(r(i, -1)), h_k(r(i, -2)), h_k(r(i, -3)), h_k(r(i, -4))\},$$

функциональные символы  $r$  и  $h_k$ , и их интерпретации на  $\mathcal{A}^*$  определены в предыдущем примере.

В данном примере шаблон термов  $Tf$  задает пять термов, интерпретациями которых на  $\mathcal{A}^*$  являются следующие функции:  $x_1(i) = D(i, k)$ ,  $x_2(i) = D(i - 1, k)$ , ...,  $x_5(i) = D(i - 4, k)$ . Таким образом, шаблон предикатов  $Tp$  задает множество атомарных формул, интерпретации которых на  $\mathcal{A}^*$  представляют собой предикаты вида  $D(i - b_1, k) < D(i - b_2, k)$ , где  $b_1, b_2 \in \{0, -1, \dots, -4\}$ .

Теперь, используя понятие шаблона предикатов, мы можем определить понятие *класса гипотез*.

*Класс гипотез* – это пара  $Th = \langle \{Tp_1, \dots, Tp_m\}, P_0^\varepsilon \rangle$ , где  $Tp_i$  – шаблоны предикатов,  $P_0^\varepsilon$  – целевая литера,  $P_0$  – атомарная формула,  $\varepsilon \in \{0, 1\}$  – обозначает наличие отрицания формулы.

Класс гипотез  $Th = \langle \{Tp_1, \dots, Tp_m\}, P_0^\varepsilon \rangle$  определяет множество формул вида

$$\forall i_1, \dots, i_k (P_{i_1}^\varepsilon \& \dots \& P_{i_k}^\varepsilon \rightarrow P_0^\varepsilon), k \geq 0 \quad (1)$$

где  $i_1, \dots, i_k$  – переменные сорта  $s_{obj}$ ,  $P_i \in [Tp_1] \cup \dots \cup [Tp_m]$ ,  $i = 1, \dots, n$ . Формулы вида (1) будем называть *правилами*.

В дальнейшем помимо указанной записи мы также будем использовать более удобную форму записи формул (1) через индивидные константы. Вместо кванторов всеобщности и связанных ими переменных введем константы  $z_1, \dots, z_k$ . Тогда формулы вида (1) преобразуются в формулы вида

$$P_1^\varepsilon \& \dots \& P_n^\varepsilon \rightarrow P_0^\varepsilon. \quad (2)$$

Смысл формулы (2) состоит в том, что при любой фиксированной замене индивидных констант на объекты из таблицы  $D$ , из истинности посылки должна следовать истинность заключения.

Данное выше определение понятия класса гипотез позволяет разработать интерактивных способ задания классов гипотез, проверяемых на данных. Для этого достаточно разработать интерактивную систему конструирования термов, атомарных формул, шабло-

нов термов и шаблонов предикатов из заданного набора функциональных и предикатных символов. Подобная интерактивная система задания классов гипотез была реализована в программной системе «Discovery» [1, 2, 4].

Предложенный способ задания классов гипотез также указывает способ вычисления различных гипотез для заданного класса. Для этого достаточно реализовать вычислительные процедуры для базового набора функциональных и предикатных символов. Вычисление различных формул вида (1) может быть сведено к вычислению соответствующих базовых символов, из которых построены формулы. Данный способ вычисления гипотез используется в реализованной версии системы «Discovery».

## §2. Общая формулировка метода предсказания

Пусть  $Reg(Th)$  – множество закономерностей, полученное методом обнаружения закономерностей [1, 2], по заданному классу гипотез  $Th = \langle \{Tp_1, \dots, Tp_m\}, P_0^e \rangle$  на обучающем множестве,  $\mathbb{D} \subset \mathbb{A}$  случайно выбранном из генеральной совокупности объектов  $\mathbb{A}$ .

В данном параграфе приводится общая формулировка метода предсказания, использующего множество обнаруженных закономерностей  $Reg(Th)$ .

Пусть  $P(Th) = \{P_1, \dots, P_n\}$  – множество всех атомарных формул, которые мы можем получить с помощью шаблонов предикатов [2]  $\{Tp_1, \dots, Tp_m\}$ , входящих в класс гипотез  $Th$  [2]. Пусть из генеральной совокупности объектов  $\mathbb{A}$  случайно выбран некоторый новый объект  $b$ . В задачах предсказания считается, что истинностные значения некоторой части атомарных формул  $P^H \subset P(Th)$  на объектах  $\mathbb{D}$  и  $b$  нам известны. Требуется, используя знания закономерностей из  $Reg(Th)$ , по известным значениям истинности атомарных формул  $P^H$  на объектах  $\mathbb{D} \cup b$  предсказать неизвестные значения истинности остальных атомарных формул  $P^H = P(Th) \setminus P^H$  на этих же объектах  $\mathbb{D} \cup b$ . Таким образом, задача предсказания состоит в том, чтобы по модели  $pr_0 = \langle \mathbb{D}, P(Th) \rangle$ , на которой проводилось обучение, и модели  $pr^H = \langle \mathbb{D} \cup b, P^H \rangle$  восстановить модель  $pr = \langle \mathbb{D} \cup b, P(Th) \rangle$ , используя закономерности из  $Reg(Th)$ .

Обозначим через  $PS$  множество всех возможных моделей  $pr = \langle \mathbb{D} \cup b, P(Th) \rangle$ , являющихся восстановлениями моделей  $pr_0, pr^H$ . Так как множество  $Reg(Th)$  содержит статистические закономерности, то разные восстановления могут иметь разную вероятность.

Таким образом, метод предсказания должен состоять в том, чтобы по закономерностям  $Reg(Th)$  и моделям  $pr_0, pr^n$  вычислить для каждой модели  $pr \in PS$  некоторую оценку её вероятности  $v(pr)$ . Для некоторых моделей  $pr \in PS$  оценка вероятности  $v(pr)$  может быть не определена, так как может, например, оказаться, что для неё нет применимых к ней закономерностей.

**Определение 1.** *Методом предсказания* будем называть алгоритм  $AP: \langle Reg(Th), pr_0, pr^n \rangle \rightarrow v$ , преобразующий тройки  $\langle Reg(Th), pr_0, pr^n \rangle$  в частично определенное отображение  $v: PS \rightarrow [0, 1]$ .

### §3. Метод предсказания

Уточним в чем состоит смысл распространения моделей  $pr_0, pr^n$  до моделей  $pr \in PS$ . Если известны все вероятности, то для любой модели  $pr \in PS$  можно подсчитать вероятность  $\wp(pr)$  того, что при случайном выборе объекта  $b$  из  $\mathbb{A}$  мы в результате эксперимента над  $\mathbb{D} \cup b$  получим модель изоморфную  $pr$ . Поэтому для получения наиболее точного предсказания алгоритм  $AP$  должен стремиться получить оценки вероятности  $v(pr)$  наиболее близкие к вероятности  $\wp(pr)$ .

Для восстановления модели  $pr$  надо определить значения истинности всех атомарных формул из  $P^n$ , на всех наборах объектов, включающих хотя бы одно вхождение объекта  $b$ . Для этой цели могут быть использованы те закономерности из  $Reg(Th)$  в заключение которых стоит атомарная формула из  $P^n$ . Разобьем это множество закономерностей  $Reg(Th)$  на три группы:

$Reg_1$  – множество закономерностей, включающее закономерности, содержащие только одноместные атомарные формулы, содержащие одну индивидуальную постоянную.

$Reg_2$  – множество закономерностей, заключение которых содержит только одну индивидуальную постоянную, а посылка содержит, по крайней мере две различные индивидуальные постоянные.

$Reg_3$  – множество закономерностей, у которых в заключении есть хотя бы две различные индивидуальные постоянные.

Для произвольного правила  $R \in Reg(Th)$ ,  $R = P_1^{\epsilon_1} \& \dots \& P_n^{\epsilon_n} \rightarrow P_0^{\epsilon_0}$  будем обозначать через  $D_n = P_1^{\epsilon_1} \& \dots \& P_n^{\epsilon_n}$  конъюнкцию литер посылки, а через  $D_c$  – заключение пра-

вила  $P_0^{e_0}$ . Будем также обозначать через  $z(D)$  множество индивидуальных констант, входящих в формулу  $D$ .

Для осуществления предсказания необходимо в первую очередь определить для каждой закономерности  $R \in \mathbf{Reg}$  множество моделей  $PS(R) \subset PS$ , которое будет являться прогнозом для данной закономерности.

Если  $R = (D_H \rightarrow D_C) \in \mathbf{Reg}_1(Th)$ , то проверим истинность формулы  $D_H$  при подстановке в нее объекта  $b$ . Если  $D_H$  истинна, то данная закономерность может быть использована для предсказания. Прогнозом закономерности  $R$  будет являться множество  $PS(R) = PS(D_C)$  тех моделей, на которых формула  $D_C$  истинна.

Закономерности множеств  $\mathbf{Reg}_2$ ,  $\mathbf{Reg}_3$  принципиально отличаются от закономерностей  $\mathbf{Reg}_1$  тем, что в них есть несколько индивидуальных постоянных. Поэтому, подставляя объект  $b$  вместо одной индивидуальной постоянной, мы должны подставить некоторые объекты и вместо других индивидуальных постоянных. Закономерность в этом случае говорит об определенной связи объекта  $b$  с другими объектами. Поэтому закономерности множеств  $\mathbf{Reg}_2$ ,  $\mathbf{Reg}_3$  могут быть различным способом использованы для предсказания.

Если закономерность  $R = (D_H \rightarrow D_C)$  принадлежит  $\mathbf{Reg}_2(Th)$ , то разобьем случайным образом обучающее множество  $\mathbb{D}$  на  $l$  наборов объектов по  $k$  объектов в каждом, где  $l = \frac{m}{k}$ ,  $m = \overline{\mathbb{D}}$ ,  $k$  – количество индивидуальных постоянных во множестве  $z(D_H) \setminus z(D_C)$ .

Будем последовательно подставлять эти наборы вместо индивидуальных постоянных и определять значения истинности формулы  $D_H$ . Если формула  $D_H$  истинна хотя бы на одном наборе объектов, то данную закономерность можно использовать для предсказания. В противном случае по этой закономерности предсказание сделать нельзя. Прогнозом данной закономерности, как и в предыдущем случае, будет множество  $PS(R) = PS(D_C)$ .

Закономерности из  $\mathbf{Reg}_3$  принципиально отличаются от закономерностей из  $\mathbf{Reg}_1$  и  $\mathbf{Reg}_2$  тем, что в них предсказывается не истинность некоторого отношения, зависящего от одной индивидуальной постоянной, а предсказывается определенное отношение между одной индивидуальной постоянной и некоторыми другими индивидуальными постоянными.

Пусть  $R = (D_H \rightarrow D_C) \in \mathbf{Reg}_3$ . Обозначим через  $\Pi = \{z_1, \dots, z_p\}$ ,  $\Pi \subset z(D_C)$  множество индивидуальных постоянных из  $D_C$ , которые входят в  $D_H$ , но не входят в атомарные формулы из  $P''$ . Предсказываемый объект  $b$  можно подставлять вместо любой индивидуальной постоянной из  $\Pi$ .



Подставив объект  $b$  вместо индивидуальной постоянной  $z \in \Pi$ , мы получим закономерность  $R_b^z = ((D_\Pi)_b^z \rightarrow (D_C)_b^z)$ . Пусть  $z_1, \dots, z_k$  – остальные индивидуальные постоянные, входящие в закономерность  $R_b^z$ . Разобьем случайным образом всё множество  $\mathbb{D}$  на  $l$  наборов объектов по  $k$  объектов в каждом, где  $l = \frac{m}{k}$ ,  $m = \overline{\mathbb{D}}$ . Пусть  $\{\langle a_1 \rangle, \dots, \langle a_l \rangle\}$  – полученное множество наборов  $\langle a_i \rangle = \{a_1^i, \dots, a_k^i\}$ ,  $a_j^i \in \mathbb{D}$ ,  $i = 1, \dots, l$ ,  $j = 1, \dots, k$ . Подставим эти наборы вместо соответствующих индивидуальных постоянных и определим значения истинности формул  $(D_\Pi)_b^z$ ,  $(D_C)_b^z$ . Если формула  $(D_\Pi)_b^z$  истинна хотя бы на одном наборе объектов, то закономерность  $R_b^z$  можно использовать для предсказания. Множество моделей, являющихся прогнозом закономерности  $R_b^z$  определим следующим образом:

$PS(R_b^z) = \bigcup_{i=1}^l PS((D_C)_b^z(\langle a_i \rangle))$ , где  $PS((D_C)_b^z(\langle a_i \rangle))$  – множество моделей из  $PS$ , в которых формула  $(D_C)_b^z$  истинна на наборе объектов  $\langle a \rangle$ .

Подставляя последовательно объект  $b$  вместо каждой индивидуальной константы из  $\Pi$  мы получим закономерности  $R_b^{z_1}, \dots, R_b^{z_p}$ . Если хотя бы одна из этих закономерностей будет применима для предсказания, то закономерность  $R$  можно использовать для предсказания. В качестве прогноза закономерности  $R$  определим множество моделей

$$PS(R) = \bigcup_{i=1}^p PS(R_b^{z_i}).$$

Чтобы оценить точность прогноза  $\nu(pr)$  некоторой модели  $pr \in PS$ , метод предсказания должен основываться на прогнозах всех закономерностей из  $Reg(Th)$ , способных предсказать модель  $pr$ .

Обозначим через  $Reg(pr)$  множество закономерностей из  $Reg(Th)$  таких, что для любого  $R \in Reg(pr)$ ,  $pr \in PS(R)$ .

Для получения итоговой оценки точности прогноза  $\nu(pr)$  для модели  $pr$  необходимо задать частично определенную функцию  $\lambda : \{\langle pr, Reg(pr) \rangle\} \rightarrow [0, 1]$ , которая для каждой модели  $pr \in PS$  на основании множества закономерностей  $Reg(pr)$ , предсказывающих данную модель, вычисляет оценку точности прогноза  $\nu(pr)$ . Тогда  $\nu(pr) = \lambda(\langle pr, Reg(pr) \rangle)$ .

В зависимости от специфики решаемой задачи функция  $\lambda$  может быть определена по-разному. Приведем несколько примеров задания функции  $\lambda$ . Обозначим через

$Reg^+(pr) \subset Reg(pr)$  множество закономерностей применимых для осуществления прогноза.

1. По максимальной вероятности:

$$\lambda(\langle pr, Reg(pr) \rangle) = \max_{R \in Reg^+(pr)} \{\wp(R)\},$$

где  $\wp(R)$  – условная вероятность правила  $R$ .

2. По относительному количеству сработавших правил:

$$\lambda(\langle pr, Reg(pr) \rangle) = \frac{n(Reg^+(pr))}{n(Reg(pr))}, \text{ если } n(Reg(pr)) \neq 0,$$

$$\lambda(\langle pr, Reg(pr) \rangle) = 0, \text{ если } n(Reg(pr)) = 0,$$

где  $n(Reg(pr))$  и  $n(Reg^+(pr))$  – количество правил во множествах  $Reg(pr)$  и  $Reg^+(pr)$ .

3. По средней вероятности:

$$\lambda(\langle pr, Reg(pr) \rangle) = \frac{\sum_{R \in Reg^+(pr)} \wp(R)}{n(Reg^+(pr))}, \text{ если } n(Reg^+(pr)) \neq 0,$$

$$\lambda(\langle pr, Reg(pr) \rangle) = 0, \text{ если } n(Reg^+(pr)) = 0.$$

4. По средневзвешенной вероятности:

$$\lambda(\langle pr, Reg(pr) \rangle) = \frac{\sum_{R \in Reg^+(pr)} \wp(R)}{\sum_{R \in Reg(pr)} \wp(R)}, \text{ если } \sum_{R \in Reg(pr)} \wp(R) \neq 0,$$

$$\lambda(\langle pr, Reg(pr) \rangle) = 0, \text{ если } \sum_{R \in Reg(pr)} \wp(R) = 0.$$

#### §4. Метод предсказания, основанный на оценке максимальной вероятности [2].

**4.1 Вероятностные оценки закономерностей.** Сначала для всех трех видов закономерностей  $Reg_1, Reg_2$  и  $Reg_3$  на обучающем материале подсчитаем некоторые вероятностные оценки, необходимые для получения предсказаний.

Для закономерностей из  $Reg_1$  подсчитаем нижнюю доверительную границу для условной вероятности  $\wp(D_C | D_H)$ . При подстановке вместо единственной индивидуальной постоянной объектов из  $\mathbb{D}$ , формулы  $D_H, D_C$  будут принимать определенные значения истинности. Подсчитаем на объектах  $\mathbb{D}$  частоту  $h(D_C | D_H)$  условного события  $D_C | D_H$ . Используем один из известных методов построения доверительных интервалов для условной вероятности. Для фиксированного доверительного уровня  $\beta$  по частоте можно опреде-

литель нижнюю доверительную границу  $\underline{h}^\beta(D_C | D_\Pi)$  для условной вероятности, обладающую свойством

$$\wp(\wp(D_C | D_\Pi) \geq \underline{h}^\beta(D_C | D_\Pi)) \geq 1 - \beta. \quad (3)$$

Для закономерностей из  $Reg_2$ , также как и для закономерностей из  $Reg_1$ , вычислим нижнюю доверительную границу для условной вероятности  $p(D_C | D_\Pi)$ . Для этого, как и в предыдущем случае, достаточно вычислить частоту  $h(D_C | D_\Pi)$ . Для двух и более индивидуальных постоянных процедура вычисления частоты отличается от предыдущего случая. Разобьем случайным образом множество объектов  $\mathbb{D}$  на два множества, которые обозначим соответственно через  $A$  и  $C$ ,  $\mathbb{D} = A \cup C$ ,  $A \cap C = \emptyset$ . Объекты, имитирующие объект  $b$ , будем брать из множества  $C$ , а остальные объекты, подставляемые в закономерность, будем брать из множества  $A$ .

Пусть  $z(D_C) = z$  и  $z(D_\Pi) = \{z, z_1, \dots, z_k\}$ . Подсчитает частоту  $h(D_C | D_\Pi)$ . Будем подставлять последовательно все объекты из  $C$  вместо индивидуальной постоянной  $z$ . Для каждого подставленного объекта случайным образом выберем набор объектов  $a_1, \dots, a_k$  из множества  $A$ . Подставим его вместо индивидуальных постоянных  $z_1, \dots, z_k$ . Формулы  $D_C$  и  $D_\Pi$  примут определенные значения истинности. Подставив последовательно все объекты из  $C$  вместо индивидуальной постоянной  $z$ , можно подсчитать частоту  $h(D_C | D_\Pi)$ . Интерпретация условного события  $D_C | D_\Pi$  определяется как вероятность того, что формула  $D_C$  будет истинна при подстановке в нее объекта  $b$  вместо индивидуальной постоянной  $z$  и при подстановке случайно выбранных из множества  $A$  объектов вместо остальных индивидуальных постоянных. При заданном доверительном уровне  $\beta$  по частоте  $h(D_C | D_\Pi)$  можно вычислить нижнюю доверительную границу условной вероятности, удовлетворяющую неравенству (3).

Возьмем произвольную закономерность из  $Reg_3$ . Объекты, с которыми предсказываемый объект  $b$  должен находиться в некотором отношении, будем, также как и в предыдущем случае, брать из множества  $C$ . Пусть  $D_C = P^{\varepsilon_0}(z_1, \dots, z_{m_0})$ . Обозначим через  $\Pi \subset \{z_1, \dots, z_{m_0}\}$  множество индивидуальных постоянных из  $D_C$ , которые в  $D_\Pi$  не входят в атомарные формулы из  $P^\Pi$ . По определению множества  $Reg_3$ , множество  $\Pi$  не пусто. Предсказываемый объект  $b$  можно подставлять вместо любой индивидуальной постоянной из  $\Pi$ .

Для данной закономерности и индивидуальной постоянной  $z \in \Pi$  подсчитаем следующую оценку. Подставим объект  $b \in C$  вместо индивидуальной постоянной  $z$ . Пусть  $z_1, \dots, z_k$  – остальные индивидуальные постоянные, входящие в закономерность. Вместо этих индивидуальных постоянных будем подставлять объекты из множества  $A$ . Для этого разобьем случайным образом все множество  $A$  на  $l$  наборов объектов по  $k$  объектов в каждом, где  $l = \frac{m_1}{k}$ ,  $m_1$  – количество элементов во множестве  $A$ . Подставив любой из этих наборов вместо соответствующих индивидуальных постоянных, мы можем определить значения истинности формул  $D_\Pi$  и  $D_C$ . Подставляя последовательно все наборы в закономерность для данного объекта  $b$  можно подсчитать частоту  $h_z^b(D_C | D_\Pi)$ . При случайной подстановке наборов объектов множества  $A$  вместо индивидуальных постоянных  $z_1, \dots, z_k$ , частота  $h_z^b(D_C | D_\Pi)$  будет одномерной случайной величиной. Эта случайная величина имеет неизвестное нам распределение  $H_z^b$ . Будем подставлять вместо объекта  $b$  объекты из множества  $C$ . Получим  $m_2$  выборочных значений частоты  $h_z^{c_1}, \dots, h_z^{c_{m_2}}$ ,  $C = \{c_1, \dots, c_{m_2}\}$ . Найти доверительные границы для частоты по данным выборочным значениям при неизвестной функции распределения можно при помощи порядковых статистик. Так как частоты  $h_z^{c_i}$  равны отношению целых чисел, то случайная величина  $h_z^c$  дискретна.

Приведем необходимые результаты из порядковых статистик [6, с. 128-130]. Пусть  $h_{(1)} < h_{(2)} < \dots < h_{(m)}$  – порядковые статистики в выборке объема  $m$  из генеральной совокупности с неизвестной непрерывной функцией распределения  $H$ . Вероятность того, что новое значение случайной величины будет находиться в пределах  $h_{(r)} < h < h_{(s)}$ , равна  $H(h_{(s)}) - H(h_{(r)})$ . Эта вероятность как функция случайных величин сама будет случайной величиной. Используя порядковые статистики  $h_{(r)}$  и  $h_{(s)}$ , можно получить следующие толерантные интервалы для распределения  $H$

$$1 - \beta = \wp(H(h_{(s)}) - H(h_{(r)}) \geq \gamma) = \sum_{i=0}^{s-r-1} \binom{m}{i} \gamma^i (1 - \gamma)^{m-i}. \quad (4)$$

Из данного равенства следует, что вероятность того, что с вероятностью большей либо равной  $\gamma$  имеет место неравенство  $h_{(r)} < h < h_{(s)}$ , равна  $1 - \beta$ . Используя эту формулу, можно для значений вероятностей  $\gamma$  и  $\beta$ , и объема  $m$  выборки найти такие порядковые статистики  $h_{(r)}$  и  $h_{(s)}$ , чтобы выполнялось неравенство  $\wp(H(h_{(s)}) - H(h_{(r)}) \geq \gamma) \geq 1 - \beta$ .

Используя результаты Тьюки [7] можно распространить это неравенство и на дискретное распределение  $H$ . Если обозначить через  $H(h_{(s)} + \theta)$  и  $H(h_{(r)} - \theta)$  пределы функции распределения  $H$  соответственно справа и слева в точках  $h_{(s)}$  и  $h_{(r)}$ , то, согласно [7], должно выполняться неравенство

$$\wp(H(h_{(s)} + \theta) - H(h_{(r)} - \theta) \geq \gamma) \geq 1 - \beta. \quad (5)$$

В неравенстве (5) номера  $s$  и  $r$  можно вычислить, используя формулу (4). Расположим выборочные значения частоты  $h_z^{c_i}$ ,  $i = 1, \dots, m_2$  в порядке возрастания. Получим  $m_2$  порядковых статистик  $h_{(1)} < h_{(2)} < \dots < h_{(m_2)}$ . Фиксируем некоторые вероятности  $\gamma$  и  $\beta$ . Для этих значений вероятности и данного числа  $m_2$  порядковых статистик можно по формуле (4) вычислить номера  $s$  и  $r$  порядковых статистик  $h_{(s)}$  и  $h_{(r)}$ , для которых будет выполняться неравенство (5). Эти статистики определяют толерантный интервал  $[h_{(r)}, h_{(s)}]$ , который обозначим через  $[\underline{h}^z, \bar{h}^z]$ .

Подсчитаем для каждой закономерности из  $Reg_3$  и для каждой индивидуальной постоянной  $z \in \Pi$  толерантный интервал  $[\underline{h}^z, \bar{h}^z]$ . На этом подсчет статистических оценок для закономерностей из  $Reg_1$ ,  $Reg_2$ ,  $Reg_3$  закончен.

**4.2 Предсказание.** Рассмотрим сначала случай, когда множество закономерностей  $Reg(Th)$  принадлежит либо  $Reg_1(Th)$ , либо  $Reg_2(Th)$ .

Пусть формула  $D_{\Pi}$  истинна. Разобьем множество  $PS$  на два множества  $PS(D_C)$  и  $PS(\bar{D}_C)$ , включающие соответственно модели, на которых формула  $D_C$  истинна, и на которых она ложна. В соответствии с закономерностью, при истинности  $D_{\Pi}$  предсказываются те модели из  $PS$ , которые принадлежат множеству  $PS(D_C)$ . Это множество, рассматриваемое как событие при случайном выборе объекта  $b$ , имеет некоторую вероятность  $\wp(PS(D_C))$ , оценкой которой является величина  $\underline{h}^{\beta}(D_C | D_{\Pi})$ . Поэтому определим отображение  $\nu$  следующим образом:  $\nu(pr) = \theta$ , если  $pr \in PS(\bar{D}_C)$ , и  $\nu(pr) = \underline{h}^{\beta}(D_C | D_{\Pi})$ , если  $pr \in PS(D_C)$ . Для такого предсказания в силу неравенства (3) будет выполнено соотношение

$$\wp(\wp(PS(D_C)) \geq \underline{h}^{\beta}(D_C | D_{\Pi})) \geq 1 - \beta. \quad (6)$$

Рассмотрим случай, когда множество  $Reg_1(Th) \cup Reg_2(Th)$  состоит из закономерностей  $D_{\Pi_1} \rightarrow D_{C_1}, \dots, D_{\Pi_k} \rightarrow D_{C_k}$ , предсказывающих одну и ту же литеру

$P_0^{\varepsilon_0}(z_1, \dots, z_{m_0}) = D_{C_1} = \dots = D_{C_k}$ . Проведем процедуру предсказания данной литеры по всем закономерностям, как описано выше. Для тех закономерностей, у которых посылка  $D_{\Pi_i}$ ,  $i = 1, 2, \dots, k$  будет истинной, получим предсказания одного и того же множества  $PS(D_C)$  с оценками вероятности  $\underline{h}^{\beta}(D_{C_i} | D_{\Pi_i}), \dots, \underline{h}^{\beta}(D_{C_k} | D_{\Pi_k})$ . Способ получения результирующей оценки должен определяться дополнительными предположениями, которые можно сделать относительно связи этих закономерностей.

#### 4.3 Дополнительные предположения о независимости закономерностей.

**Определение 2.** Две закономерности  $D_{\Pi}^1 \rightarrow P_0^{\varepsilon}$  и  $D_{\Pi}^2 \rightarrow P_0^{\varepsilon}$  будем называть *независимыми*, если выполнены следующие условия

1.  $\wp(D_{\Pi}^1 \& D_{\Pi}^2) = \wp(D_{\Pi}^1)p(D_{\Pi}^2)$ ,
2.  $\wp((D_{\Pi}^1 \& \bar{P}_0^{\varepsilon}) \& (D_{\Pi}^2 \& \bar{P}_0^{\varepsilon})) = \wp((D_{\Pi}^1 \& \bar{P}_0^{\varepsilon}))\wp((D_{\Pi}^2 \& \bar{P}_0^{\varepsilon}))$ ,

где  $\bar{P}_0^{\varepsilon}$  – отрицание литеры  $P_0^{\varepsilon}$ .

**Лемма 1.** Если закономерности  $D_{\Pi}^1 \rightarrow P_0^{\varepsilon}$  и  $D_{\Pi}^2 \rightarrow P_0^{\varepsilon}$  независимы, то условная вероятность закономерности  $D_{\Pi}^1 \& D_{\Pi}^2 \rightarrow P_0^{\varepsilon}$  равна

$$\wp(P_0^{\varepsilon} | D_{\Pi}^1 \& D_{\Pi}^2) = 1 - (1 - \wp(P_0^{\varepsilon} | D_{\Pi}^1))(1 - \wp(P_0^{\varepsilon} | D_{\Pi}^2)). \quad (7)$$

**Доказательство.** Разделим условие 2 независимости на  $\wp(D_{\Pi}^1 \& D_{\Pi}^2)$ . Получим

$$\frac{\wp((D_{\Pi}^1 \& \bar{P}_0^{\varepsilon}) \& (D_{\Pi}^2 \& \bar{P}_0^{\varepsilon}))}{\wp(D_{\Pi}^1 \& D_{\Pi}^2)} = \frac{\wp(D_{\Pi}^1 \& \bar{P}_0^{\varepsilon})\wp(D_{\Pi}^2 \& \bar{P}_0^{\varepsilon})}{\wp(D_{\Pi}^1 \& D_{\Pi}^2)}. \quad (8)$$

Левая часть равенства равна  $\wp(\bar{P}_0^{\varepsilon} | D_{\Pi}^1 \& D_{\Pi}^2)$ . Заменяем вероятность  $\wp(D_{\Pi}^1 \& D_{\Pi}^2)$  в правой части равенства на  $\wp(D_{\Pi}^1)\wp(D_{\Pi}^2)$ , согласно условию 1 независимости. Тогда правая часть равенства (8) равна  $\wp(\bar{P}_0^{\varepsilon} | D_{\Pi}^1)\wp(\bar{P}_0^{\varepsilon} | D_{\Pi}^2)$ . Поэтому равенство (8) переходит в равенство

$$\wp(\bar{P}_0^{\varepsilon} | D_{\Pi}^1 \& D_{\Pi}^2) = \wp(\bar{P}_0^{\varepsilon} | D_{\Pi}^1)\wp(\bar{P}_0^{\varepsilon} | D_{\Pi}^2).$$

Так как

$$\begin{aligned} \wp(\bar{P}_0^{\varepsilon} | D_{\Pi}^1 \& D_{\Pi}^2) &= 1 - \wp(P_0^{\varepsilon} | D_{\Pi}^1 \& D_{\Pi}^2), \\ \wp(\bar{P}_0^{\varepsilon} | D_{\Pi}^1) &= 1 - \wp(P_0^{\varepsilon} | D_{\Pi}^1), \\ \wp(\bar{P}_0^{\varepsilon} | D_{\Pi}^2) &= 1 - \wp(P_0^{\varepsilon} | D_{\Pi}^2), \end{aligned}$$

то отсюда получаем утверждение леммы ■.

**Лемма 2.** Если закономерности  $D_{\Pi}^1 \rightarrow P_0^{\varepsilon}$  и  $D_{\Pi}^2 \rightarrow P_0^{\varepsilon}$  независимы, и их нижние доверительные границы условных вероятностей  $p(P_0^{\varepsilon} | D_{\Pi}^1)$  и  $p(P_0^{\varepsilon} | D_{\Pi}^2)$  равны соответственно  $\underline{h}^{\beta}(P_0^{\varepsilon} | D_{\Pi}^1)$ ,  $\underline{h}^{\beta}(P_0^{\varepsilon} | D_{\Pi}^2)$ , то нижняя доверительная граница для условной вероятности  $\wp(P_0^{\varepsilon} | D_{\Pi}^1 \& D_{\Pi}^2)$  равна  $\underline{h}^{2\beta}(P_0^{\varepsilon} | D_{\Pi}^1 \& D_{\Pi}^2) = 1 - (1 - \underline{h}^{\beta}(P_0^{\varepsilon} | D_{\Pi}^1))(1 - \underline{h}^{\beta}(P_0^{\varepsilon} | D_{\Pi}^2))$ .

**Доказательство.** Из условия и леммы 1 следует, что имеют место следующие отношения:

$$\wp(\wp(P_0^\varepsilon | D_{II}^1) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^1)) \geq 1 - \beta,$$

$$\wp(\wp(P_0^\varepsilon | D_{II}^2) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^2)) \geq 1 - \beta,$$

$$\wp(P_0^\varepsilon | D_{II}^1 \& D_{II}^2) = 1 - (1 - \wp(P_0^\varepsilon | D_{II}^1))(1 - \wp(P_0^\varepsilon | D_{II}^2)).$$

Из первых двух неравенств получаем следующие неравенства:

$$\wp(\wp(P_0^\varepsilon | D_{II}^1) < \underline{h}^\beta(P_0^\varepsilon | D_{II}^1)) < \beta,$$

$$\wp(\wp(P_0^\varepsilon | D_{II}^2) < \underline{h}^\beta(P_0^\varepsilon | D_{II}^2)) < \beta.$$

Применяя теорему о сложении вероятностей, получим

$$\wp((\wp(P_0^\varepsilon | D_{II}^1) < \underline{h}^\beta(P_0^\varepsilon | D_{II}^1)) \vee (\wp(P_0^\varepsilon | D_{II}^2) < \underline{h}^\beta(P_0^\varepsilon | D_{II}^2))) < 2\beta,$$

$$\wp((\wp(P_0^\varepsilon | D_{II}^1) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^1)) \& (\wp(P_0^\varepsilon | D_{II}^2) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^2))) \geq 1 - 2\beta. \quad (9)$$

Нетрудно видеть, что из условия

$$(\wp(P_0^\varepsilon | D_{II}^1) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^1)) \& (\wp(P_0^\varepsilon | D_{II}^2) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^2))$$

вытекает следующее неравенство

$$1 - (1 - \wp(P_0^\varepsilon | D_{II}^1))(1 - \wp(P_0^\varepsilon | D_{II}^2)) \geq 1 - (1 - \underline{h}^\beta(P_0^\varepsilon | D_{II}^1))(1 - \underline{h}^\beta(P_0^\varepsilon | D_{II}^2))$$

Из условия и леммы 1 следует, что левая часть этого неравенства равна  $\wp(P_0^\varepsilon | D_{II}^1 \& D_{II}^2)$ .

Обозначим через  $\underline{h}^{2\beta}(P_0^\varepsilon | D_{II}^1 \& D_{II}^2)$  величину  $1 - (1 - \underline{h}^\beta(P_0^\varepsilon | D_{II}^1))(1 - \underline{h}^\beta(P_0^\varepsilon | D_{II}^2))$  правую часть неравенства, тогда вероятность неравенства

$$\wp(\wp(P_0^\varepsilon | D_{II}^1 \& D_{II}^2) \geq \underline{h}^{2\beta}(P_0^\varepsilon | D_{II}^1 \& D_{II}^2))$$

равна вероятности, стоящей в левой части неравенства (9)

$$\wp((\wp(P_0^\varepsilon | D_{II}^1) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^1)) \& (\wp(P_0^\varepsilon | D_{II}^2) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^2))).$$

Отсюда следует, что  $\wp(\wp(P_0^\varepsilon | D_{II}^1 \& D_{II}^2) \geq \underline{h}^{2\beta}(P_0^\varepsilon | D_{II}^1 \& D_{II}^2)) \geq 1 - 2\beta$  ■.

**Определение 3.** Закономерности  $D_{II}^1 \rightarrow P_0^\varepsilon$ ,  $D_{II}^2 \rightarrow P_0^\varepsilon, \dots, D_{II}^n \rightarrow P_0^\varepsilon$  будем называть *независимыми*, если закономерности  $D_{II}^1 \rightarrow P_0^\varepsilon$  и  $D_{II}^2 \rightarrow P_0^\varepsilon$  независимы, закономерности  $D_{II}^1 \& D_{II}^2 \rightarrow P_0^\varepsilon$  и  $D_{II}^3 \rightarrow P_0^\varepsilon$  независимы и так далее, закономерности  $D_{II}^1 \& \dots \& D_{II}^{n-1} \rightarrow P_0^\varepsilon$  и  $D_{II}^n \rightarrow P_0^\varepsilon$  независимы.

**Лемма 3.** Для независимых закономерностей  $D_{II}^1 \rightarrow P_0^\varepsilon, \dots, D_{II}^n \rightarrow P_0^\varepsilon$  выполняются следующие соотношения, обобщающие результаты предыдущих лемм

$$1. \quad \wp(P_0^\varepsilon | D_{II}^1 \& \dots \& D_{II}^n) = 1 - \prod_{i=1}^n (1 - \wp(P_0^\varepsilon | D_{II}^i)),$$

$$2. \wp(\wp(P_0^\varepsilon | D_{II}^1 \& \dots \& D_{II}^n) \geq \underline{h}^{n\beta}(P_0^\varepsilon | D_{II}^1 \& \dots \& D_{II}^n)) \geq 1 - n\beta,$$

$$\underline{h}^{n\beta}(P_0^\varepsilon | D_{II}^1 \& \dots \& D_{II}^n) = 1 - \prod_{i=1}^n (1 - \underline{h}^\beta(P_0^\varepsilon | D_{II}^i)).$$

**Доказательство.**

1) Доказательство первого равенства проведем по индукции. Случай  $n = 2$  доказан в лемме 2.

Пусть закономерности  $D_{II}^1 \& \dots \& D_{II}^k \rightarrow P_0^\varepsilon$  и  $D_{II}^{k+1} \rightarrow P_0^\varepsilon$  независимы. Тогда применяя к ним лемму 2 получим

$$\wp(P_0^\varepsilon | D_{II}^1 \& \dots \& D_{II}^{k+1}) = 1 - (1 - \wp(P_0^\varepsilon | D_{II}^1 \& \dots \& D_{II}^k))(1 - \wp(P_0^\varepsilon | D_{II}^{k+1})).$$

По индуктивному предположению

$$\wp(P_0^\varepsilon | D_{II}^1 \& \dots \& D_{II}^k) = 1 - \prod_{i=1}^k (1 - \wp(P_0^\varepsilon | D_{II}^i)).$$

Подставляя это выражение в предыдущую формулу, получим равенство 1 леммы.

2) Доказательство пункта 2 леммы проводится аналогично предыдущему пункту, только надо иметь в виду, что результат леммы 2 не измениться, если нижние доверительные границы  $\underline{h}^\beta(P_0^\varepsilon | D_{II}^i)$ ,  $i = 1, 2$  будут иметь различный доверительный уровень  $\beta$  ■.

Выясним соотношение свойства независимости двух закономерностей  $D_{II}^1 \rightarrow P_0^\varepsilon$ ,  $D_{II}^2 \rightarrow P_0^\varepsilon$  и критерия для отбора закономерностей, применяемого в методе обнаружения закономерностей. Покажем, что свойство независимости означает независимость вклада этих закономерностей в предсказание  $P_0^\varepsilon$ .

**Лемма 4.** Если закономерности  $D_{II}^1 \rightarrow P_0^\varepsilon$  и  $D_{II}^2 \rightarrow P_0^\varepsilon$  независимы, то формулы  $D_{II}^1$  и  $D_{II}^2$  «существенны» в закономерности  $D_{II}^1 \& D_{II}^2 \rightarrow P_0^\varepsilon$ , т.е. для них выполняется приведённое ниже условие вероятностных закономерностей [4].

**Доказательство.** Из независимости следует, что

$$\begin{aligned} \wp(P_0^\varepsilon | D_{II}^1 \& D_{II}^2) &= 1 - (1 - \wp(P_0^\varepsilon | D_{II}^1))(1 - \wp(P_0^\varepsilon | D_{II}^2)) = \\ &= \wp(P_0^\varepsilon | D_{II}^1) + \wp(P_0^\varepsilon | D_{II}^2) - \wp(P_0^\varepsilon | D_{II}^1)\wp(P_0^\varepsilon | D_{II}^2). \end{aligned}$$

Так как вероятность находится в интервале  $[0, 1]$ , то

$$\wp(P_0^\varepsilon | D_{II}^1 \& D_{II}^2) > \wp(P_0^\varepsilon | D_{II}^1), \quad \wp(P_0^\varepsilon | D_{II}^1 \& D_{II}^2) > \wp(P_0^\varepsilon | D_{II}^2),$$

что является условием «существенности» формул  $D_{II}^1$  и  $D_{II}^2$  в закономерности ■.

Таким образом, если закономерности  $D_{II}^1 \rightarrow P_0^\varepsilon$  и  $D_{II}^2 \rightarrow P_0^\varepsilon$  независимы, то их можно улучшить добавлением в посылку одной закономерности посылки из другой зако-



номерности. Если существенность формулы  $D_{II}^1$  для закономерности  $D_{II}^2 \rightarrow P_0^\varepsilon$  и формулы  $D_{II}^2$  для закономерности  $D_{II}^1 \rightarrow P_0^\varepsilon$  может быть установлена на обучающем материале, как описано в методе обнаружения закономерностей, то должна будет обнаружена и закономерность  $D_{II}^1 \& D_{II}^2 \rightarrow P_0^\varepsilon$ . Поэтому предположение о независимости закономерностей означает, что они улучшают друг друга, но по обучающему материалу это нельзя обнаружить вследствие его ограниченности.

Предположение о независимости позволяет нам по закономерностям  $D_{II}^1 \rightarrow P_0^\varepsilon, \dots, D_{II}^n \rightarrow P_0^\varepsilon$  определить общую закономерность  $D_{II}^1 \& \dots \& D_{II}^n \rightarrow P_0^\varepsilon$  и вычислить для неё оценку вероятности. Поэтому отображение  $\nu$  будет иметь вид, приведенный в лемме 3. Для такого предсказания будет выполнено вероятностное неравенство, приведенное в лемме 3.

Если предположение о независимости сделать нельзя, то предсказание по закономерностям  $D_{II}^1 \rightarrow P_0^\varepsilon, \dots, D_{II}^n \rightarrow P_0^\varepsilon$  можно сделать другим способом. Для этого можно выбрать максимальную оценку условной вероятности. Но сама операция выбора изменит эту оценку. Поэтому надо учесть это изменение.

**Лемма 5.** Если среди оценок  $\underline{h}^\beta(P_0^\varepsilon | D_{II}^1), \dots, \underline{h}^\beta(P_0^\varepsilon | D_{II}^n)$  оценка  $\underline{h}^\beta(P_0^\varepsilon | D_{II}^k)$  является максимальной, и мы включим в процедуру нахождения лучшей закономерности выбор закономерности  $D_{II}^k \rightarrow P_0^\varepsilon$  среди остальных по критерию максимальной оценки, то оценка вероятности изменится и станет равной  $\underline{h}^{n\beta}(P_0^\varepsilon | D_{II}^k)$ , т.е.

$$\wp(\wp(P_0^\varepsilon | D_{II}^k) \geq \underline{h}^{n\beta}(P_0^\varepsilon | D_{II}^k)) \geq 1 - n\beta.$$

**Доказательство.** Проводя такие же рассуждения, что и при выводе неравенства (9) можно получить неравенство

$$\wp((\wp(P_0^\varepsilon | D_{II}^1) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^1)) \& \dots \& (\wp(P_0^\varepsilon | D_{II}^n) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^n))) \geq 1 - n\beta.$$

Из конъюнкции неравенств следует, что такое неравенство должно выполняться отдельно для каждого члена, в том числе и того, у которого оценка  $\underline{h}^\beta(P_0^\varepsilon | D_{II}^k)$  максимальна. Поэтому вероятность, стоящая в левой части неравенства, только увеличивается, и само неравенство сохранится, если мы конъюнкцию неравенств заменим одним неравенством  $\wp(P_0^\varepsilon | D_{II}^k) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^k)$ , в котором оценка  $\underline{h}^\beta(P_0^\varepsilon | D_{II}^k)$  – максимальна. Тогда получим неравенство  $\wp(\wp(P_0^\varepsilon | D_{II}^k) \geq \underline{h}^\beta(P_0^\varepsilon | D_{II}^k)) \geq 1 - n\beta$ . Остается только индекс  $\beta$  у оценки  $\underline{h}^\beta(P_0^\varepsilon | D_{II}^k)$  заменить на индекс  $n\beta$ , поскольку в силу неравенства доверительный уровень равен  $n\beta$  ■.

Действуя таким образом, мы свели все закономерности к одной  $D_{\Pi}^k \rightarrow P_0^\varepsilon$  и вычислили для нее оценку. Поэтому отображение  $\nu$  определяется также как и в предыдущих случаях.

Рассмотрим случай, когда множества  $Reg_1(Th)$  и  $Reg_2(Th)$  состоят из закономерностей  $D_{\Pi_1}^1 \rightarrow P_0^\varepsilon, \dots, D_{\Pi_1}^{n_1} \rightarrow P_0^\varepsilon$ , предсказывающих литеру  $P_0^\varepsilon$ , и закономерностей  $D_{\Pi_2}^1 \rightarrow \bar{P}_0^\varepsilon, \dots, D_{\Pi_2}^{n_2} \rightarrow \bar{P}_0^\varepsilon$ , предсказывающих отрицание этой же литеры. Осуществим, принимая предположение о независимости (лемма 3) или не принимая никаких предположений (лемма 5) предсказание литеры  $P_0^\varepsilon$  по первым закономерностям и предсказание литеры  $\bar{P}_0^\varepsilon$  по вторым. Получим соответствующие оценки. На этом можно остановиться, получив предсказание двух различных множеств  $PS(P_0^\varepsilon)$  и  $PS(\bar{P}_0^\varepsilon)$ , но можно получить и результирующее предсказание.

Предположение, аналогичное предположению о независимости, в данном случае принять нельзя, так как предсказания литеры  $P_0^\varepsilon$  противоречат предсказаниям отрицания этой же литеры  $\bar{P}_0^\varepsilon$ . В этом случае можно действовать также как в лемме 5. Для этого достаточно заметить, что в доказательстве леммы 5 могут использоваться не только закономерности, предсказывающие  $P_0^\varepsilon$ , но и закономерности, предсказывающие  $\bar{P}_0^\varepsilon$ . Тогда выберем из всех закономерностей  $Reg_1(Th) \cup Reg_2(Th)$  одну с максимальной оценкой условной вероятности. Отображение  $\nu$  определяется также как и в предыдущем случае.

#### 4.4 Предсказание по закономерностям $Reg_1(Th) \cup Reg_2(Th)$ в общем случае.

Пусть  $(P_0^{\varepsilon_0})_1, (\bar{P}_0^{\varepsilon_0})_1, \dots, (P_0^{\varepsilon_0})_l, (\bar{P}_0^{\varepsilon_0})_l$  – все литеры или/и их отрицания предсказываемые закономерностями  $Reg_1(Th) \cup Reg_2(Th)$ . Найдем оценки условных вероятностей для этих литер, как указано выше. Пусть  $(P_0^{\varepsilon_0})_1, (\bar{P}_0^{\varepsilon_0})_1, \dots, (P_0^{\varepsilon_0})_l, (\bar{P}_0^{\varepsilon_0})_l$  – все литеры, для которых получены соответствующие оценки. Эти литеры определяют множества  $PS((P_0^{\varepsilon_0})_1), PS((\bar{P}_0^{\varepsilon_0})_1), \dots, PS((P_0^{\varepsilon_0})_l), PS((\bar{P}_0^{\varepsilon_0})_l)$ . По этим множествам нужно получить предсказание для пересечения любого числа из этих множеств. На основании каких дополнительных предположений это можно сделать? Наиболее естественным нам представляется следующее предположение о независимости (здесь независимость понимается в другом смысле и относится к другим закономерностям).

**Определение 4.** Закономерности  $D_{\Pi}^1 \rightarrow (P_0^{\varepsilon_0})_1, \dots, D_{\Pi}^n \rightarrow (P_0^{\varepsilon_0})_n$  будем называть *независимыми*, если имеет место равенство

$$\wp((P_0^{\varepsilon_0})_1 \& \dots \& (P_0^{\varepsilon_0})_n \mid D_{\Pi}^1 \& \dots \& D_{\Pi}^n) = \prod_{i=1}^n \wp((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i). \quad (10)$$

**Лемма 6.** Если закономерности  $D_{\Pi}^1 \rightarrow (P_0^{\varepsilon_0})_1, \dots, D_{\Pi}^k \rightarrow (P_0^{\varepsilon_0})_k$  независимы согласно определению 4, и для каждой из них выполнено неравенство

$$\wp(\wp((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i) \geq \underline{h}^{\beta_i}((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i)) \geq 1 - \beta_i, \quad i = 1, 2, \dots, k,$$

то имеет место неравенство

$$\wp(\wp((P_0^{\varepsilon_0})_1 \& \dots \& (P_0^{\varepsilon_0})_k \mid D_{\Pi}^1 \& \dots \& D_{\Pi}^k) \geq \underline{h}^{\gamma}((P_0^{\varepsilon_0})_1 \& \dots \& (P_0^{\varepsilon_0})_k \mid D_{\Pi}^1 \& \dots \& D_{\Pi}^k)) \geq 1 - \gamma,$$

$$\gamma = \sum_{i=1}^k \beta_i,$$

$$\underline{h}^{\gamma}((P_0^{\varepsilon_0})_1 \& \dots \& (P_0^{\varepsilon_0})_k \mid D_{\Pi}^1 \& \dots \& D_{\Pi}^k) = \prod_{i=1}^k \underline{h}^{\beta_i}((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i).$$

**Доказательство.** Из условия следует, что для каждого  $i = 1, 2, \dots, k$  выполняется неравенство  $\wp(\wp((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i) < \underline{h}^{\beta_i}((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i)) < \beta_i$ . Отсюда следует, что

$$\wp(\wp(\wp((P_0^{\varepsilon_0})_1 \mid D_{\Pi}^1) \geq \underline{h}^{\beta_1}((P_0^{\varepsilon_0})_1 \mid D_{\Pi}^1)) \& \dots \& (\wp((P_0^{\varepsilon_0})_k \mid D_{\Pi}^k) \geq \underline{h}^{\beta_k}((P_0^{\varepsilon_0})_k \mid D_{\Pi}^k))) \geq 1 - \sum_{i=1}^k \beta_i.$$

Из истинности конъюнкции, стоящей в левой части неравенства следует, что

$$\prod_{i=1}^k \wp((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i) \geq \prod_{i=1}^k \underline{h}^{\beta_i}((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i).$$

Отсюда следует, что

$$\wp(\prod_{i=1}^k \wp((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i) \geq \prod_{i=1}^k \underline{h}^{\beta_i}((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i)) \geq 1 - \sum_{i=1}^k \beta_i.$$

Заменяя первое произведение на условную вероятность, согласно равенству (10), получим утверждение леммы ■.

Пусть у нас есть пересечение множеств, взятых из подмножества множества  $\{PS((P_0^{\varepsilon_0})_1), \dots, PS((P_0^{\varepsilon_0})_k)\}$ , например,  $PS((P_0^{\varepsilon_0})_1) \cap \dots \cap PS((P_0^{\varepsilon_0})_1)$ . Определим отображение  $\nu$  следующим образом:  $\nu(pr) = 0$ , если  $pr \notin PS((P_0^{\varepsilon_0})_1) \cap \dots \cap PS((P_0^{\varepsilon_0})_1)$ , и

$$\nu(pr) = \prod_{i=1}^k \underline{h}^{\beta_i}((P_0^{\varepsilon_0})_i \mid D_{\Pi}^i), \quad \text{если } pr \in PS((P_0^{\varepsilon_0})_1) \cap \dots \cap PS((P_0^{\varepsilon_0})_1).$$

Такое предсказание, в предположении независимости, имеет оценку, приведенную в лемме 6. Другие возможности использования предсказаний приведены в лемме 10.

#### 4.5 Предсказание по закономерностям $Reg_3(Th)$ .

Рассмотрим сначала случай, когда  $Reg_3(Th)$  состоит из одной закономерности  $D_{\Pi} \rightarrow P_0^{\varepsilon}$ , а множество  $\Pi$  для данной закономерности состоит из индивидуальной постоян-

ной  $z$ . Для такой закономерности в был получен толерантный интервал  $[\underline{h}^z, \bar{h}^z]$ , для которого выполняется неравенство (5). Для каждой модели  $pr \in PS$ , подставляя вместо индивидуальной постоянной  $z$  объект  $b$  (см. вывод формулы (4)), можно подсчитать частоту  $h_z^b(pr)$ . Определим в качестве предсказания подмножество  $PS(D_{II} \rightarrow P_0^\varepsilon) \subset PS$  тех моделей  $pr$ , для которых  $\underline{h}^z \leq h_z^b(pr) \leq \bar{h}^z$ . В силу неравенства (5) для заданных  $\beta$  и  $\gamma$  будет иметь место неравенство

$$\wp(\wp(pr \in PS(D_{II} \rightarrow P_0^\varepsilon)) \geq \gamma) \geq 1 - \beta. \quad (11)$$

Отображение  $\nu$  в данном случае можно определить следующим образом:  $\nu(pr) = \theta$ , если  $pr \notin PS(D_{II} \rightarrow P_0^\varepsilon)$  и  $\nu(pr) = \gamma$ , если  $pr \in PS(D_{II} \rightarrow P_0^\varepsilon)$ .

Рассмотрим случай произвольного множества закономерностей из  $Reg_3(Th)$ . Пусть есть закономерности  $D_{II}^1 \rightarrow (P_0^{\varepsilon_0})_1, \dots, D_{II}^k \rightarrow (P_0^{\varepsilon_0})_k$  из  $Reg_3(Th)$ . Для получения результирующего предсказания можно также, либо ввести предположение о независимости, либо не делать дополнительных предположений.

**Определение 5.** Закономерности  $D_{II}^1 \rightarrow (P_0^{\varepsilon_0})_1, \dots, D_{II}^k \rightarrow (P_0^{\varepsilon_0})_k$  будем называть независимыми относительно индивидуальных постоянных  $z_1 \in \Pi_1, \dots, z_k \in \Pi_k$ , если для случайных величин  $h_{z_1}^b, \dots, h_{z_k}^b$ , имеет место равенство

$$\wp((\underline{h}^{z_1} \leq h_{z_1}^b \leq \bar{h}^{z_1}) \& \dots \& (\underline{h}^{z_k} \leq h_{z_k}^b \leq \bar{h}^{z_k})) = \prod_{i=1}^k \wp(\underline{h}^{z_i} \leq h_{z_i}^b \leq \bar{h}^{z_i}).$$

**Лемма 7.** Если закономерности  $D_{II}^1 \rightarrow (P_0^{\varepsilon_0})_1, \dots, D_{II}^k \rightarrow (P_0^{\varepsilon_0})_k$  независимы относительно индивидуальных постоянных  $z_1, \dots, z_k$  и для них выполнены неравенства

$$\wp(\wp(\underline{h}^{z_i} \leq h_{z_i}^b \leq \bar{h}^{z_i}) \geq \gamma_i) \geq \beta_i, \quad i = 1, \dots, k,$$

то имеет место неравенство

$$\wp(\wp((\underline{h}^{z_1} \leq h_{z_1}^b \leq \bar{h}^{z_1}) \& \dots \& (\underline{h}^{z_k} \leq h_{z_k}^b \leq \bar{h}^{z_k})) \geq \gamma) \geq 1 - \beta, \quad (12)$$

$$\gamma = \prod_{i=1}^k \gamma_i, \quad \beta = \sum_{i=1}^k \beta_i.$$

**Доказательство.** Проводится также, как и доказательство леммы 6.

Каждой закономерности соответствует предсказание подмножества  $PS(D_{II}^i \rightarrow (P_0^{\varepsilon_0})_i) \subset PS$ ,  $i = 1, \dots, k$ . Используя лемму 7, можно, в предположении независимости, получить предсказание и соответствующие оценки для пересечения произвольного числа этих множеств. Например, если предположение о независимости можно сде-

для всех  $k$  закономерностей, то отображение  $\nu$  можно определить следующим образом:

$$\nu(pr) = 0, \text{ если } pr \notin PS(D_{\Pi}^1 \rightarrow (P_0^{\varepsilon_0})_1) \cap \dots \cap PS(D_{\Pi}^k \rightarrow (P_0^{\varepsilon_0})_k) \text{ и}$$

$$\nu(pr) = \gamma = \prod_{i=1}^k \gamma_i, \text{ если } pr \in PS(D_{\Pi}^1 \rightarrow (P_0^{\varepsilon_0})_1) \cap \dots \cap PS(D_{\Pi}^k \rightarrow (P_0^{\varepsilon_0})_k).$$

Для так определённого предсказания будет выполнено неравенство (12).

Предположим, что мы не можем принять предположение о независимости. Тогда мы можем выбрать одно из предсказаний  $PS(D_{\Pi}^e \rightarrow (P_0^{\varepsilon_0})_i)$ , являющееся лучшим по какому-нибудь критерию. В предыдущем случае таким критерием был максимум нижней доверительной границы, что вполне естественно для закономерностей  $Reg_1(Th)$  и  $Reg_2(Th)$ . Для закономерностей  $Reg_3(Th)$  может быть много различных критериев, зависящих от типов шкал, от соотношения желаемой точности и надёжности предсказания и т.д. Поэтому в данном случае необходимо определить некоторую функцию  $F$  на множестве  $PS$ , определяющую критерий качества предсказания.

Пусть есть некоторые закономерности  $D_{\Pi}^1 \rightarrow (P_0^{\varepsilon_0})_1, \dots, D_{\Pi}^k \rightarrow (P_0^{\varepsilon_0})_k$  и индивидуальные постоянные  $z_1 \in \Pi_1, \dots, z_k \in \Pi_k$ . По ним можно найти предсказания  $PS(D_{\Pi}^1 \rightarrow (P_0^{\varepsilon_0})_1), \dots, PS(D_{\Pi}^k \rightarrow (P_0^{\varepsilon_0})_k)$  удовлетворяющие неравенству (11). В качестве результирующего предсказания можно взять предсказание, на котором значение критерия качества  $F$  максимально. Примеры таких функций приведены ниже, а также в следующем параграфе. Найдем оценку вероятности такого предсказания, которое будет включать в себя процедуру выбора лучшего предсказания в соответствии со значением критерия  $F$ .

**Лемма 8.** Если для закономерностей  $D_{\Pi}^1 \rightarrow (P_0^{\varepsilon_0})_1, \dots, D_{\Pi}^k \rightarrow (P_0^{\varepsilon_0})_k$  относительно переменных  $z_1 \in \Pi_1, \dots, z_k \in \Pi_k$  выполнены неравенства

$$\wp(\wp(pr \in PS(D_{\Pi}^i \rightarrow (P_0^{\varepsilon_0})_i)) \geq \gamma_i) \geq 1 - \beta_i, \quad i = 1, 2, \dots, k,$$

то для закономерности с номером 1, выбираемой в соответствии со значением критерия  $F$ , выполнено неравенство

$$\wp(\wp(pr \in PS(D_{\Pi}^1 \rightarrow (P_0^{\varepsilon_0})_1)) \geq \gamma_1) \geq 1 - \sum_{i=1}^k \beta_i. \quad (13)$$

**Доказательство.** Проводиться аналогично доказательству леммы 5.

Закономерности из леммы 8 не обязательно все различны, поэтому для одной закономерности можно брать сразу несколько индивидуальных постоянных  $z_1, \dots, z_e \in \Pi$ . Таким образом, если не делать никаких предположений относительно закономерностей из

$Reg_3(Th)$ , то можно получить следующее предсказание: выбрать из всех (из части) закономерностей  $Reg_3(Th)$  по всем индивидуальным постоянным, входящим во множество  $\Pi$  этих закономерностей, ту закономерность  $D_{\Pi}^l \rightarrow (P_0^{\varepsilon_0})_l$  и индивидуальную постоянную  $z$ , дающую наилучшее предсказание в смысле выбранного критерия  $F$ . Для полученного предсказания будет выполнено неравенство (13). Отображение  $v$  тогда определяется следующим образом:

$$v(pr) = 0, \text{ если } pr \notin PS(D_{\Pi}^l \rightarrow (P_0^{\varepsilon_0})_l) \text{ и}$$

$$v(pr) = \gamma_e, \text{ если } pr \in PS(D_{\Pi}^l \rightarrow (P_0^{\varepsilon_0})_l).$$

Леммы 7, 8 дают возможность также получить оценку предсказания для случая, когда предположение о независимости можно сделать только относительно части закономерностей из  $Reg_3(Th)$ . В этом случае для закономерностей и индивидуальных постоянных, для которых можно сделать предположение о независимости, вычисление оценки осуществляется в соответствии с леммой 7, а после этого результирующее предсказание получается в соответствии с леммой 8.

Приведем наиболее естественные, с нашей точки зрения, критерии качества предсказания  $F$  для разных шкал.

Если предсказываемое отношение  $P_0^{\varepsilon}$  является отношением порядка  $\leq$ , то для одних и тех же значений  $\gamma$  и  $\beta$  то предсказание лучше, которое содержит в себе больше точек материала обучения. Точкам материала обучения  $a \in \mathbb{D}$  соответствуют такие предсказания  $pr_a \in PS$ , для которых  $a \leq b$  либо  $b \leq a$ . Поэтому наилучшим предсказанием с точки зрения этого критерия будет такое предсказание  $Q \subset PS$ , которое содержит максимальное число моделей из множества  $\{pr_{a_1}, \dots, pr_{a_m}\}$ ,  $\mathbb{D} = \{a_1, \dots, a_m\}$ . Значение критерия  $F(Q)$  — есть полученное максимальное число.

Если предсказываемая величина измеряется в шкале отношений, то имеет смысл говорить о величине интервала предсказания. Каждой модели  $pr \in PS$  можно поставить в соответствие интервал  $I_{pr}$  тех значений, соответствующих объекту  $b$ , которые определяются моделью  $pr$ . Критерий  $F$  тогда можно определить как сумму всех интервалов, входящих в предсказание  $Q \subset PS$ :  $F(Q) = \sum_{pr \in Q} I_{pr}$ . По максимуму этого критерия можно определить наилучшее предсказание.

#### 4.6 Общий случай произвольных множеств $Reg_1, Reg_2, Reg_3$ .

Получим предсказание для закономерностей  $Reg_1$  и  $Reg_2$ , как указано в леммах 5, 6 и по закономерностям из  $Reg_3$ , как указано в леммах 7, 8. Рассмотрим, также, как и предыдущих случаях, два способа получения результирующего предсказания: с использованием некоторого дополнительного предположения о независимости и без него. Рассмотрим первый способ.

**Определение 6.** Пусть есть некоторые закономерности  $D_{\Pi 1}^1 \rightarrow (P_0^{\varepsilon_0})_1, \dots, D_{\Pi 1}^n \rightarrow (P_0^{\varepsilon_0})_n$  из  $Reg_1 \cup Reg_2$  и закономерности  $D_{\Pi 2}^1 \rightarrow (P_0^{\varepsilon_0})_1, \dots, D_{\Pi 2}^k \rightarrow (P_0^{\varepsilon_0})_k$  из  $Reg_3$ , рассматриваемые относительно индивидуальных постоянных  $z_1 \in \Pi_1, \dots, z_k \in \Pi_k$ . Событие

$$h = ((\underline{h}^{z_1} \leq h_{z_1}^b \leq \bar{h}^{z_1}) \& \dots \& (\underline{h}^{z_k} \leq h_{z_k}^b \leq \bar{h}^{z_k}))$$

и условное событие

$$((P_0^{\varepsilon_0})_1 \& \dots \& (P_0^{\varepsilon_0})_n \mid D_{\Pi 1}^n \& \dots \& D_{\Pi 1}^n)$$

будем называть независимыми, если

$$\wp((P_0^{\varepsilon_0})_1 \& \dots \& (P_0^{\varepsilon_0})_n \& h \mid D_{\Pi 1}^n \& \dots \& D_{\Pi 1}^n) = \wp((P_0^{\varepsilon_0})_1 \& \dots \& (P_0^{\varepsilon_0})_n \mid D_{\Pi 1}^n \& \dots \& D_{\Pi 1}^n) \wp(h).$$

**Лемма 9.** Если для закономерности  $D_{\Pi 1}^n \& \dots \& D_{\Pi 1}^n \rightarrow (P_0^{\varepsilon_0})_1 \& \dots \& (P_0^{\varepsilon_0})_n$  справедлива оценка

$$\wp(\wp((P_0^{\varepsilon_0})_1 \& \dots \& (P_0^{\varepsilon_0})_n \mid D_{\Pi 1}^n \& \dots \& D_{\Pi 1}^n) \geq \underline{h}^{\beta_1}) \geq 1 - \beta_1,$$

приведенная в леммах 5, 6, а для закономерностей  $D_{\Pi 2}^1 \rightarrow (P_0^{\varepsilon_0})_1, \dots, D_{\Pi 2}^k \rightarrow (P_0^{\varepsilon_0})_k$  справедлива оценка

$$\wp(\wp(pr \in PS(D_{\Pi 2}^1 \rightarrow (P_0^{\varepsilon_0})_1) \cap \dots \cap PS(D_{\Pi 2}^k \rightarrow (P_0^{\varepsilon_0})_k)) \geq \gamma) \geq 1 - \beta_2,$$

приведенная в леммах 7, 8, и соответствующие события независимы, в соответствии с определением 8, то справедлива следующая оценка

$$\wp(\wp((P_0^{\varepsilon_0})_1 \& \dots \& (P_0^{\varepsilon_0})_n \& h \mid D_{\Pi 1}^n \& \dots \& D_{\Pi 1}^n) \geq \gamma^1) \geq 1 - \beta, \quad \gamma^1 = \underline{h}^{\beta_1} \gamma, \quad \beta = \beta_1 + \beta_2.$$

**Доказательство.** Проводится аналогично доказательству леммы 6.

Если предположение о независимости сделать нельзя, то предсказание можно получить путем выбора наилучшего из них в соответствии со значением критерия качества  $F$ . Для этого критерий  $F$  надо определить для множества  $PS(D_{\Pi 2}^1 \rightarrow (P_0^{\varepsilon_0})_1) \cap \dots \cap PS(D_{\Pi 2}^k \rightarrow (P_0^{\varepsilon_0})_k)$  и для множества  $PS((P_0^{\varepsilon_0})_1) \cap \dots \cap PS((P_0^{\varepsilon_0})_n)$ . Оценка выбранного предсказания находится также, как и в лемме 5.

Отображение  $\nu$  определяется также, как и в предыдущих случаях.

Рассмотрим дополнительные возможности получения предсказания. Закономерности из  $Reg_1 \cup Reg_2 \cup Reg_3$  определяют различные подмножества:  $PS((P_0^{e_0})_i)$  или  $PS(D_H^j \rightarrow (P_0^{e_0})_j)$ . Все рассмотренные предсказания состоят либо из этих подмножеств, либо из различных их пересечений. Для того, чтобы расширить возможности получения предсказаний, рассмотрим предсказания в виде объединения каких-либо подмножеств множества  $PS$ .

**Лемма 10.** Если множества  $PS_1, \dots, PS_n$ ;  $PS_i \in PS$ ,  $i = 1, 2, \dots, n$  попарно не пересекаются, и для каждого из них справедлива оценка  $\wp(\wp(PS_i) \geq \gamma_i) \geq 1 - \beta_i$ ,  $i = 1, 2, \dots, n$ , то для объединения этих множеств справедлива оценка

$$\wp(\wp(\bigcup_{i=1}^n PS_i) \geq \sum_{i=1}^n \gamma_i) \geq 1 - \sum_{i=1}^n \beta_i.$$

**Доказательство.** Из условия вытекает следующее неравенство

$$\wp((\wp(PS_1) \geq \gamma_1) \& \dots \& (\wp(PS_n) \geq \gamma_n)) \geq 1 - \sum_{i=1}^n \beta_i.$$

Если высказывание, стоящее в левой части неравенства истинно, то

$$\sum_{i=1}^n \wp(PS_i) \geq \sum_{i=1}^n \gamma_i.$$

Но так как множества  $PS_i$ ,  $i = 1, 2, \dots, n$  попарно не пересекаются, то

$$\sum_{i=1}^n \wp(PS_i) = \wp(\bigcup_{i=1}^n PS_i) \quad \blacksquare.$$

Таким образом, мы рассмотрели различные возможности определения алгоритма  $AP$ .

## §5. Метод принятия решений

В практических задачах полученные предсказания, как правило, используются для принятия определенных решений. К примеру, в финансовых задачах полученные предсказания о дальнейшем движении цен могут быть использованы для принятия решения о покупке или продажи соответствующей ценной бумаги. Рассмотрим общий механизм принятия решений на основе предсказаний.

Пусть нам известно множество вариантов решений  $S = \{s_1, \dots, s_n\}$ . Предположим также, что нам известно, какое одно конкретное решение должно быть выбрано в том случае, если бы нам было известно, какая модель из  $PS$  окажется истинной. Это означает, что для каждого решения  $s \in S$  можно указать множество моделей  $PS(s) \subset PS$ , для кото-



рых должно быть выбрано данное решение. Будем называть множества  $PS(s_i)$ ,  $i = 1, \dots, n$  вариантами исхода. Отметим также, что поскольку для каждой модели может быть выбрано только одно решение, то множества  $PS(s_i)$ ,  $i = 1, \dots, n$  попарно не пересекаются.

Поскольку реально нам не известна модель, которая на самом деле окажется истинной, то задачей принятия решения является выбор одного варианта решения, основываясь на прогнозах различных вариантов исходов. Таким образом, метод принятия решений мы можем определить как функцию  $Dec : \{\langle \lambda(PS(s_1)), \dots, \lambda(PS(s_n)) \rangle\} \rightarrow \{S \cup \emptyset\}$ , которая для каждого набора оценок вероятностей исходов  $\langle \lambda(PS(s_1)), \dots, \lambda(PS(s_n)) \rangle$  возвращает один выбранный вариант решения  $s \in S$  либо  $\emptyset$  в том случае, если решение не может быть принято.

Прежде чем переходить к построению метода принятия решений, необходимо задать функцию оценки вероятностей исходов  $\lambda : \{PS(s_1), \dots, PS(s_n)\} \rightarrow [0, 1]$ , которая для любого исхода  $PS(s) \subset PS$  вычисляет результирующую оценку вероятности данного исхода, основываясь на оценках  $v(pr)$  всех моделей  $pr \in PS(s)$ , входящих в этот исход  $PS(s)$ .

Способы определения функции  $\lambda$  могут быть различны и зависят от специфики решаемой задачи. Приведем два примера определения функции  $\lambda$ .

1. По максимальной оценке:  $\lambda(PS(s)) = \max_{pr \in PS(s)} \{v(pr)\}$ .

2. По средней оценке:  $\lambda(PS(s)) = \frac{\sum_{pr \in PS(s)} v(pr)}{n(PS(s))}$ , если  $n(PS(s)) \neq 0$ , и  $\lambda(PS(s)) = 0$ ,

если  $n(PS(s)) = 0$ . Здесь  $n(PS(s))$  – количество моделей во множестве  $PS(s)$ .

Для определения метода принятия решения в различных задачах также могут быть использованы различные способы задания функции  $Dec$ . В качестве одного достаточно универсального способа определения функции  $Dec$  можно предложить следующий вариант, учитывающий при принятии решения согласованность прогнозов различных вариантов исходов.

Для каждого варианта исхода  $PS(s_i)$ ,  $i = 1, \dots, n$  рассчитывается показатель согласованности его прогноза по формуле  $Ctrl_i = \kappa(PS(s_i)) - \max_{j \neq i} \{\kappa(PS(s_j))\}$ , т.е. как разность между оценкой вероятности данного исхода и максимальной оценкой вероятности остальных исходов. В качестве окончательного решения выбирается решение, соответствующее исходу, показатель согласованности которого строго больше заданного порога

$0 < \delta < 1$ , т.е.  $Dec = s_k$ , где  $k = \arg \max_{i=1, \dots, n} \{Ctr_i : Ctr_i > \delta\}$ . Порог  $\delta$  будем называть *порогом согласованности*. В случае, если не существует исхода, показатель согласованности которого выше указанного порога, то решение не принимается и  $Dec = \emptyset$ . Таким образом, регулируя порог согласованности  $\delta$ , можно регулировать степень уверенности при принятии решений.

## БЛАГОДАРНОСТИ

Работа поддержана грантом РФФИ 08-07-00272-а; интеграционными проектами СО РАН №№ 1, 115, а также работа выполнена при финансовой поддержке Государственного контракта 2007-4-1.4-00-04 и Совета по грантам Президента РФ и государственной поддержке ведущих научных школ (проект НШ-335.2008.1).

## ЛИТЕРАТУРА

1. Демин А.В., Витяев Е.Е. Разработка универсальной системы извлечения знаний «Discovery» и ее применения. Вестник НГУ, Новосибирск, 2008, (в печати).
2. Демин А.В., Витяев Е.Е. Реализация универсальной системы извлечения знаний «discovery» и ее применение в задачах финансового прогнозирования // Вычислительные системы, Новосибирск, 2008, настоящий сборник.
3. Витяев Е.Е. Анализ данных в языках эмпирических систем. Диссертация на соискание ученой степени кандидата технических наук, Новосибирск, 1983, р.192.
4. Витяев Е.Е. Извлечение знаний из данных. Компьютерное познание. Модели когнитивных процессов: Моногр. // Новосибирский гос. ун-т. Новосибирск, 2006. 293 с.
5. Model-theoretic logics. Eds: J. Barwise, S. Feferman, Springer-Verlag, NY, 1985.
6. Дискретная математика и математические вопросы кибернетики / Под ред. Яблонского С.В., Лупанова О.Б., т. 1. – М.: Наука, 1974. – 311 с.
7. Scheffe H., Tukey J.M. Non-parametric estimation 1. Validation of order statistics. – Ann. Math. Stat., v.16, 1945, p. 187-192.